



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Bayesian Networks (Part I)

Graphical Model Readings:

Murphy 10 – 10.2.1

Bishop 8.1, 8.2.2

HTF --

Mitchell 6.11

Matt Gormley
Lecture 22
April 10, 2017

Reminders

- Peer Tutoring
- Homework 7: Deep Learning
 - Release: Wed, Apr. 05
 - Part I due Wed, Apr. 12
 - Part II due Mon, Apr. 17

Start Early

CONVOLUTIONAL NEURAL NETS

Deep Learning Outline

- **Background: Computer Vision**
 - Image Classification
 - ILSVRC 2010 - 2016
 - Traditional Feature Extraction Methods
 - Convolution as Feature Extraction
- **Convolutional Neural Networks (CNNs)**
 - Learning Feature Abstractions
 - Common CNN Layers:
 - Convolutional Layer
 - Max-Pooling Layer
 - Fully-connected Layer (w/tensor input)
 - Softmax Layer
 - ReLU Layer
 - Background: Subgradient
 - Architecture: LeNet
 - Architecture: AlexNet
- **Training a CNN**
 - SGD for CNNs
 - Backpropagation for CNNs

Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

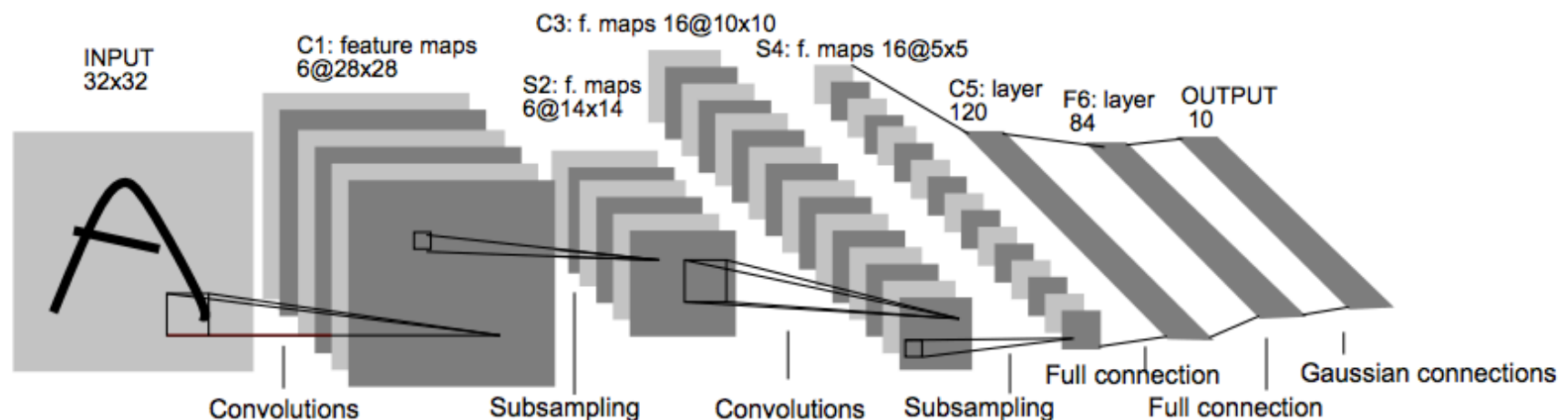


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Convolutional Layer

CNN key idea:
Treat convolution matrix as
parameters and learn them!

Input Image

0	0	0	0	0	0	0
0	1	1	1	1	1	0
0	1	0	0	1	0	0
0	1	0	1	0	0	0
0	1	1	0	0	0	0
0	1	0	0	0	0	0
0	0	0	0	0	0	0



Learned
Convolution

θ_{11}	θ_{12}	θ_{13}
θ_{21}	θ_{22}	θ_{23}
θ_{31}	θ_{32}	θ_{33}

Convolved Image

.4	.5	.5	.5	.4
.4	.2	.3	.6	.3
.5	.4	.4	.2	.1
.5	.6	.2	.1	0
.4	.3	.1	0	0

Downsampling by Averaging

- Downsampling by averaging **used to be** a common approach
- This is a special case of convolution where the weights are fixed to a uniform distribution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Convolution

$\frac{1}{4}$	$\frac{1}{4}$
$\frac{1}{4}$	$\frac{1}{4}$

Convolved Image

$\frac{3}{4}$	$\frac{3}{4}$	$\frac{1}{4}$
$\frac{3}{4}$	$\frac{1}{4}$	0
$\frac{1}{4}$	0	0

Max-Pooling

- Max-pooling is another (common) form of downsampling
- Instead of averaging, we take the max value within the same range as the equivalently-sized convolution
- The example below uses a stride of 2

Input Image

1	1	1	1	1	0
1	0	0	1	0	0
1	0	1	0	0	0
1	1	0	0	0	0
1	0	0	0	0	0
0	0	0	0	0	0

Max-
pooling

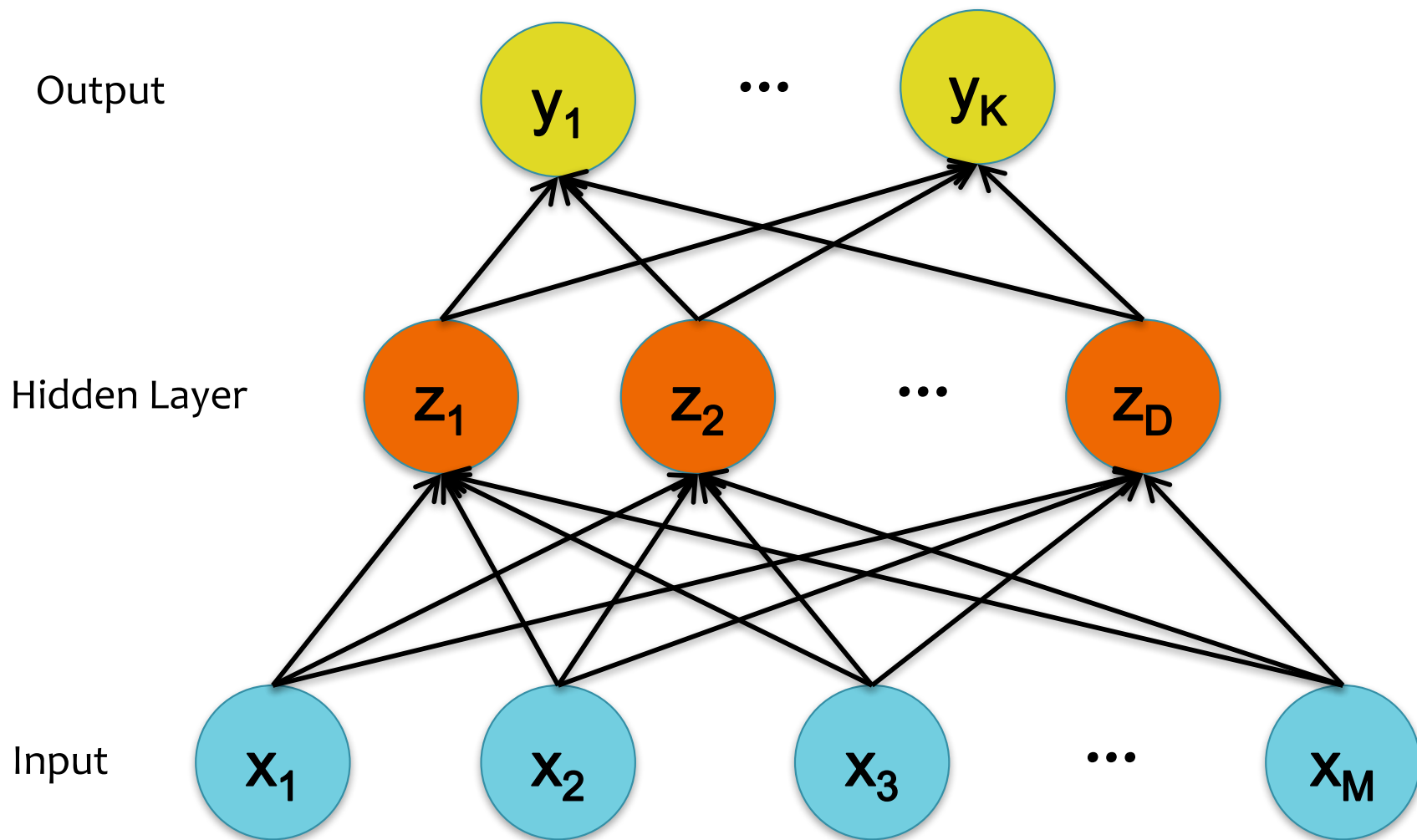
$x_{i,j}$	$x_{i,j+1}$
$x_{i+1,j}$	$x_{i+1,j+1}$

