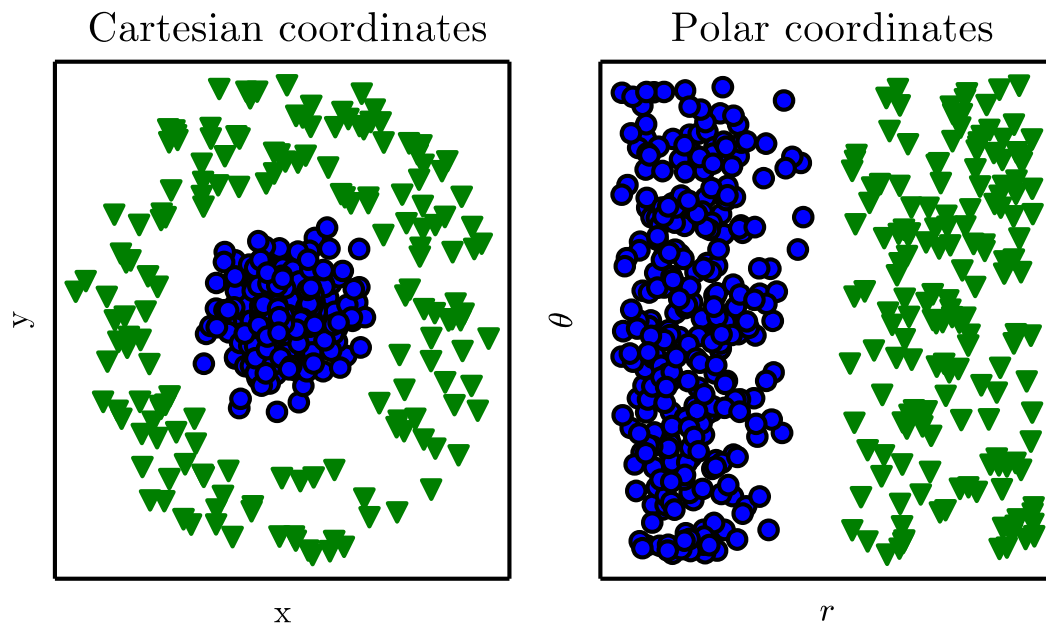# Introduction to Deep Learning
# CMPT 733

Steven Bergner

# Overview

- Renaissance of artificial neural networks
  - Representation learning vs feature engineering

- Background
  - Linear Algebra, Optimization
  - Regularization

- Construction and training of layered learners
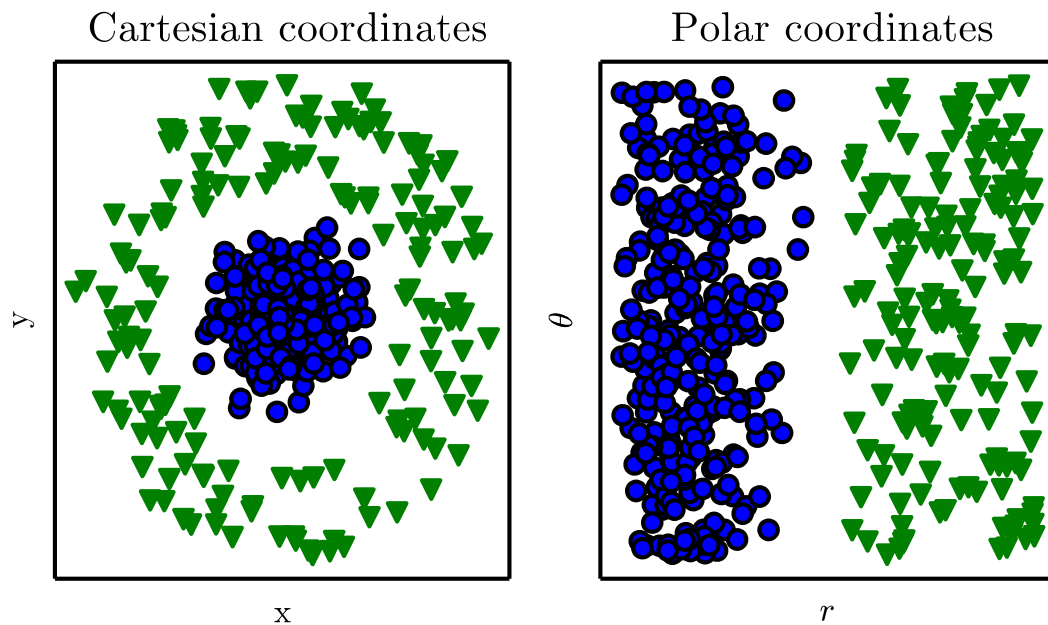
- Frameworks for deep learning

# Representations matter

Cartesian coordinates

Polar coordinates

- Transform into the right representation
- Classify points simply by threshold on radius axis

[Goodfellow, Bengio, Courville 2016]