Max-Pooled
Image

1	1	1
1	1	0
1	0	0

$$y_{ij} = \max(x_{ij}, \\ x_{i,j+1}, \\ x_{i+1,j}, \\ x_{i+1,j+1})$$

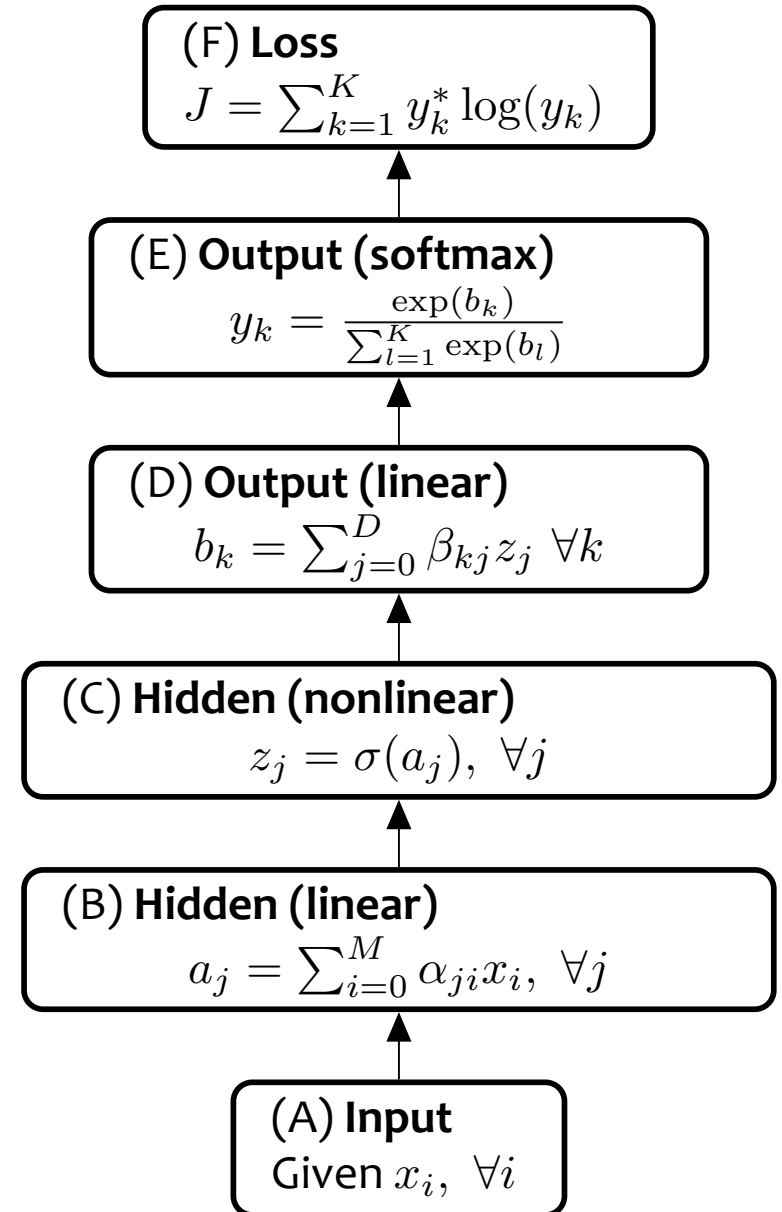
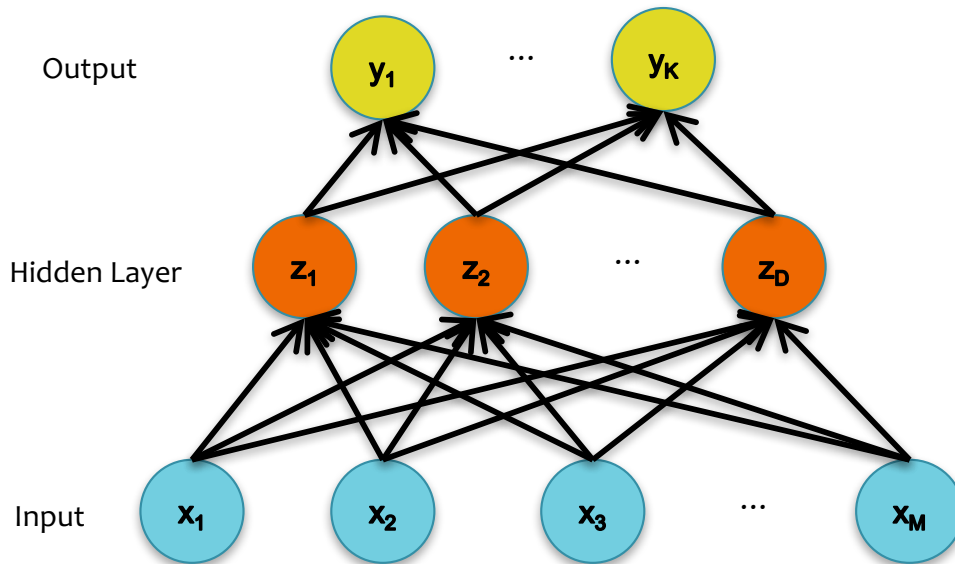
Multi-Class Output



Multi-Class Output

Softmax Layer:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^K \exp(b_l)}$$



Training a CNN

Whiteboard

- SGD for CNNs
- Backpropagation for CNNs

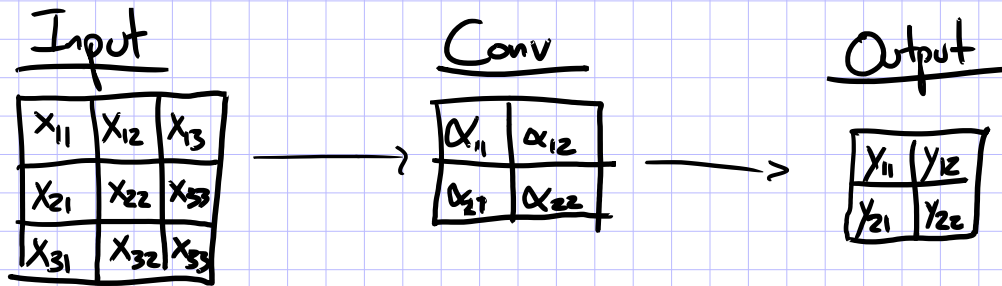
Common CNN Layers

Whiteboard

- ReLU Layer
- Background: Subgradient
- Fully-connected Layer (w/tensor input)
- Softmax Layer
- Convolutional Layer
- Max-Pooling Layer

Convolutional Layer

Ex: 1 input channel, 1 output channel



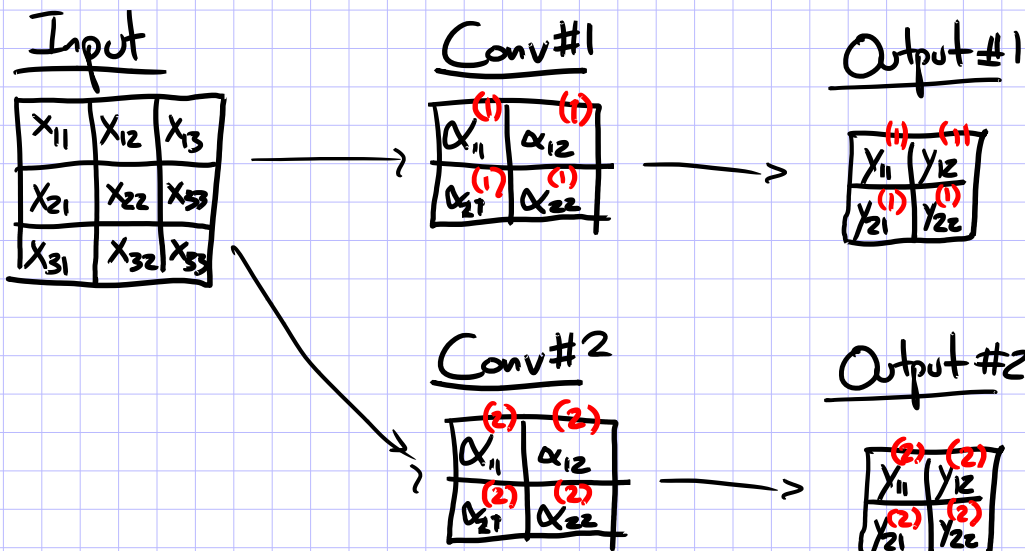
$$y_{11} = \alpha_{11}x_{11} + \alpha_{12}x_{12} + \alpha_{21}x_{21} + \alpha_{22}x_{22} + \alpha_0$$

$$y_{12} = \alpha_{11}x_{12} + \alpha_{12}x_{13} + \alpha_{21}x_{22} + \alpha_{22}x_{23} + \alpha_0$$

$$y_{21} = \alpha_{11}x_{21} + \alpha_{12}x_{22} + \alpha_{21}x_{31} + \alpha_{22}x_{32} + \alpha_0$$

$$y_{22} = \alpha_{11}x_{22} + \alpha_{12}x_{23} + \alpha_{21}x_{32} + \alpha_{22}x_{33} + \alpha_0$$

Ex: 1 input channel, 2 output channels



$$y_{11}^{(1)} = \alpha_{11}^{(1)}x_{11} + \alpha_{12}^{(1)}x_{12} + \alpha_{21}^{(1)}x_{21} + \alpha_{22}^{(1)}x_{22} + \alpha_0^{(1)}$$

$$y_{12}^{(1)} = \dots$$

$$y_{21}^{(1)} = \dots$$

$$y_{22}^{(1)} = \alpha_{11}^{(1)}x_{22} + \alpha_{12}^{(1)}x_{23} + \alpha_{21}^{(1)}x_{32} + \alpha_{22}^{(1)}x_{33} + \alpha_0^{(1)}$$

$$y_{11}^{(2)} = \alpha_{11}^{(2)}x_{11} + \alpha_{12}^{(2)}x_{12} + \alpha_{21}^{(2)}x_{21} + \alpha_{22}^{(2)}x_{22} + \alpha_0^{(2)}$$

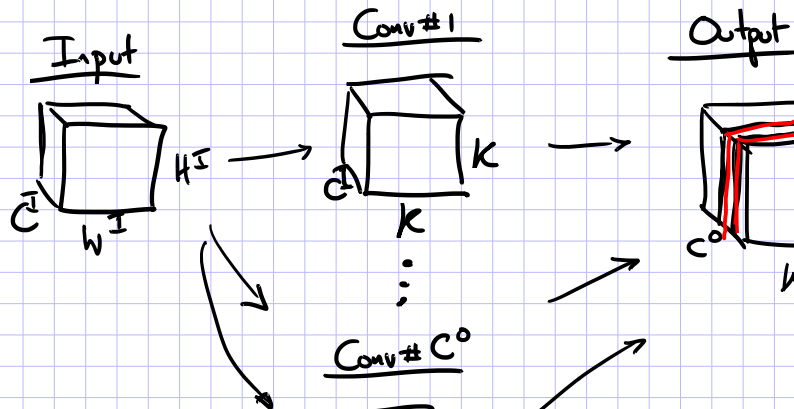
$$y_{12}^{(2)} = \dots$$

$$y_{21}^{(2)} = \dots$$

$$y_{22}^{(2)} = \alpha_{11}^{(2)}x_{22} + \alpha_{12}^{(2)}x_{23} + \alpha_{21}^{(2)}x_{32} + \alpha_{22}^{(2)}x_{33} + \alpha_0^{(2)}$$