# Representations matter

Cartesian coordinates

Polar coordinates
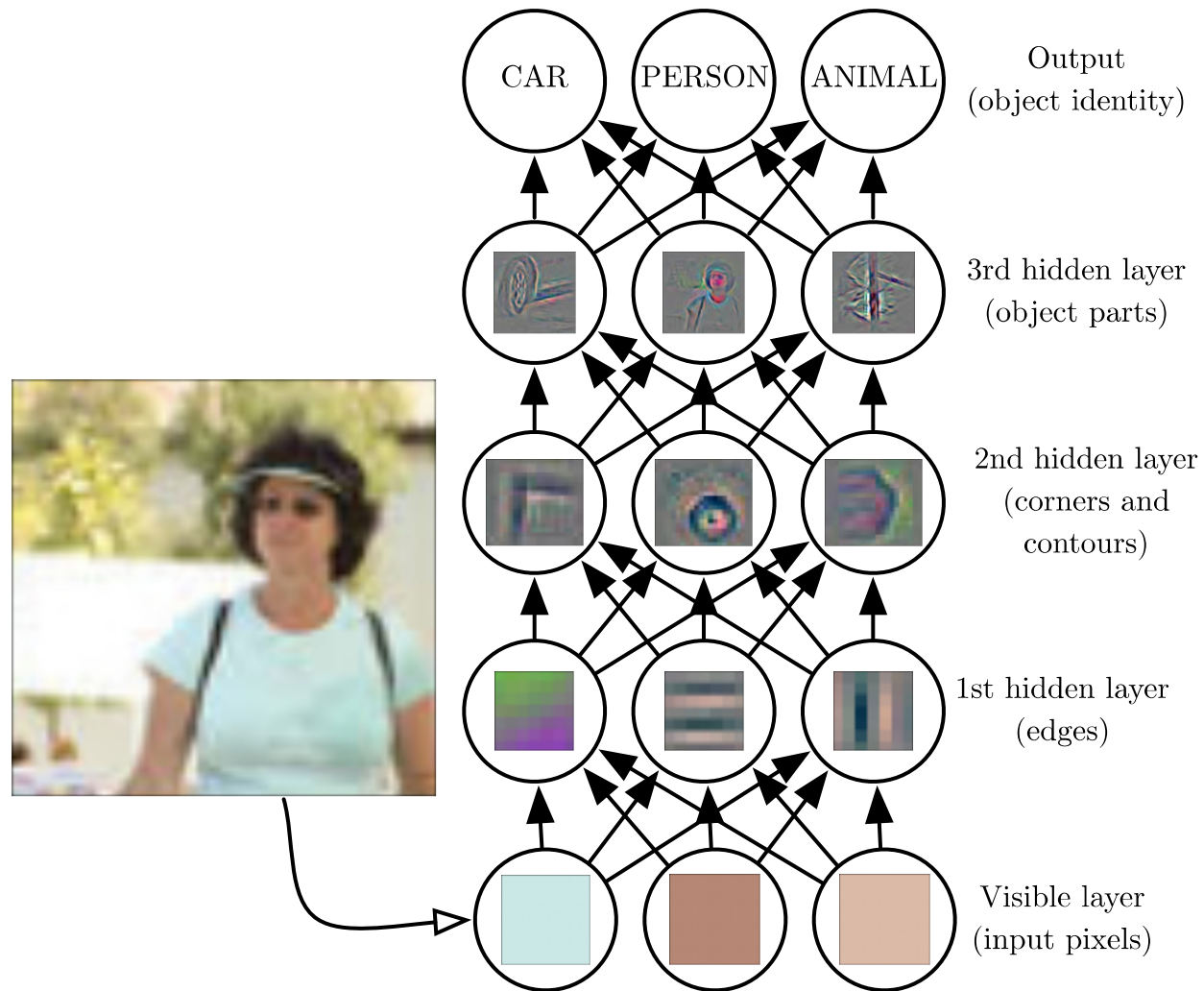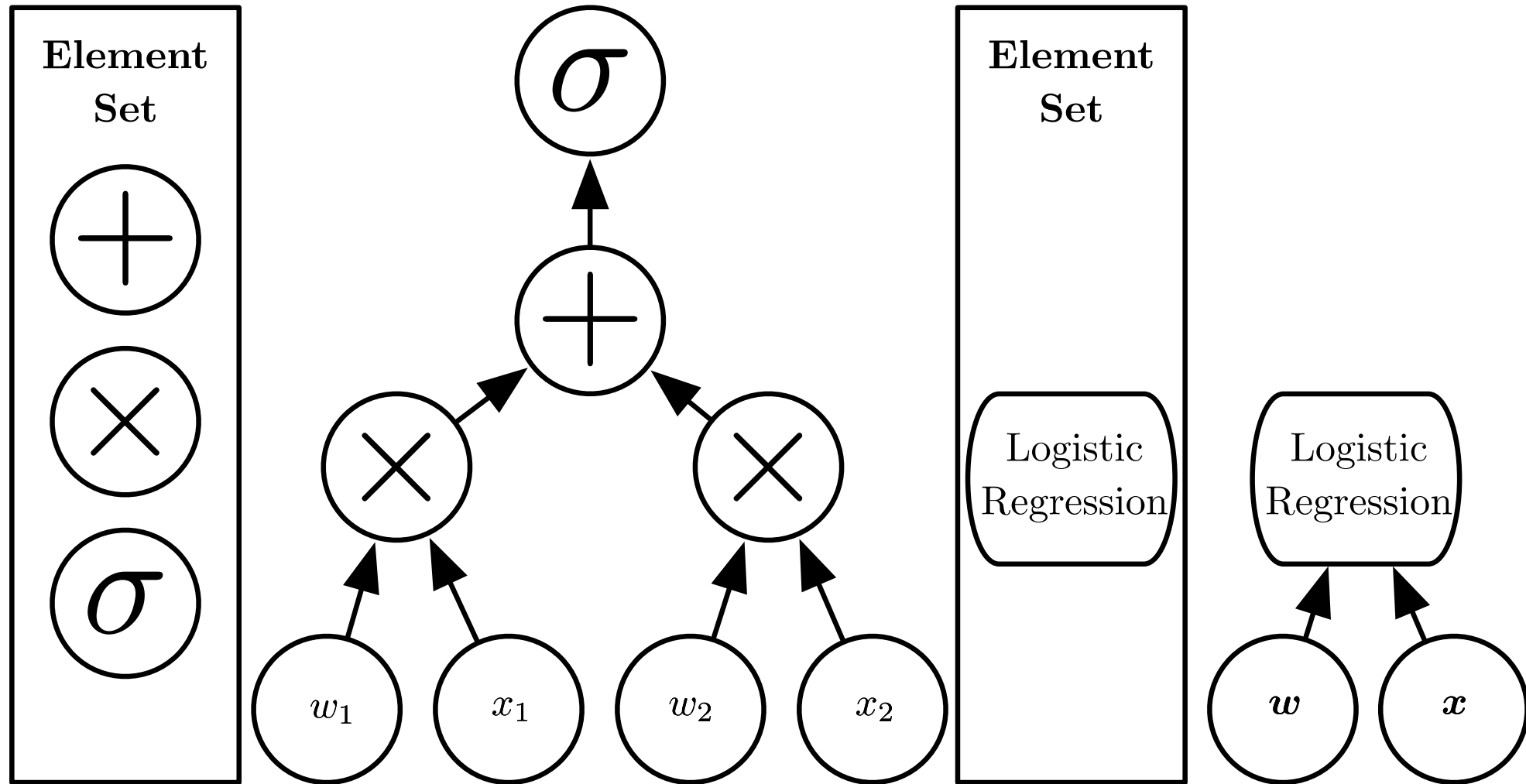
- Transform into the right representation

- Classify points simply by threshold on radius axis

- Single neuron with non-linearity can do this

[Goodfellow, Bengio, Courville 2016]

4

# Depth: layered composition



CAR    PERSON    ANIMAL    Output (object identity)

3rd hidden layer (object parts)

2nd hidden layer (corners and contours)

1st hidden layer (edges)

Visible layer (input pixels)

[Goodfellow, Bengio, Courville 2016]
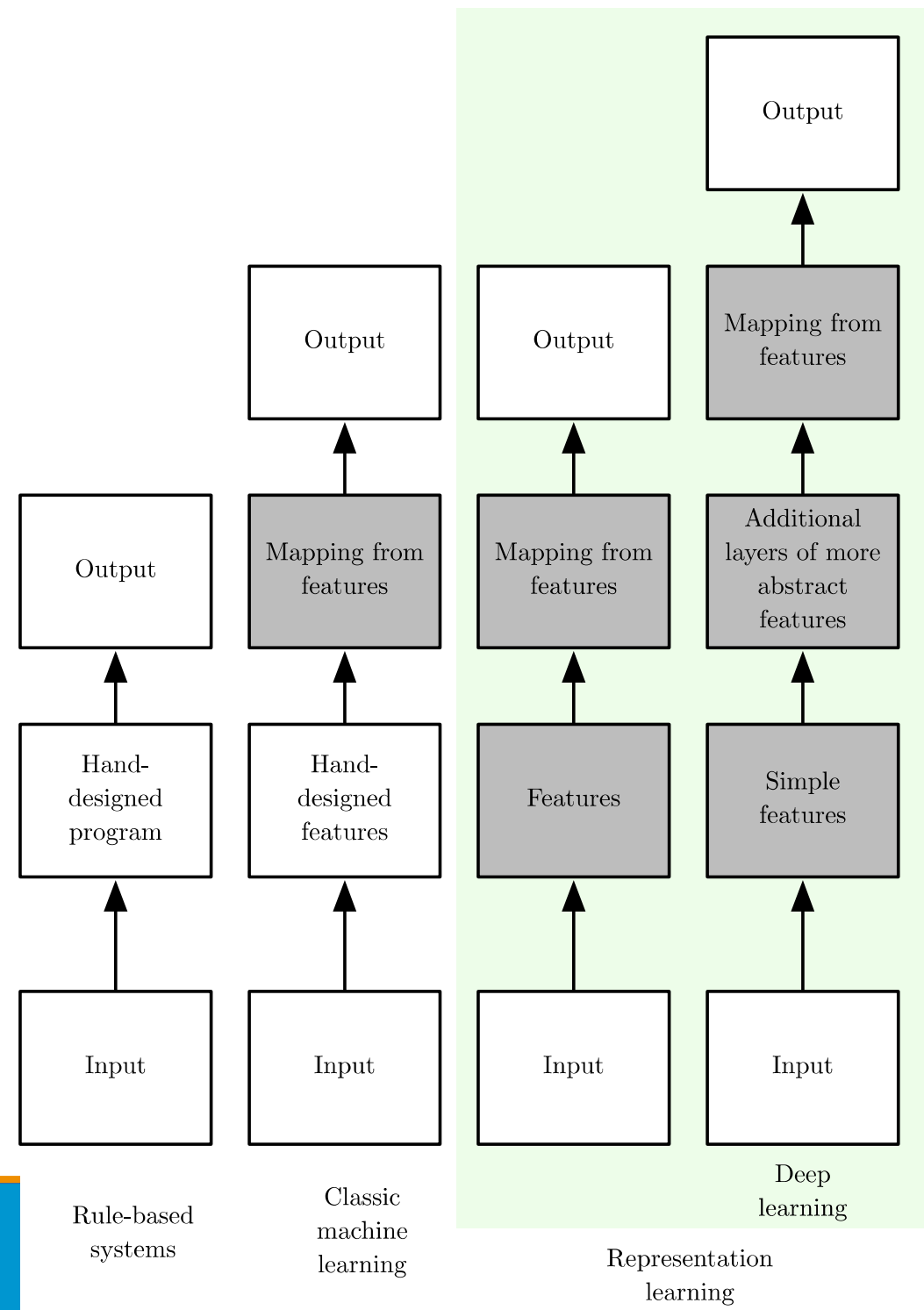
# Computational graph