Convolutional Layer

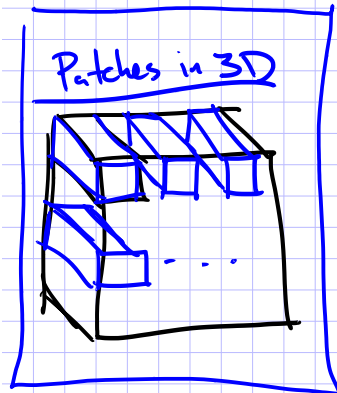
Ex: C^I input channels, C^O output channels



$$H^O = \lfloor (H^I + 2p - K) / s + 1 \rfloor$$

$$W^O = \lfloor (W^I + 2p - K) / s + 1 \rfloor$$

where p = # pixels of padding on input
 k = size of conv. matrix
 s = stride length



Forward:

$$y_{ij}^{(k)} = \alpha_0^{(k)} + \sum_{c=1}^{C^I} \sum_{q=1}^K \sum_{r=1}^K \alpha_{qr}^{(c)} x_{mn}^{(c)} \quad \text{where } m = s(i-1) + q, n = s(j-1) + r$$

Backward:

$$\frac{dJ}{d\alpha_0^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_0^{(k)}}$$

$$\frac{dJ}{d\alpha_{qr}^{(c)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{d\alpha_{qr}^{(c)}}$$

$$\frac{dJ}{dx_{mn}^{(c)}} = \sum_i \sum_j \sum_k \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(c)}}$$

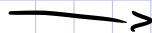
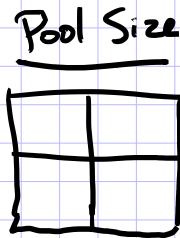
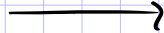
just some calculus

Max-Pooling Layer

Ex: 1 input channel, 1 output channel, stride of 1

Input

x_{11}	x_{12}	x_{13}
x_{21}	x_{22}	x_{23}
x_{31}	x_{32}	x_{33}

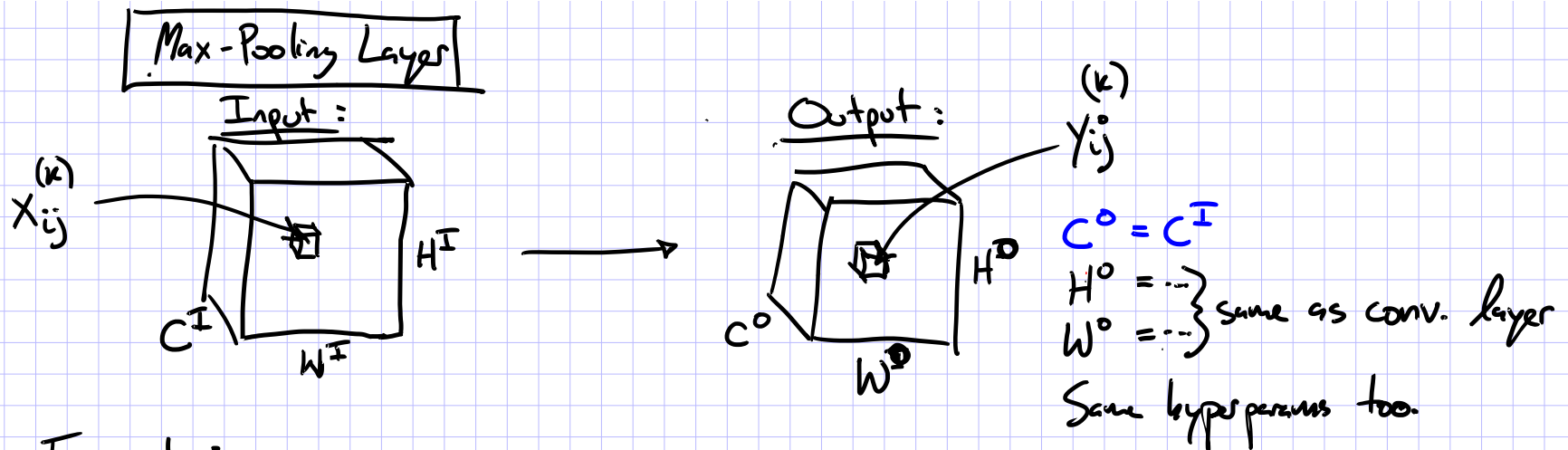


Output

y_{11}	y_{12}
y_{21}	y_{22}

$$\begin{aligned} y_{11} &= \max(x_{11}, x_{12}, x_{21}, x_{22}) \\ y_{12} &= \max(x_{12}, x_{13}, x_{22}, x_{23}) \\ y_{21} &= \max(x_{21}, x_{22}, x_{31}, x_{32}) \\ y_{22} &= \max(x_{22}, x_{23}, x_{32}, x_{33}) \end{aligned}$$

Max-Pooling Layer



Forward:

$$Y_{ij}^{(k)} = \max_{\substack{q \in \{1, \dots, k\} \\ r \in \{1, \dots, k\}}} X_{mn}^{(k)} \text{ where } m = s(i-1) + q \\ n = s(j-1) + r$$

Backward:

$$\frac{dJ}{dx_{mn}^{(k)}} = \sum_i \sum_j \frac{dJ}{dy_{ij}^{(k)}} \frac{dy_{ij}^{(k)}}{dx_{mn}^{(k)}}$$

Subderivatives

- + $\text{Max}()$ is not differentiable, but subdifferentiable.
- + There are a set of derivatives and we can just choose one for SGD.

$$y = \max(a, b)$$

$$\Rightarrow \frac{dJ}{da} = \frac{dJ}{dy} \frac{dy}{da} \text{ where } \frac{dy}{da} = \begin{cases} 1 & \text{if } a > b \\ 0 & \text{otherwise} \end{cases}$$

Convolutional Neural Network (CNN)

- Typical layers include:
 - Convolutional layer
 - Max-pooling layer
 - Fully-connected (Linear) layer
 - ReLU layer (or some other nonlinear activation function)
 - Softmax
- These can be arranged into arbitrarily deep topologies

Architecture #1: LeNet-5

PROC. OF THE IEEE, NOVEMBER 1998

7

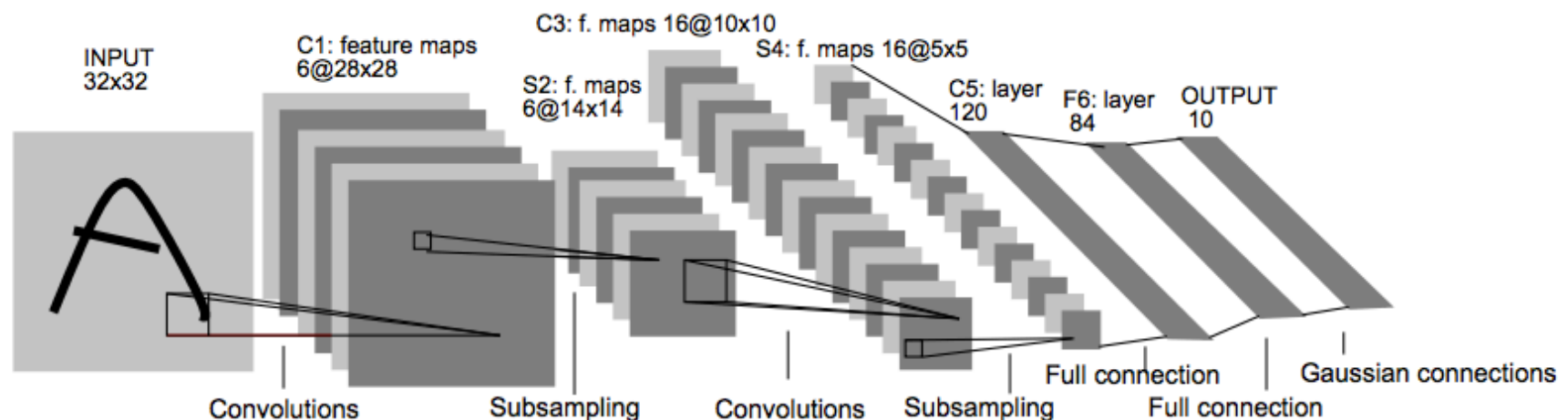


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Architecture #2: AlexNet

CNN for Image Classification

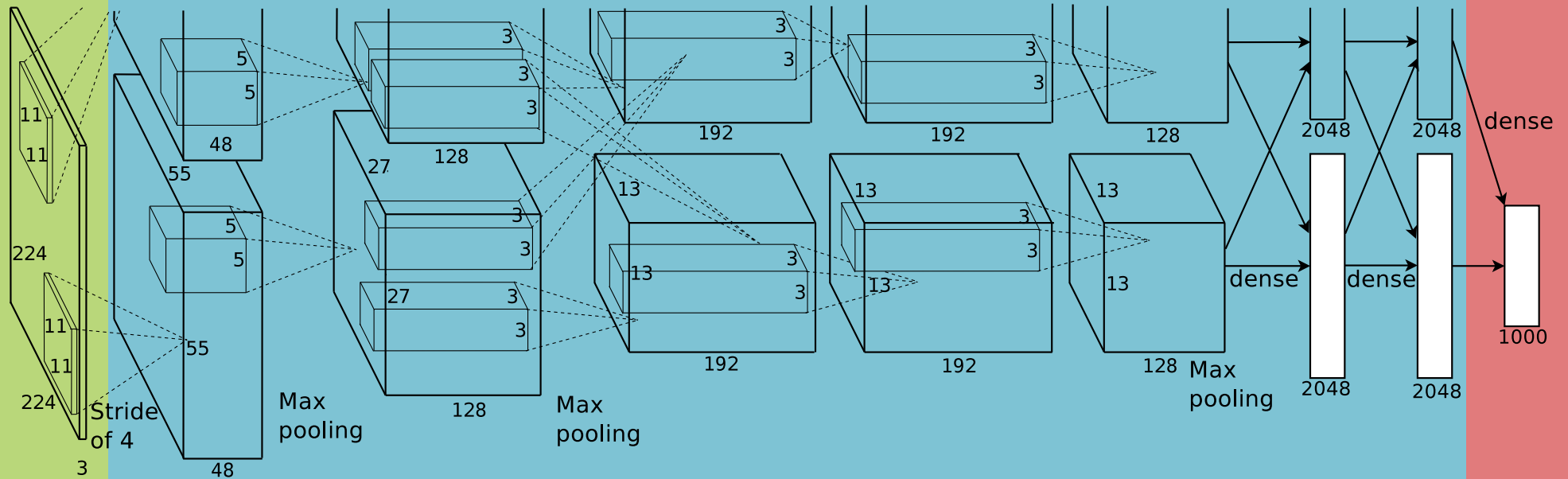
(Krizhevsky, Sutskever & Hinton, 2012)