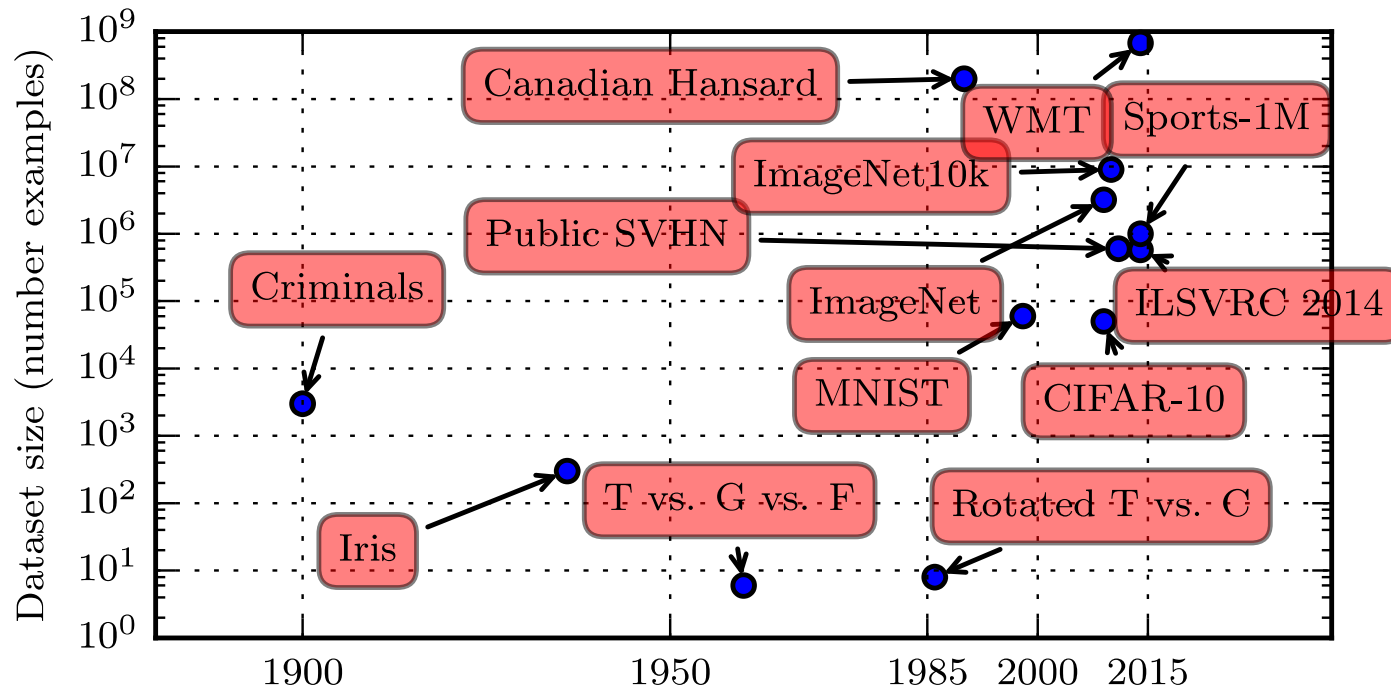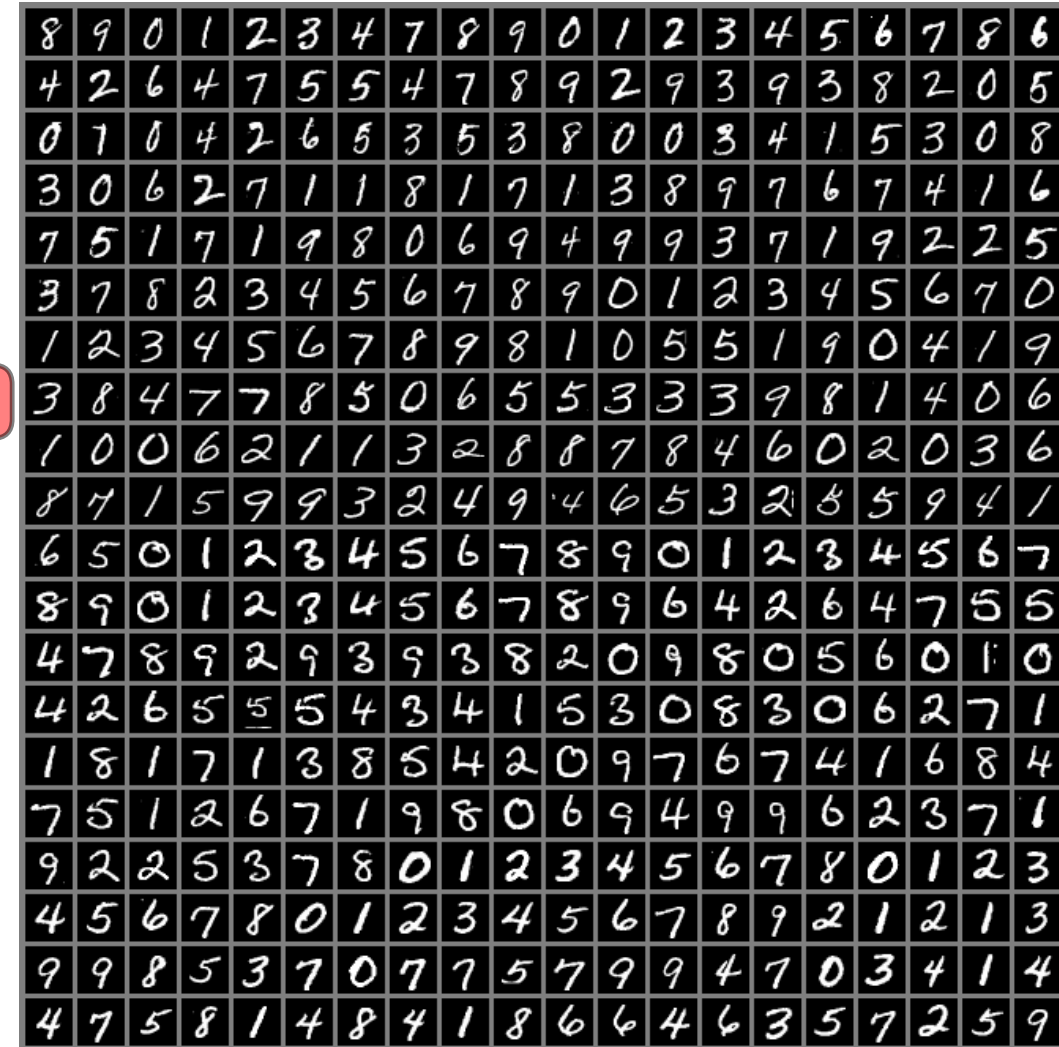[Goodfellow, Bengio, Courville 2016]

# Components of learning

- Hand designed program
  - Input $\rightarrow$ Output
- Increasingly automated
  - Simple features
  - Abstract features
  - Mapping from features

[Goodfellow, Bengio, Courville 2016]