15.3% error on ImageNet LSVRC-2012 contest

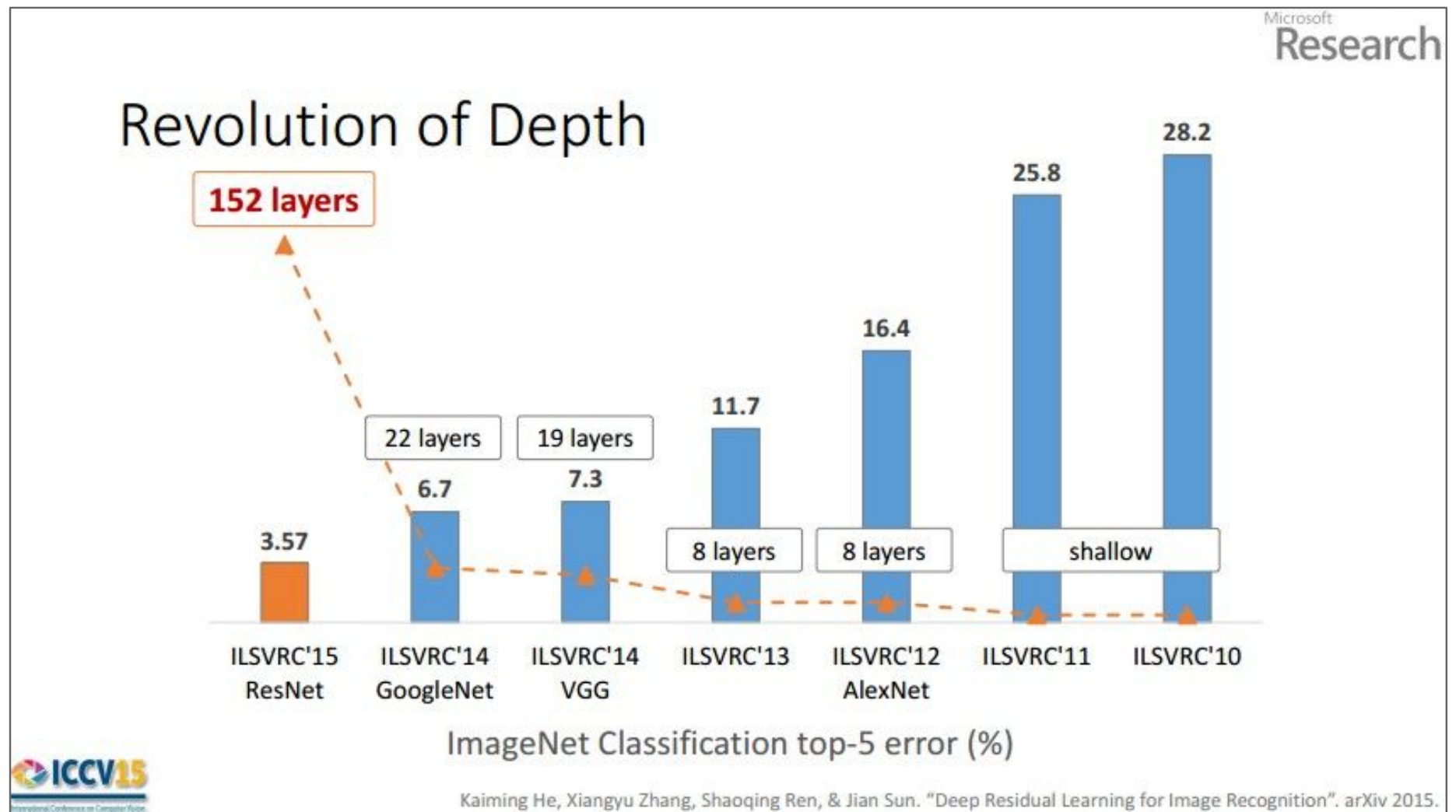
Input
image
(pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way
softmax



CNNs for Image Recognition



Mini-Batch SGD

- **Gradient Descent:**
Compute true gradient exactly from all N examples
- **Mini-Batch SGD:**
Approximate true gradient by the average gradient of K randomly chosen examples
- **Stochastic Gradient Descent (SGD):**
Approximate true gradient by the gradient of one randomly chosen example

Mini-Batch SGD

while not converged: $\theta \leftarrow \theta - \lambda \mathbf{g}$

Three variants of first-order optimization:

Gradient Descent: $\mathbf{g} = \nabla J(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \nabla J^{(i)}(\boldsymbol{\theta})$

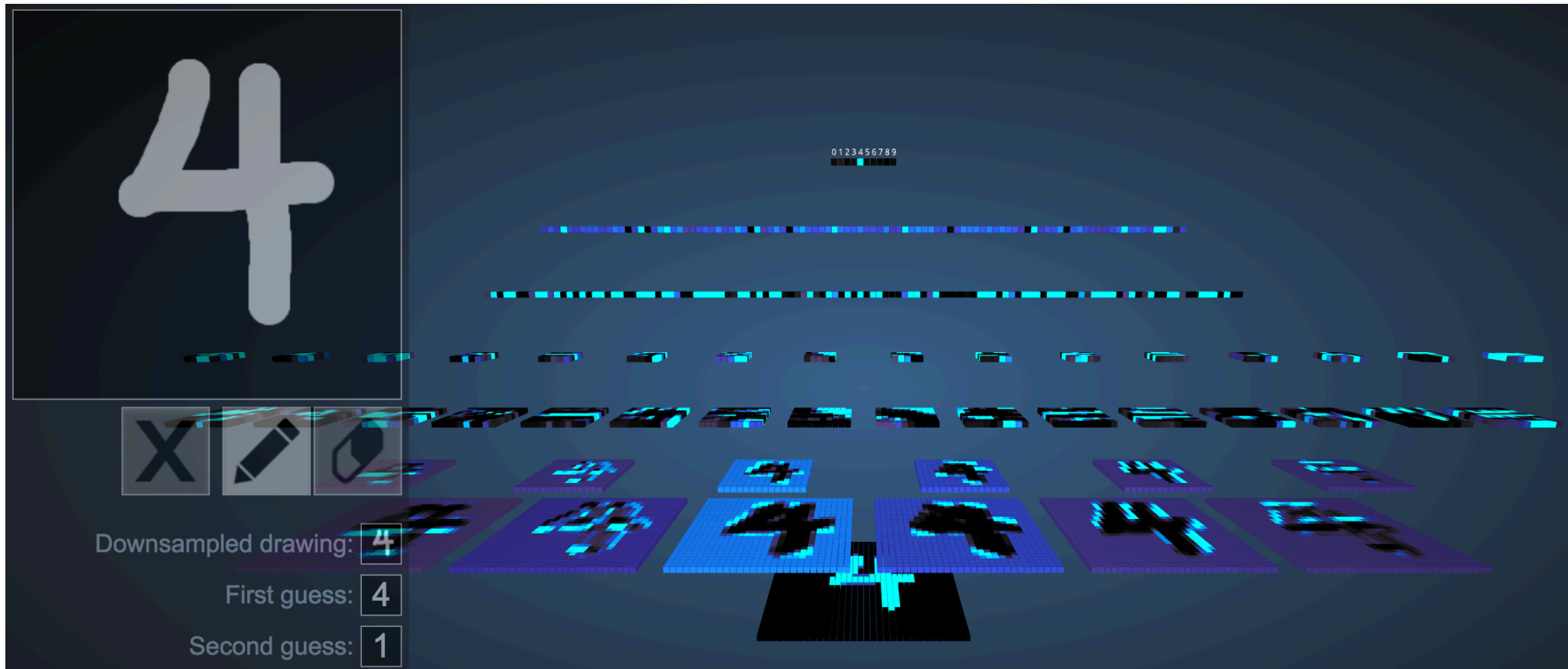
SGD: $\mathbf{g} = \nabla J^{(i)}(\boldsymbol{\theta})$ where i sampled uniformly

Mini-batch SGD: $\mathbf{g} = \frac{1}{S} \sum_{s=1}^S \nabla J^{(i_s)}(\boldsymbol{\theta})$ where i_s sampled uniformly $\forall s$

CNN VISUALIZATIONS

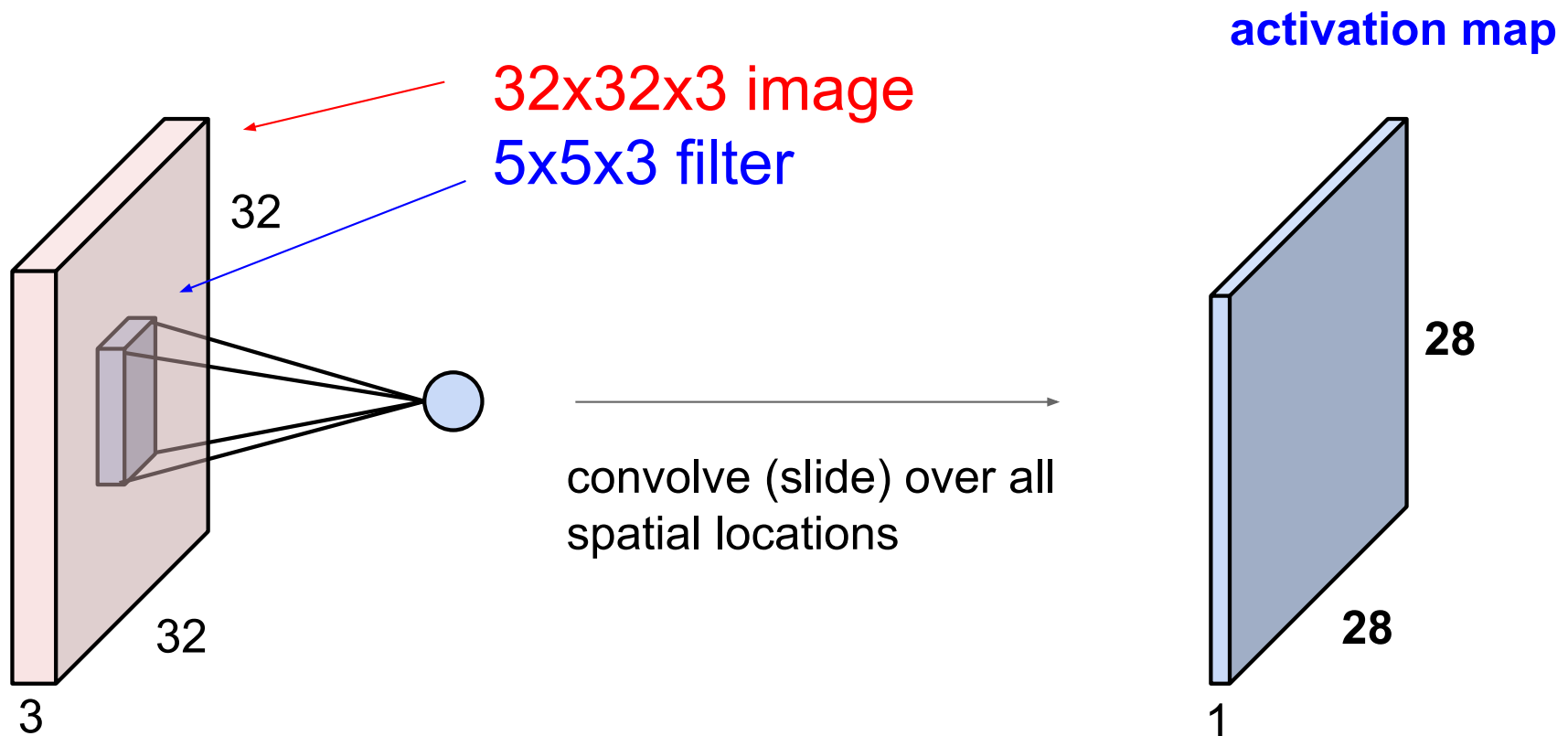
3D Visualization of CNN

<http://scs.ryerson.ca/~aharley/vis/conv/>



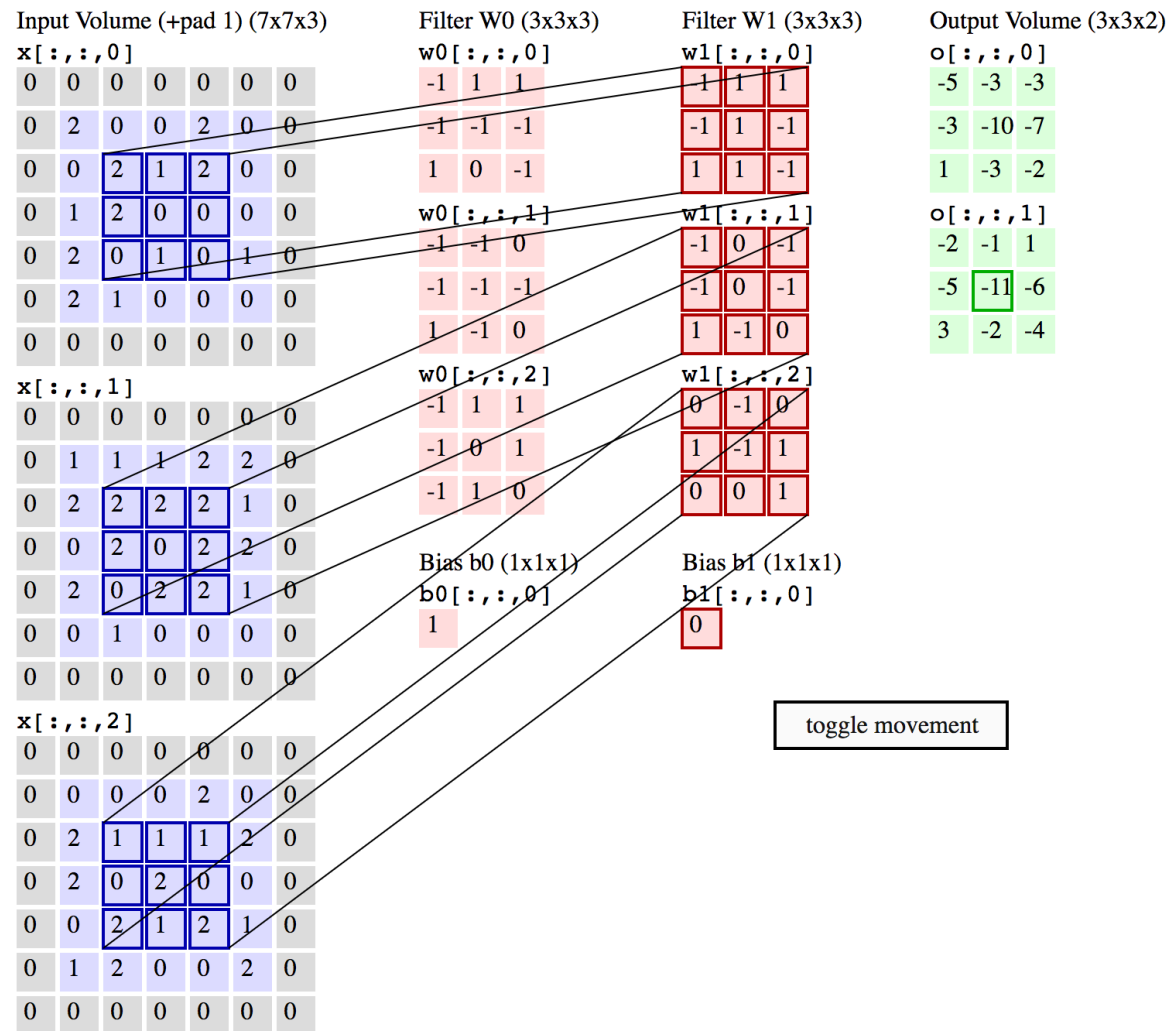
Convolution of a Color Image

- Color images consist of 3 floats per pixel for RGB (red, green blue) color values
- Convolution must also be 3-dimensional



Animation of 3D Convolution

<http://cs231n.github.io/convolutional-networks/>



MNIST Digit Recognition with CNNs (in your browser)

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

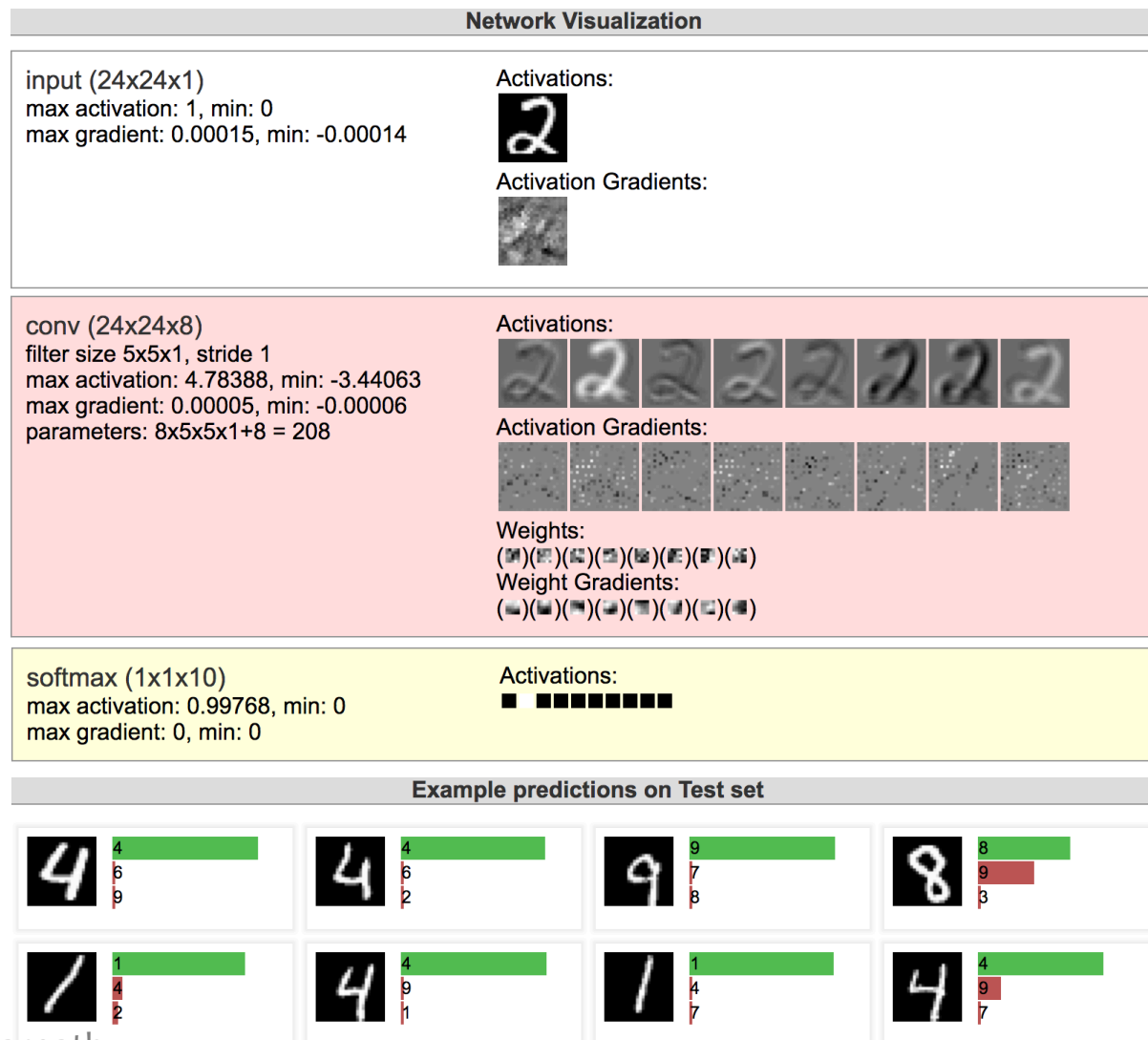


Figure from Andrej Karpathy

CNN Summary

CNNs

- Are used for all aspects of **computer vision**, and have won numerous pattern recognition competitions
- Able learn **interpretable features** at different levels of abstraction
- Typically, consist of **convolution** layers, **pooling** layers, **nonlinearities**, and **fully connected** layers