# Growing Dataset Size



MNIST dataset

# Basics
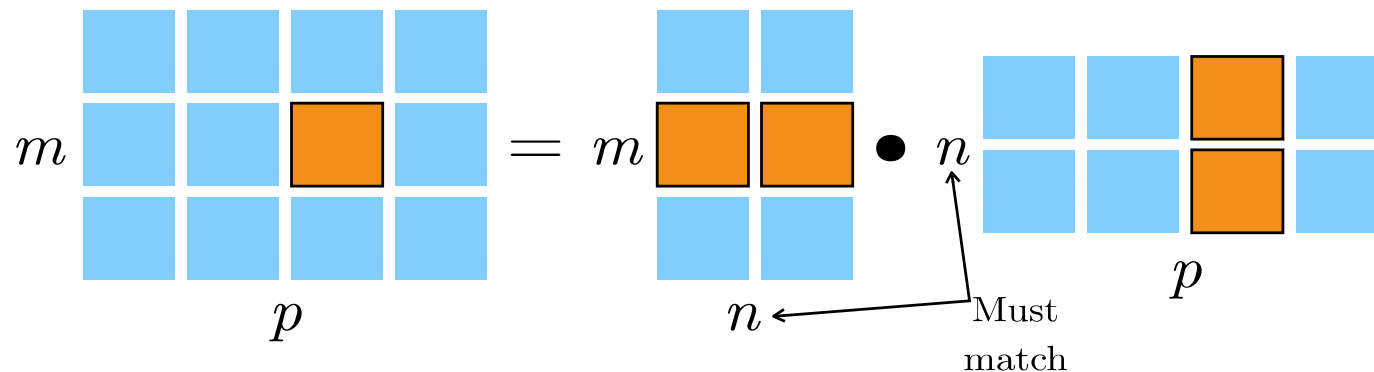
Linear Algebra and Optimization

# Linear Algebra

- Tensor is an array of numbers
    - Multi-dim: 0d scalar, 1d vector, 2d matrix/image, 3d RGB image
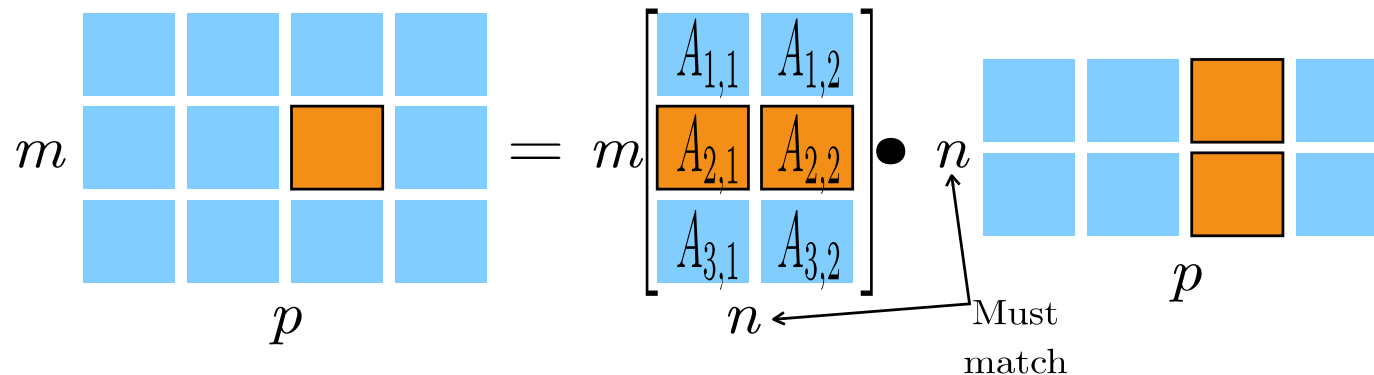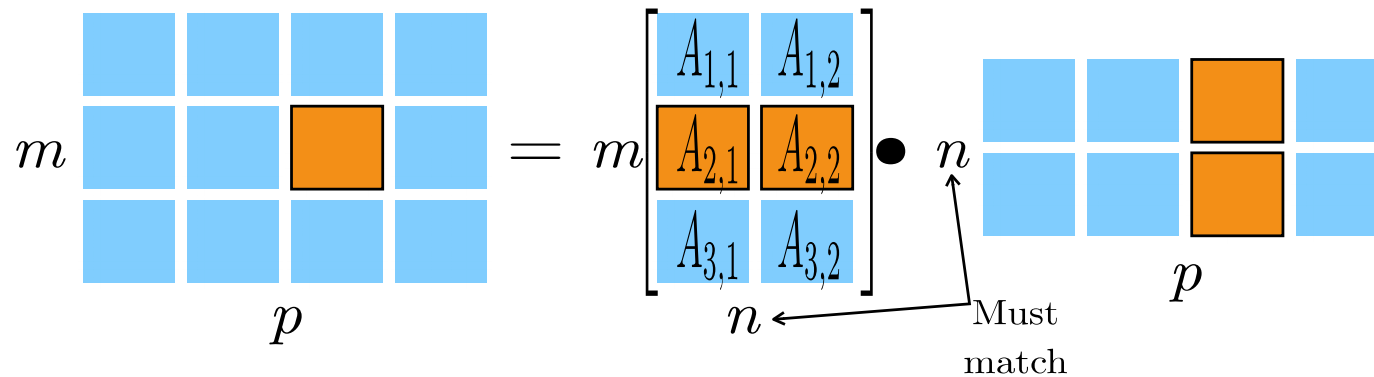- Matrix (dot) product $\quad C = AB \qquad C_{i,j} = \sum_k A_{i,k} B_{k,j}$

- Dot product of vectors A and B

    (m = p = 1 in above notation, n=2)

# Linear Algebra

- Tensor is an array of numbers

  - Multi-dim: 0d scalar, 1d vector, 2d matrix/image, 3d RGB image

- Matrix (dot) product $\boldsymbol{C} = \boldsymbol{AB}$ $\qquad C_{i,j} = \sum_k A_{i,k} B_{k,j}$



- Dot product of vectors A and B

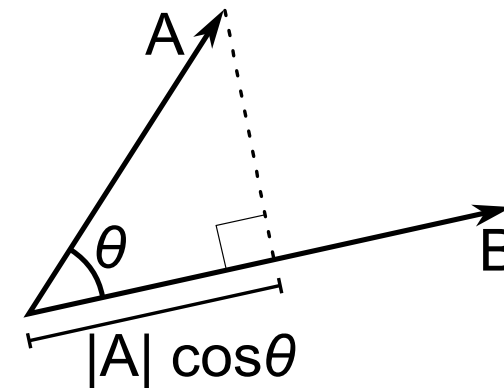  (m = p = 1 in above notation, n=2)

[Goodfellow, Bengio, Courville 2016]

# Linear Algebra

- Tensor is an array of numbers
  - Multi-dim: 0d scalar, 1d vector, 2d matrix/image, 3d RGB image

- Matrix (dot) product $\quad \boldsymbol{C} = \boldsymbol{AB} \qquad C_{i,j} = \sum_k A_{i,k} B_{k,j}$

$$m \begin{array}{|c|c|c|c|}\hline & & & \\\hline & & & \\\hline & & & \\\hline\end{array} \underset{p}{} = m \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix}_{n} \bullet \; n \begin{array}{|c|c|c|c|}\hline & & & \\\hline & & & \\\hline\end{array} \underset{p}{}$$

Must match

- Dot product of vectors A and B

  (m = p = 1 in above notation, n=2)

[Goodfellow, Bengio, Courville 2016]

# Linear Algebra

- Tensor is an array of numbers

  – Multi-dim: 0d scalar, 1d vector, 2d matrix/image, 3d RGB image

- Matrix (dot) product $\qquad C = AB \qquad C_{i,j} = \sum_k A_{i,k} B_{k,j}$



Must match

- Dot product of vectors A and B

  (m = p = 1 in above notation, n=2)

$|A| \cos\theta$

# Linear algebra: Norms

- $L^p$ norm

$$||\boldsymbol{x}||_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- Most popular norm: L2 norm, $p{=}2$

- L1 norm, $p{=}1$:  $||\boldsymbol{x}||_1 = \sum_i |x_i|.$

- Max norm, infinite $p$:  $||\boldsymbol{x}||_\infty = \max_i |x_i|.$

# Nonlinearities

- ReLU

- Softplus

- Logistic Sigmoid

[(c) public domain]

[Goodfellow, Bengio, Courville 2016]

# Approximate Optimization



This local minimum performs nearly as well as the global one, so it is an acceptable halting point.

Ideally, we would like to arrive at the global minimum, but this might not be possible.

This local minimum performs poorly and should be avoided.

$f(x)$

$x$

[Goodfellow, Bengio, Courville 2016]

# Gradient descent



Global minimum at $x = 0$.
Since $f'(x) = 0$, gradient
descent halts here.

For $x < 0$, we have $f'(x) < 0$,
so we can decrease $f$ by
moving rightward.

For $x > 0$, we have $f'(x) > 0$,
so we can decrease $f$ by
moving leftward.

$f(x) = \frac{1}{2}x^2$

$f'(x) = x$

# Critical points

| Minimum | Maximum | Saddle point |
|---------|---------|--------------|

[Goodfellow, Bengio, Courville 2016]

# Critical points

Minimum | Maximum | Saddle point



Saddle point – 1$^{st}$ and 2$^{nd}$ derivative vanish

[Goodfellow, Bengio, Courville 2016]
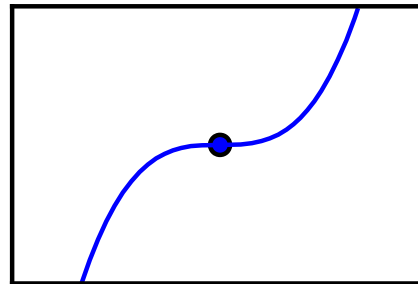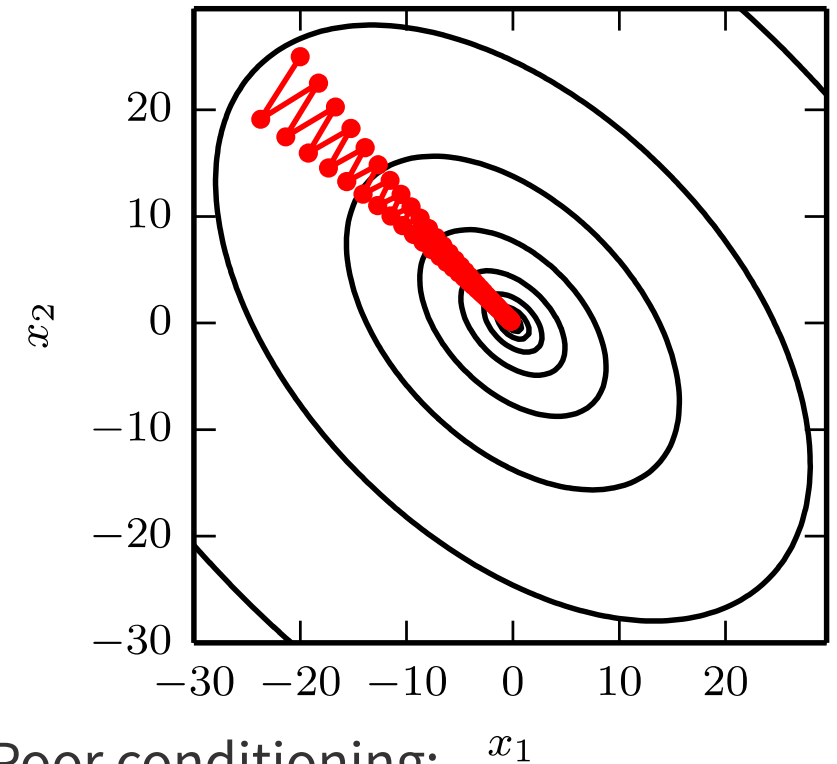
# Critical points

Minimum          Maximum          Saddle point



Saddle point – 1$^{st}$ and 2$^{nd}$ derivative vanish

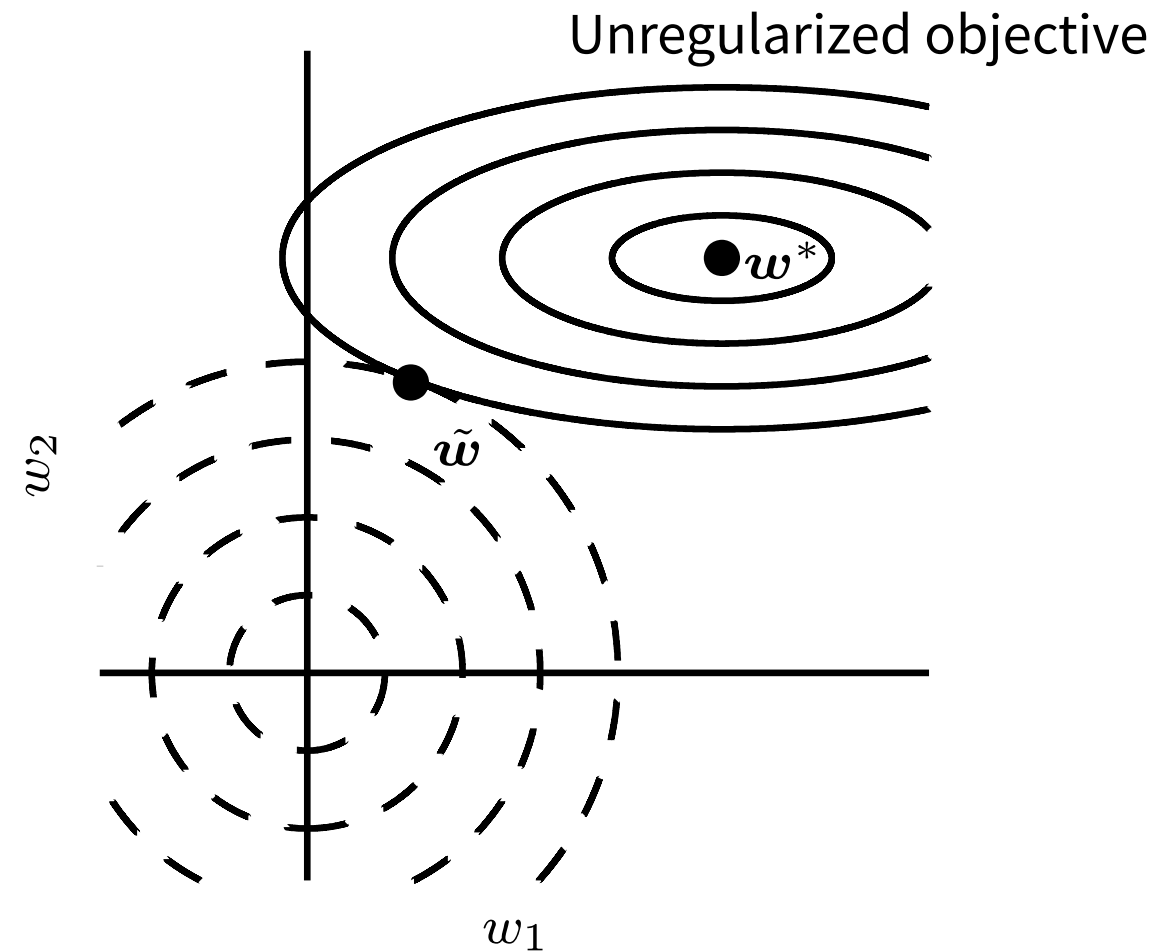Poor conditioning:
1$^{st}$ deriv large in one and small in another direction

[Goodfellow, Bengio, Courville 2016]

# Tensorflow Playground

- http://playground.tensorflow.org/
  - Try out simple network configurations


- https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html
  - Visualize linear and non-linear mappings

# Regularization

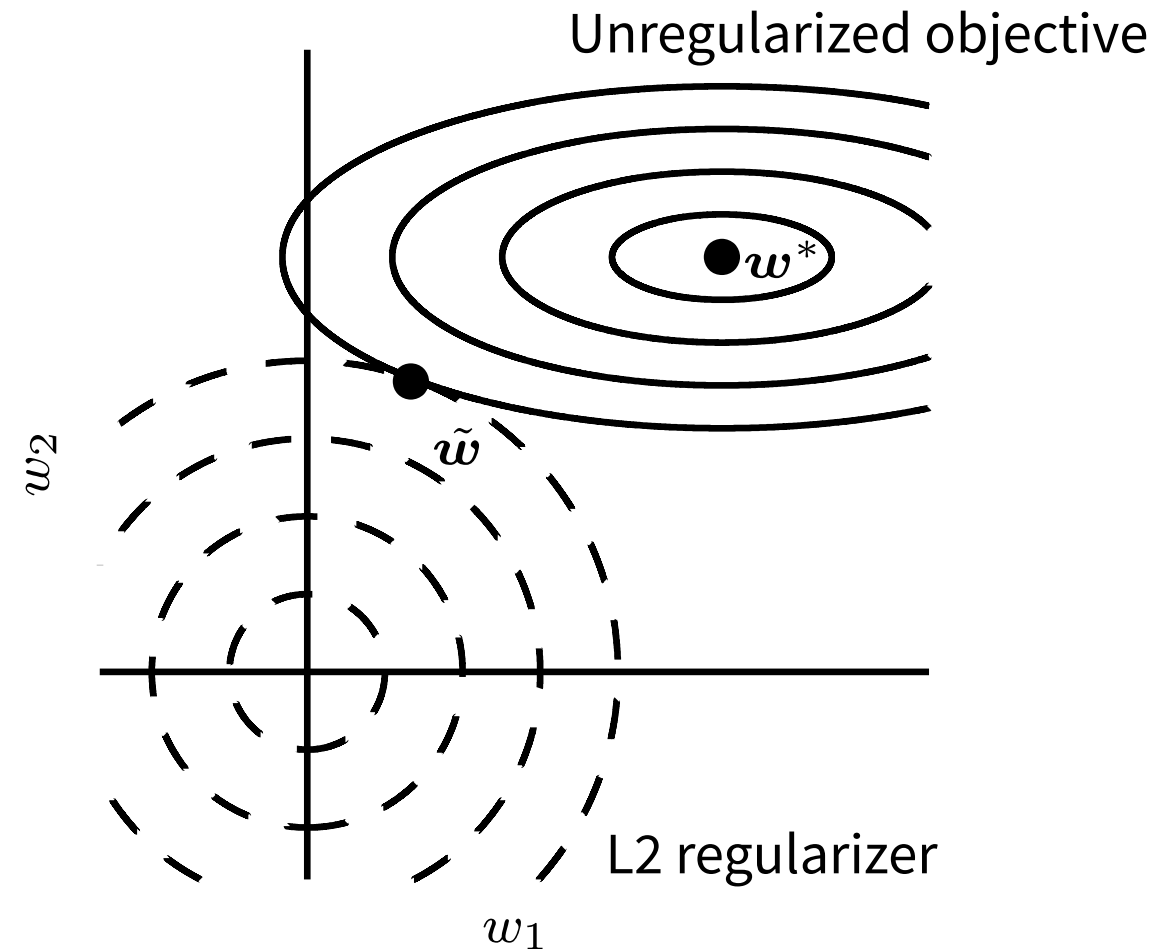Reduced generalization error without impacting training error