Other Resources:

- Readings on course website
- Andrej Karpathy, CS231n Notes
<http://cs231n.github.io/convolutional-networks/>

BAYESIAN NETWORKS

Bayes Nets Outline

- **Motivation**
 - Structured Prediction
- **Background**
 - Conditional Independence
 - Chain Rule of Probability
- **Directed Graphical Models**
 - Writing Joint Distributions
 - Definition: Bayesian Network
 - Qualitative Specification
 - Quantitative Specification
 - Familiar Models as Bayes Nets
- **Conditional Independence in Bayes Nets**
 - Three case studies
 - D-separation
 - Markov blanket
- **Learning**
 - Fully Observed Bayes Net
 - (Partially Observed Bayes Net)
- **Inference**
 - Sampling directly from the joint distribution
 - Gibbs Sampling

MOTIVATION: STRUCTURED PREDICTION

Structured Prediction

- Most of the models we've seen so far were for **classification**
 - Given observations: $\mathbf{x} = (x_1, x_2, \dots, x_K)$
 - Predict a (binary) **label**: y
- Many real-world problems require **structured prediction**
 - Given observations: $\mathbf{x} = (x_1, x_2, \dots, x_K)$
 - Predict a **structure**: $\mathbf{y} = (y_1, y_2, \dots, y_J)$
- Some *classification* problems benefit from **latent structure**

Structured Prediction Examples

- **Examples of structured prediction**

- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

- **Examples of latent structure**

- Object recognition

Dataset for Supervised Part-of-Speech (POS) Tagging

Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

Sample 1:	<div>n</div> <div>time</div>	<div>v</div> <div>flies</div>	<div>p</div> <div>like</div>	<div>d</div> <div>an</div>	<div>n</div> <div>arrow</div>	<div>} $y^{(1)}$</div> <div>} $x^{(1)}$</div>
Sample 2:	<div>n</div> <div>time</div>	<div>n</div> <div>flies</div>	<div>v</div> <div>like</div>	<div>d</div> <div>an</div>	<div>n</div> <div>arrow</div>	<div>} $y^{(2)}$</div> <div>} $x^{(2)}$</div>
Sample 3:	<div>n</div> <div>flies</div>	<div>v</div> <div>fly</div>	<div>p</div> <div>with</div>	<div>n</div> <div>their</div>	<div>n</div> <div>wings</div>	<div>} $y^{(3)}$</div> <div>} $x^{(3)}$</div>
Sample 4:	<div>p</div> <div>with</div>	<div>n</div> <div>time</div>	<div>n</div> <div>you</div>	<div>v</div> <div>will</div>	<div>v</div> <div>see</div>	<div>} $y^{(4)}$</div> <div>} $x^{(4)}$</div>

Dataset for Supervised Handwriting Recognition

Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$



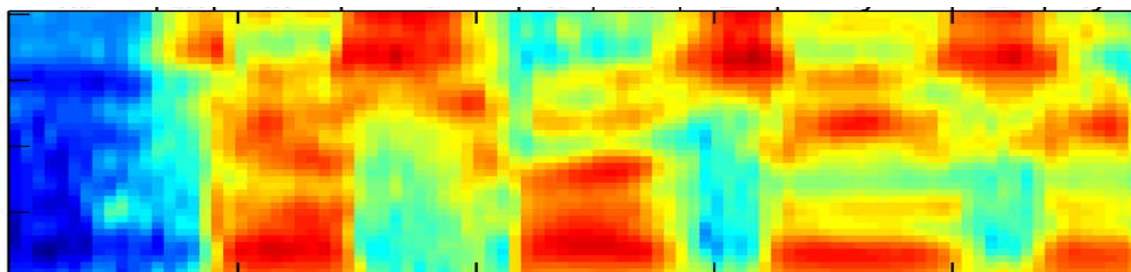
Dataset for Supervised Phoneme (Speech) Recognition

Data: $\mathcal{D} = \{x^{(n)}, y^{(n)}\}_{n=1}^N$

Sample 1:



} $y^{(1)}$

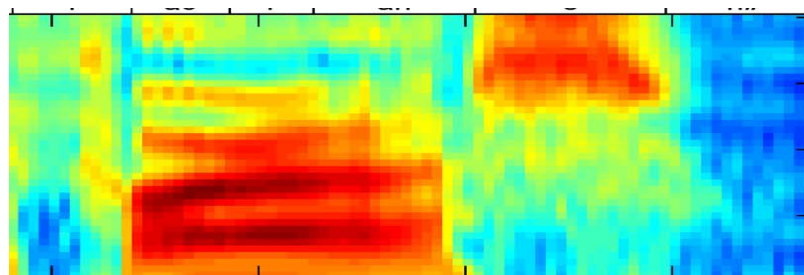


} $x^{(1)}$

Sample 2:



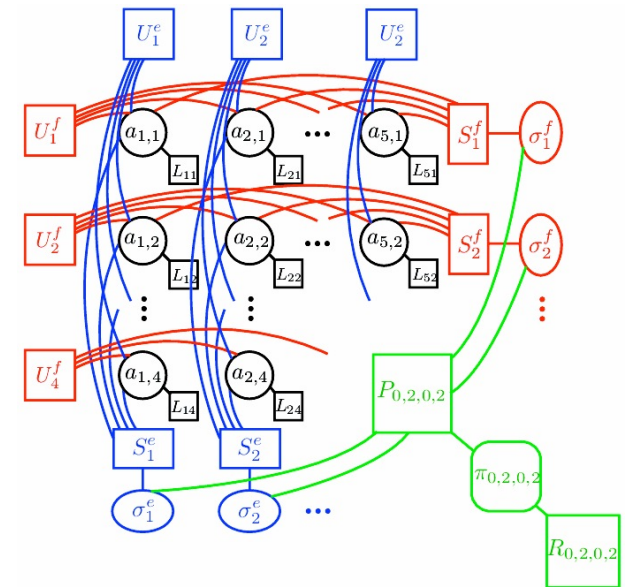
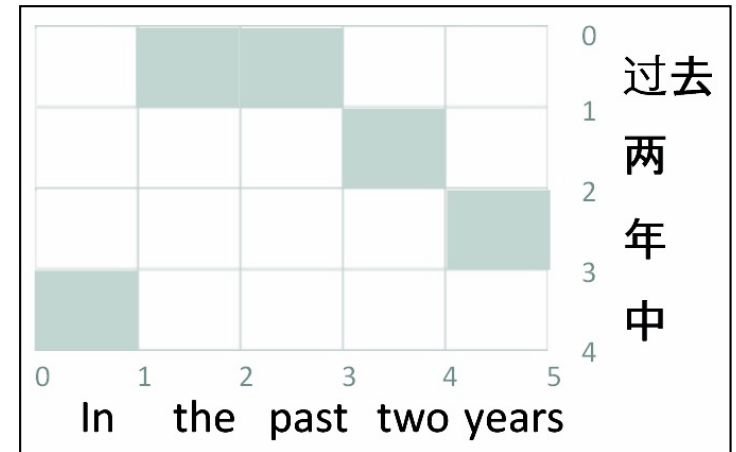
} $y^{(2)}$



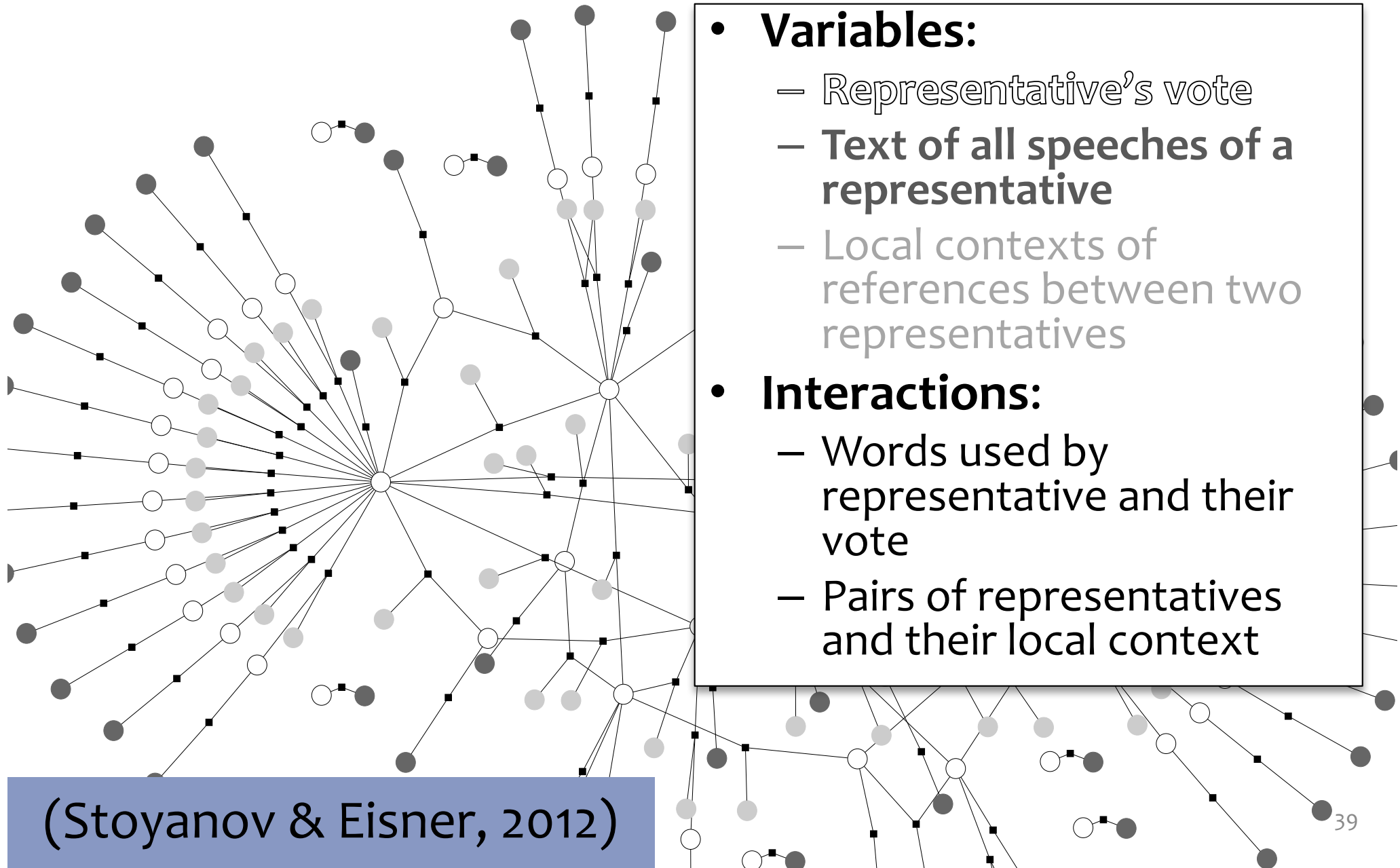
} $x^{(2)}$

Word Alignment / Phrase Extraction

- **Variables (boolean):**
 - For each (Chinese phrase, English phrase) pair, are they linked?
- **Interactions:**
 - Word fertilities
 - Few “jumps” (discontinuities)
 - Syntactic reorderings
 - “ITG constraint” on alignment
 - Phrases are disjoint (?)