# Constrained optimization



Unregularized objective

$w^*$

$\tilde{w}$

$w_2$

$w_1$

[Goodfellow, Bengio, Courville 2016]

# Constrained optimization

- Squared L2 encourages small weights



Unregularized objective

$w^*$

$w_2$

$\tilde{w}$

L2 regularizer

$w_1$

[Goodfellow, Bengio, Courville 2016]

# Constrained optimization

- Squared L2 encourages small weights

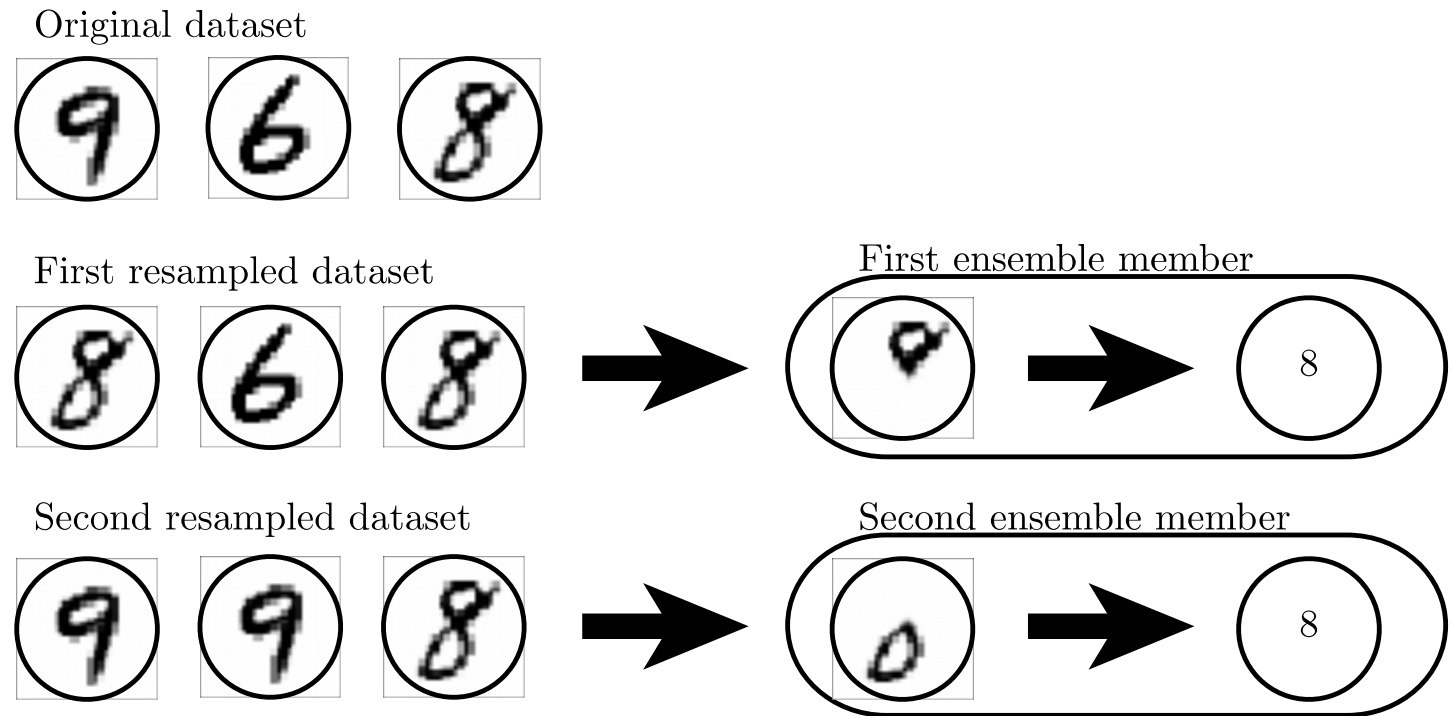- L1 encourages sparsity of model parameters (weights)



Unregularized objective

$\bullet w^*$

$\tilde{w}$

$w_2$

$w_1$

L2 regularizer

# Dataset augmentation

Affine Distortion

Noise

Elastic Deformation

Horizontal flip

Random Translation

Hue Shift

[Goodfellow, Bengio, Courville 2016]

# Learning curves

# Learning curves



- Early stopping before validation error starts to increase

# Bagging
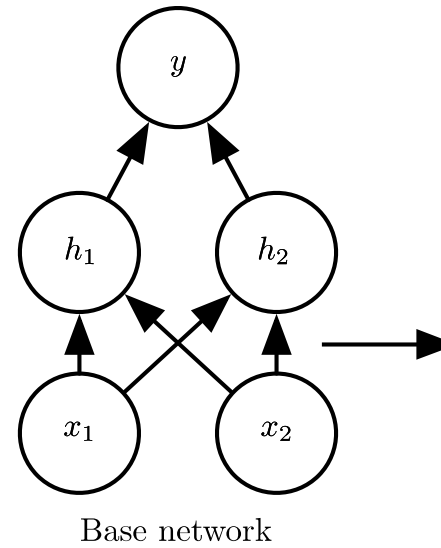
- Average multiple models trained on subsets of the data

# Bagging

- Average multiple models trained on subsets of the data
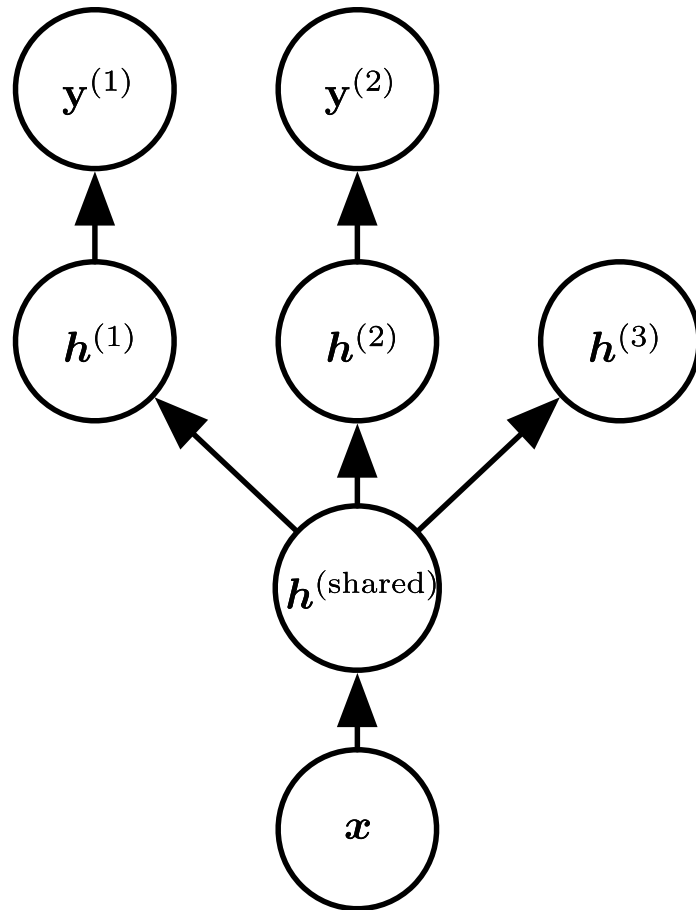- First subset: learns top loop, Second subset: bottom loop

Original dataset

First resampled dataset

First ensemble member

8

Second resampled dataset

Second ensemble member

8

# Dropout

- Random sample of connection weights is set to zero

- Train different network model each time

- Learn more robust, generalizable features

Base network

Ensemble of subnetworks

[Goodfellow, Bengio, Courville 2016]

# Multitask learning



- Shared parameters are trained with more data

- Improved generalization error due to increased statistical strength

[Goodfellow, Bengio, Courville 2016]

# Components of popular architectures
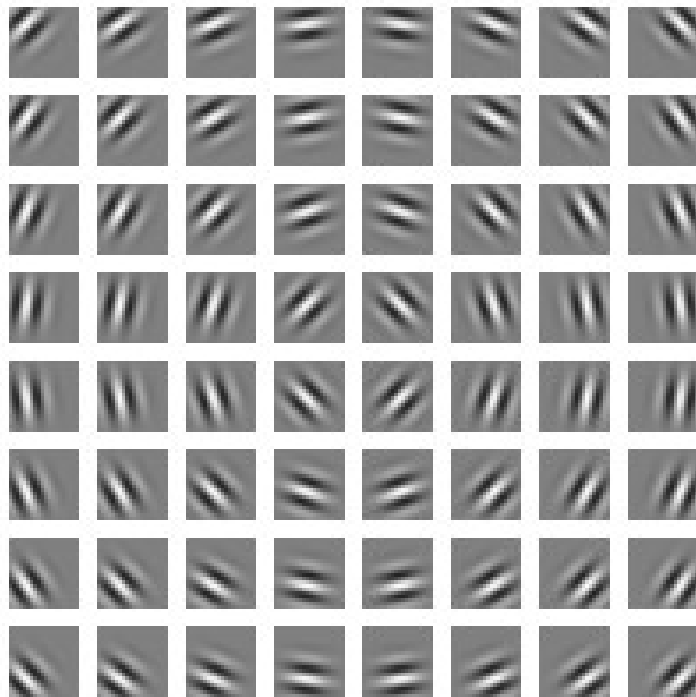
# Convolution as edge detector



Input



Output

| 1 | -1 |
|---|---|

Kernel

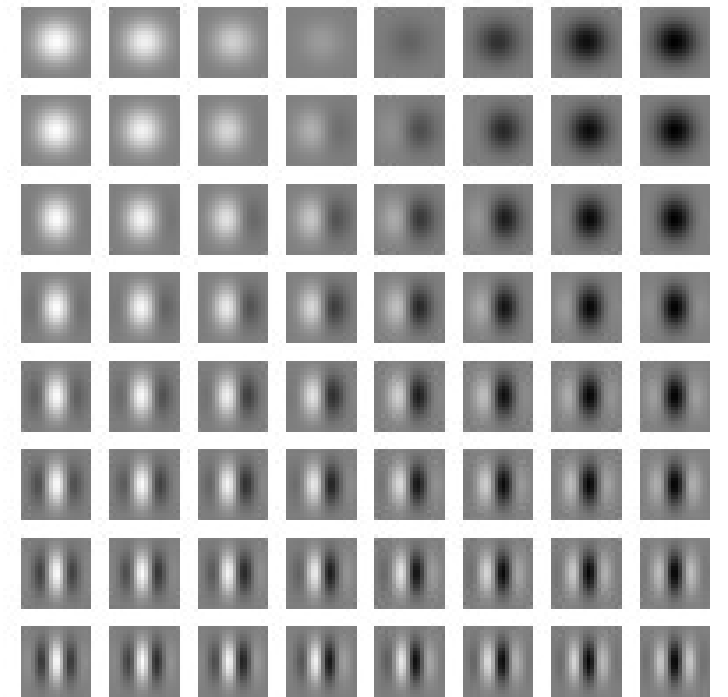[Goodfellow, Bengio, Courville 2016]

# Gabor wavelets (kernels)



[Goodfellow, Bengio, Courville 2016]

35

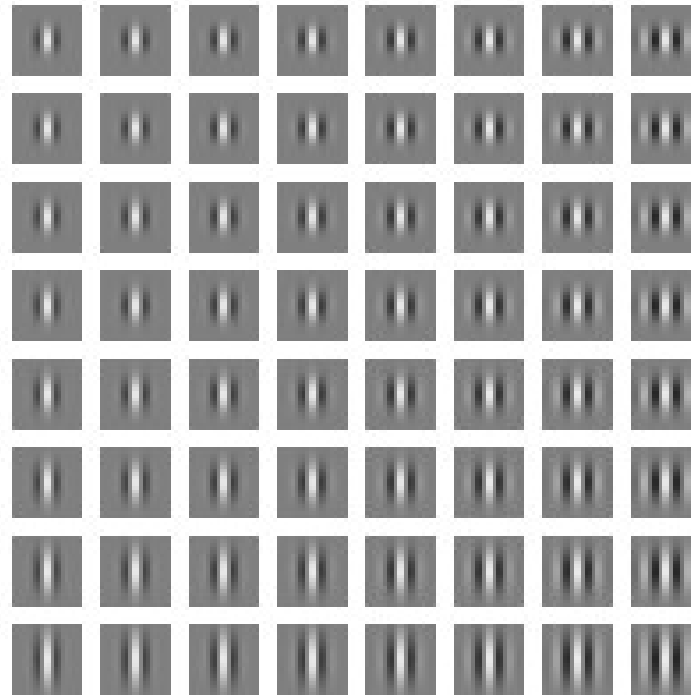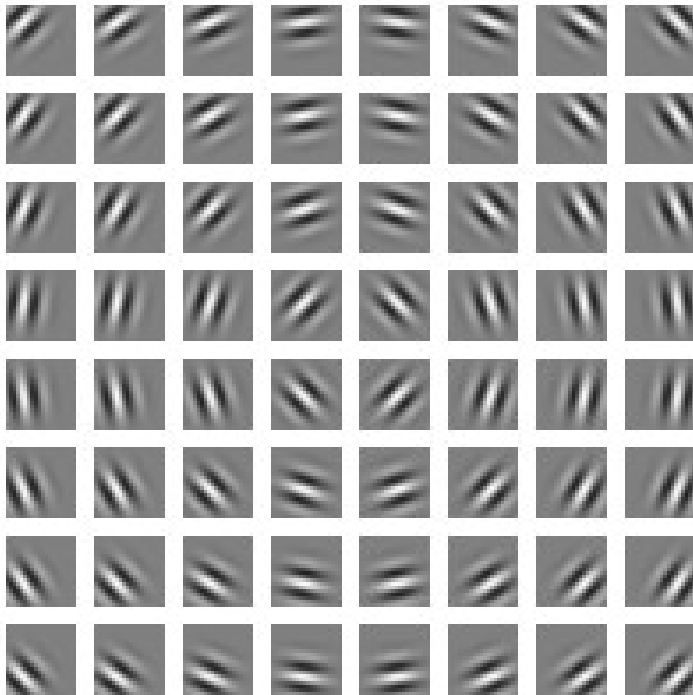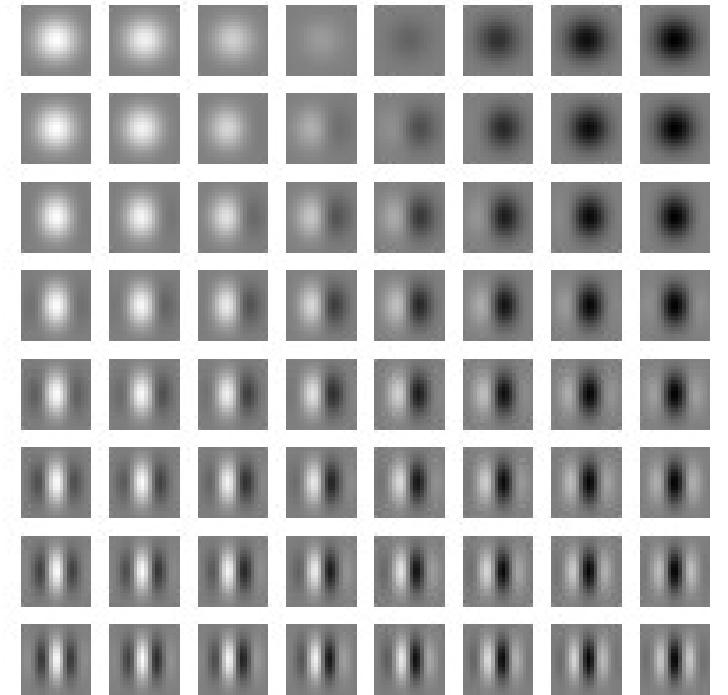# Gabor wavelets (kernels)