Congressional Voting



Structured Prediction Examples

- **Examples of structured prediction**

- Part-of-speech (POS) tagging
- Handwriting recognition
- Speech recognition
- Word alignment
- Congressional voting

- **Examples of latent structure**

- Object recognition

Case Study: Object Recognition

Data consists of images x and labels y .



pigeon

$x^{(1)}$

$y^{(1)}$



rhinoceros

$x^{(2)}$

$y^{(2)}$



leopard

$x^{(3)}$

$y^{(3)}$



llama

$x^{(4)}$

$y^{(4)}$

Case Study: Object Recognition

Data consists of images x and labels y .

- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time

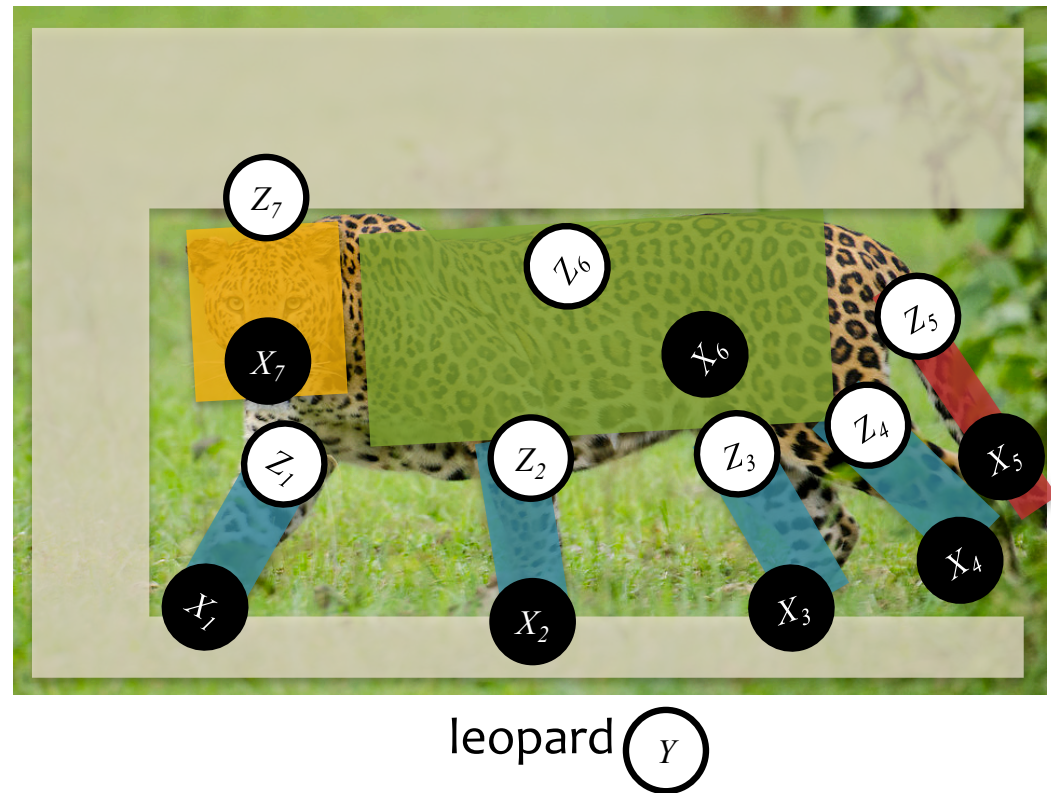


leopard

Case Study: Object Recognition

Data consists of images x and labels y .

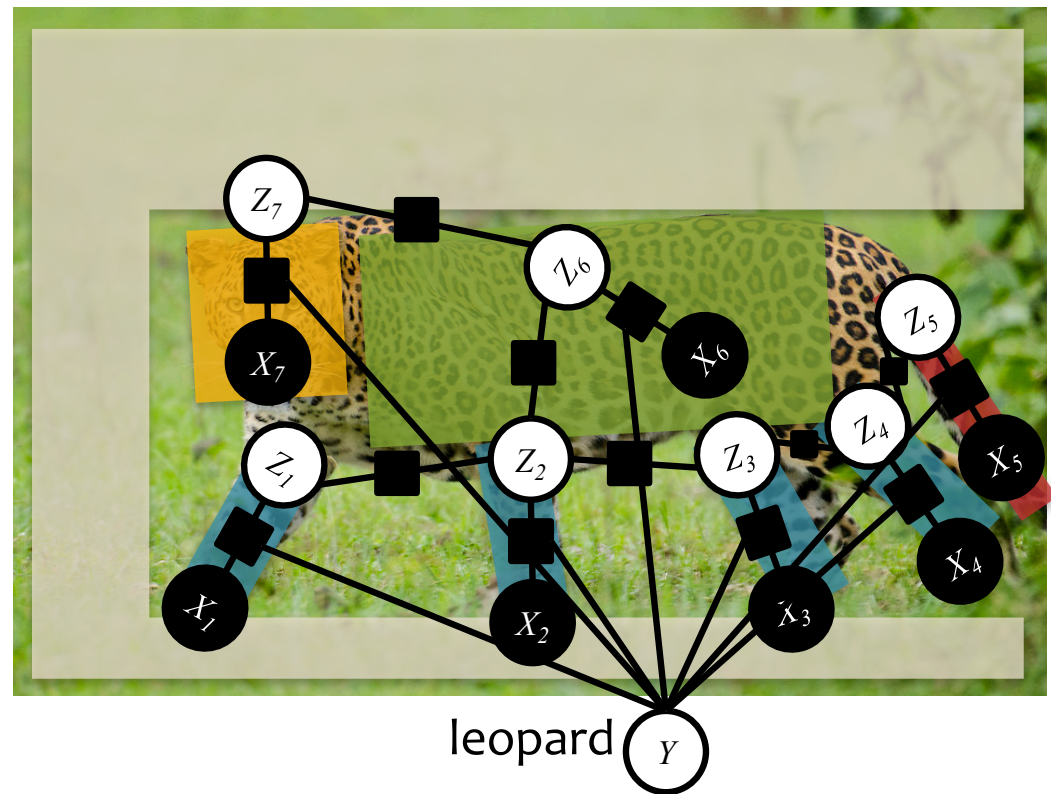
- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



Case Study: Object Recognition

Data consists of images x and labels y .

- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



Structured Prediction

Preview of challenges to come...

- Consider the task of finding the **most probable assignment** to the output

Classification

$$\hat{y} = \operatorname{argmax}_y p(y|\mathbf{x})$$

where $y \in \{+1, -1\}$

Structured Prediction

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}|\mathbf{x})$$

where $\mathbf{y} \in \mathcal{Y}$

and $|\mathcal{Y}|$ is very large

Machine Learning

The **data** inspires
the structures
we want to
predict



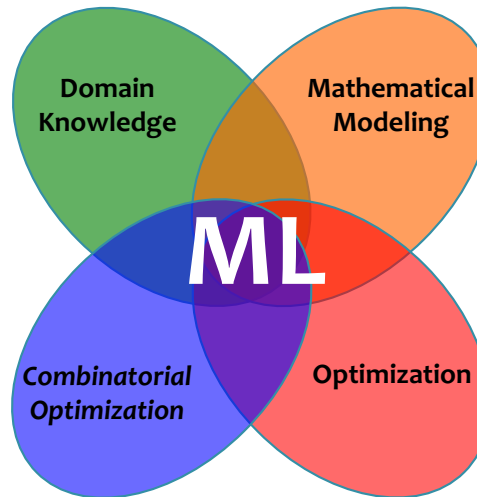
Our **model**
defines a score
for each structure

It also tells us
what to optimize



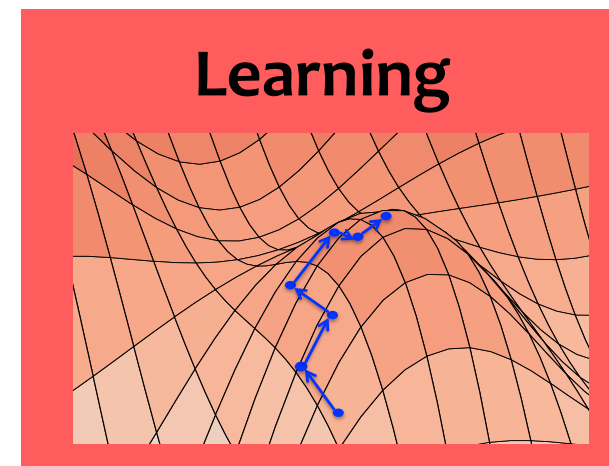
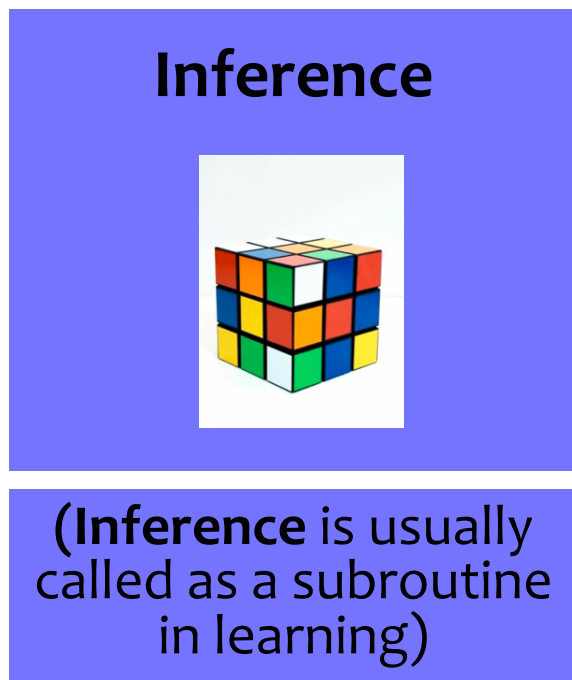
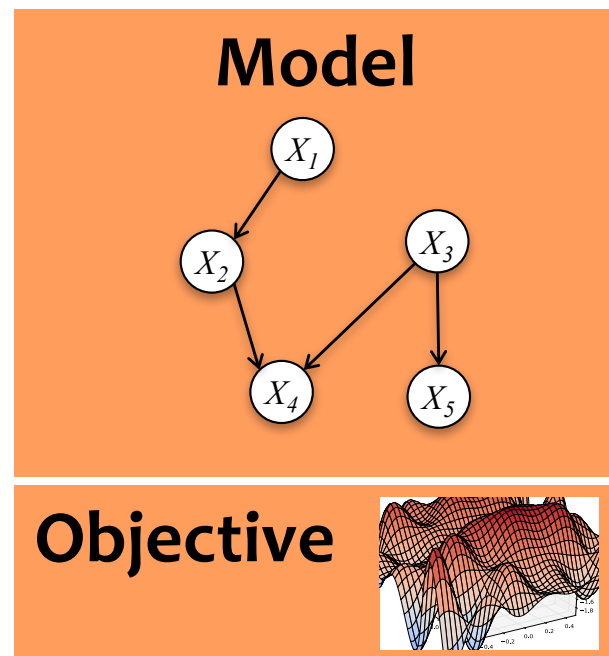
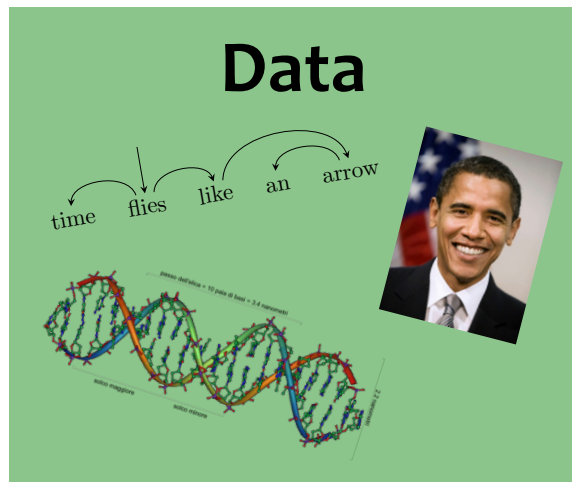
Inference finds
{best structure, marginals,
partition function} for a
new observation

(**Inference** is usually
called as a subroutine
in learning)



Learning tunes the
parameters of the
model

Machine Learning



BACKGROUND

Background

Whiteboard

- Chain Rule of Probability
- Conditional Independence

Background: Chain Rule of Probability

For random variables A and B :

$$P(A, B) = P(A|B)P(B)$$

For random variables X_1, X_2, X_3, X_4 :

$$\begin{aligned} P(X_1, X_2, X_3, X_4) = & P(X_1|X_2, X_3, X_4) \\ & P(X_2|X_3, X_4) \\ & P(X_3|X_4) \\ & P(X_4) \end{aligned}$$

Background:

Conditional Independence

Random variables A and B are conditionally independent given C if:

$$P(A, B|C) = P(A|C)P(B|C) \quad (1)$$

or equivalently:

$$P(A|B, C) = P(A|C) \quad (2)$$

We write this as:

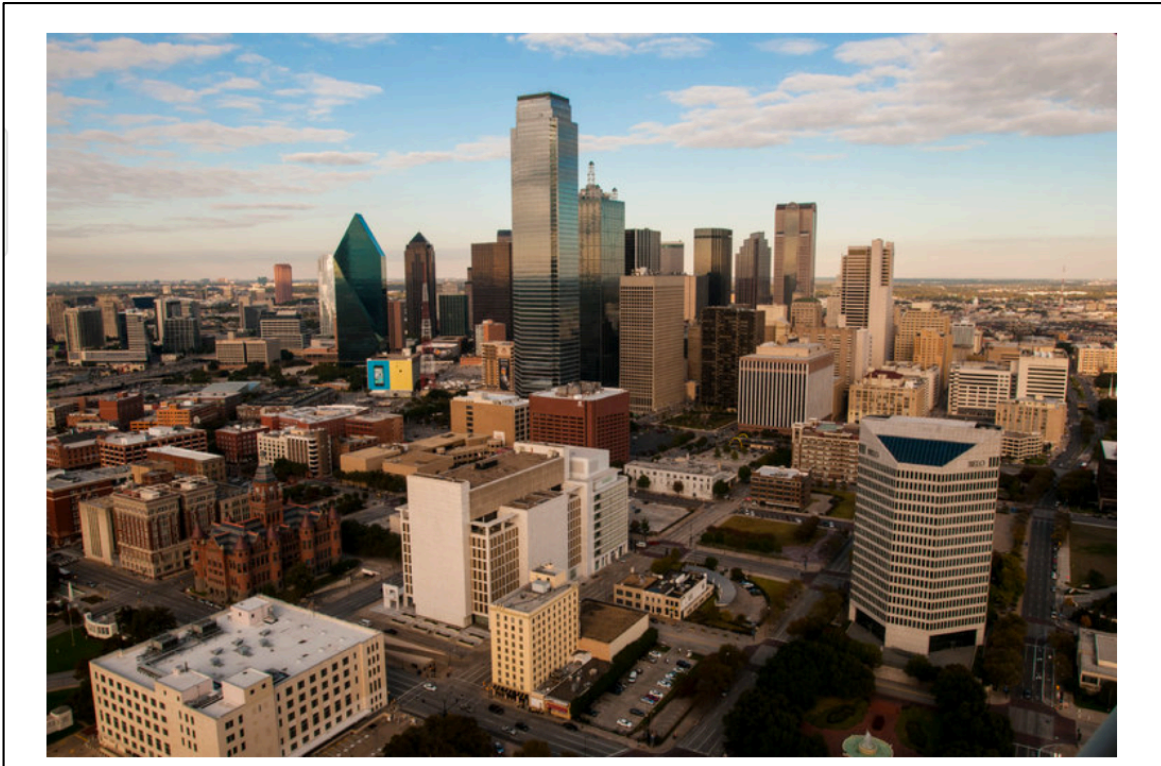
$$A \perp\!\!\!\perp B|C$$

Later we will also write: $I\langle A, \{C\}, B \rangle$

Bayesian Networks

DIRECTED GRAPHICAL MODELS

Example: Tornado Alarms



1. Imagine that you work at the 911 call center in Dallas
2. You receive six calls informing you that the Emergency Weather Sirens are going off
3. What do you conclude?

Example: Tornado Alarms

Hacking Attack Woke Up Dallas With Emergency Sirens, Officials Say

By ELI ROSENBERG and MAYA SALAM APRIL 8, 2017



Warning sirens in Dallas, meant to alert the public to emergencies like severe weather, started sounding around 11:40 p.m. Friday, and were not shut off until 1:20 a.m. Rex C. Curry for The New York Times

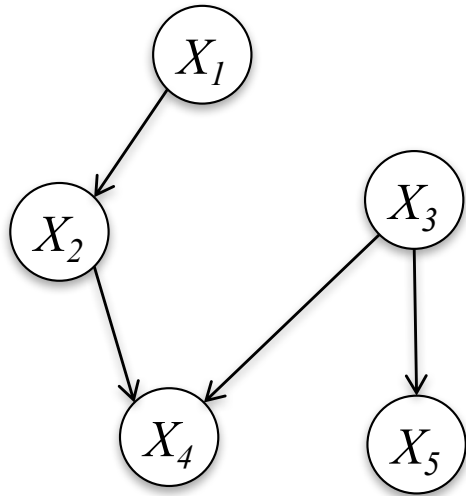
1. Imagine that you work at the 911 call center in Dallas
2. You receive six calls informing you that the Emergency Weather Sirens are going off
3. What do you conclude?

Directed Graphical Models (Bayes Nets)

Whiteboard

- Example: Tornado Alarms
- Writing Joint Distributions
 - Idea #1: Giant Table
 - Idea #2: Rewrite using chain rule
 - Idea #3: Assume full independence
 - Idea #4: Drop variables from RHS of conditionals
- Definition: Bayesian Network
- Observed Variables in Graphical Models

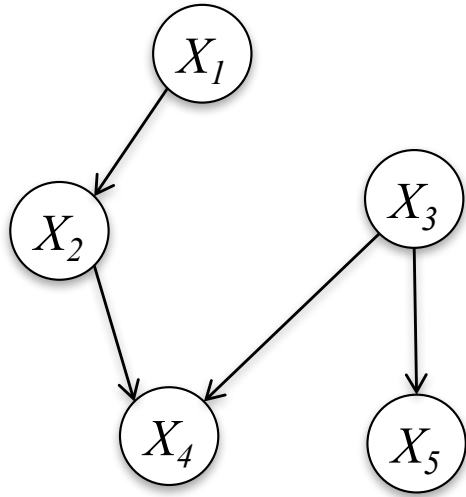
Bayesian Network



$$\begin{aligned} p(X_1, X_2, X_3, X_4, X_5) = \\ p(X_5|X_3)p(X_4|X_2, X_3) \\ p(X_3)p(X_2|X_1)p(X_1) \end{aligned}$$

Bayesian Network

Definition:



$$P(X_1 \dots X_n) = \prod_{i=1}^n P(X_i \mid \text{parents}(X_i))$$

- A Bayesian Network is a **directed graphical model**
- It consists of a graph **G** and the conditional probabilities **P**
- These two parts full specify the distribution:
 - Qualitative Specification: **G**
 - Quantitative Specification: **P**