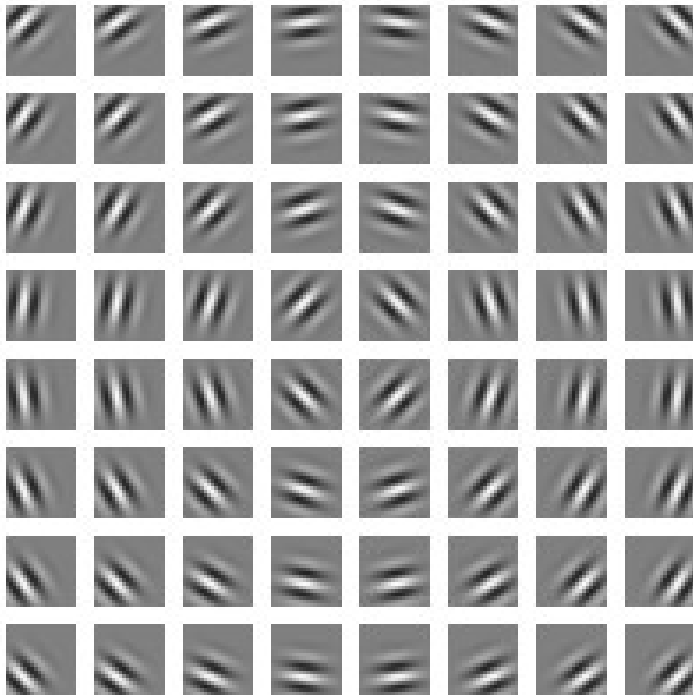Local average, first derivative
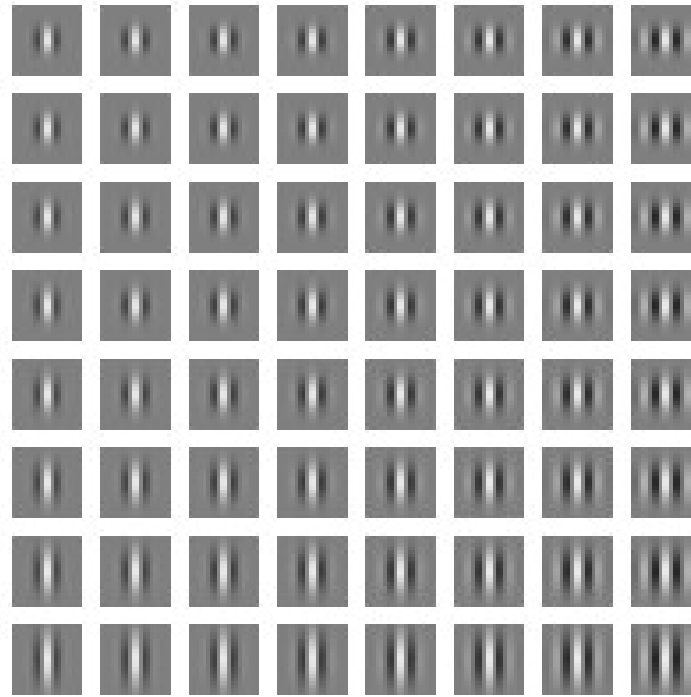
# Gabor wavelets (kernels)

Second derivative (curvature)          Local average, first derivative

[Goodfellow, Bengio, Courville 2016]
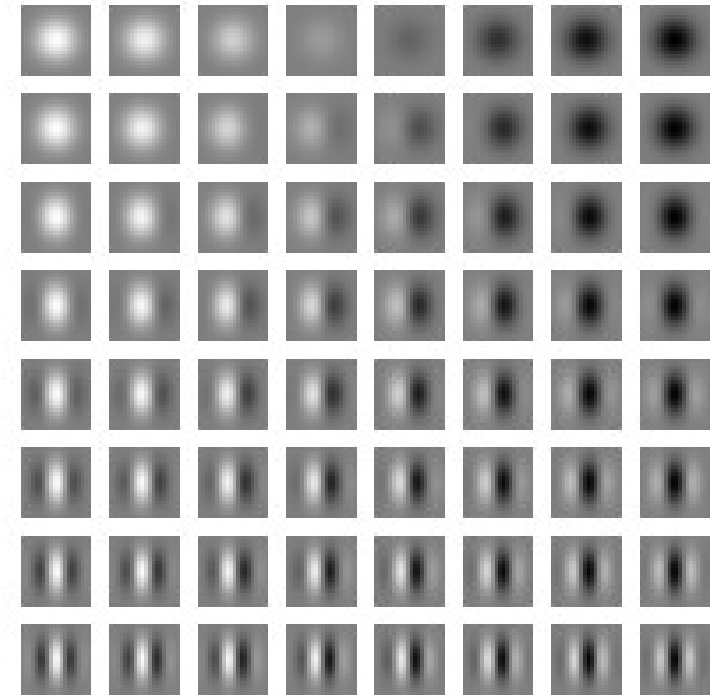
# Gabor wavelets (kernels)



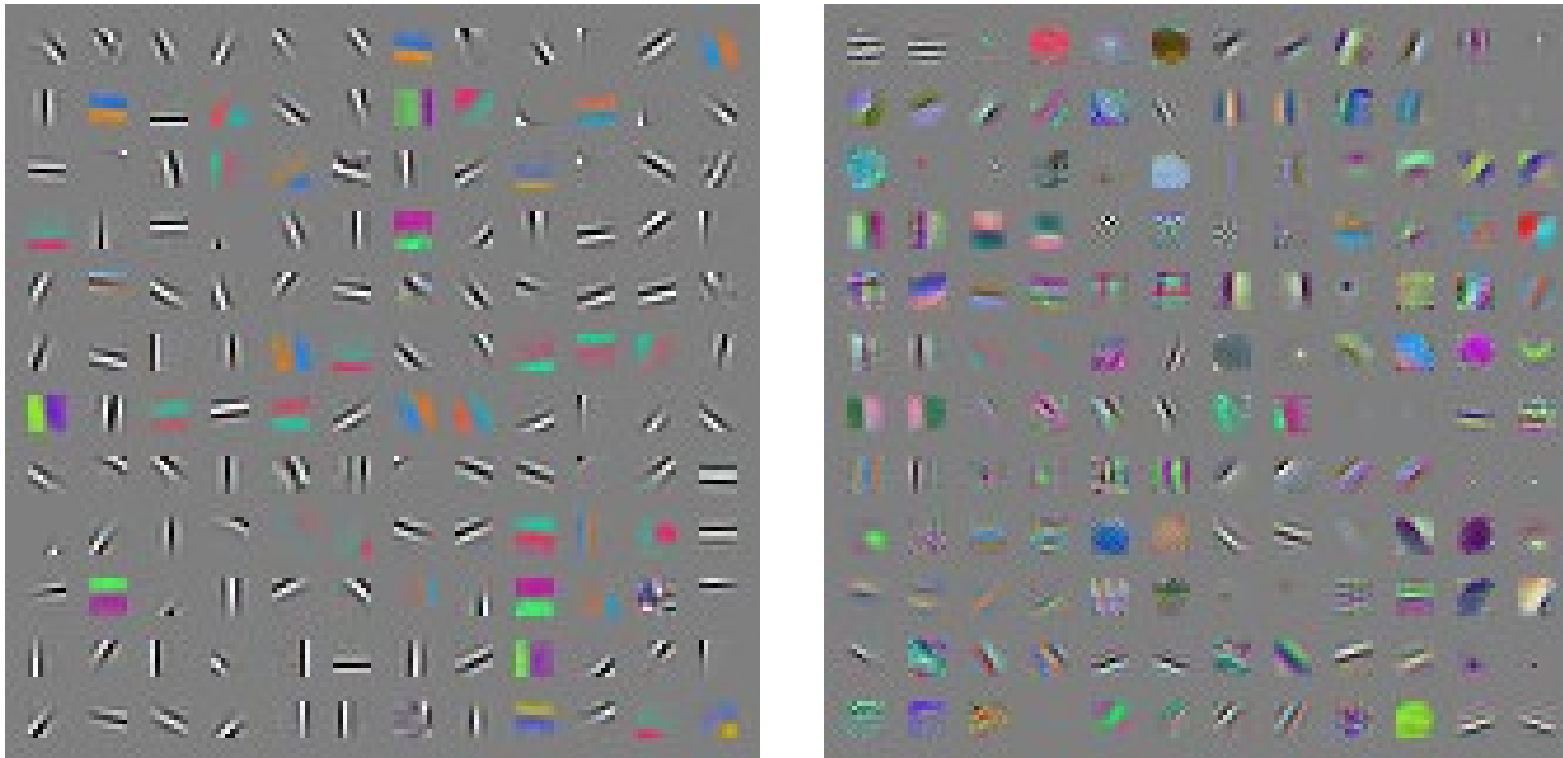Directional second derivative          Second derivative (curvature)          Local average, first derivative
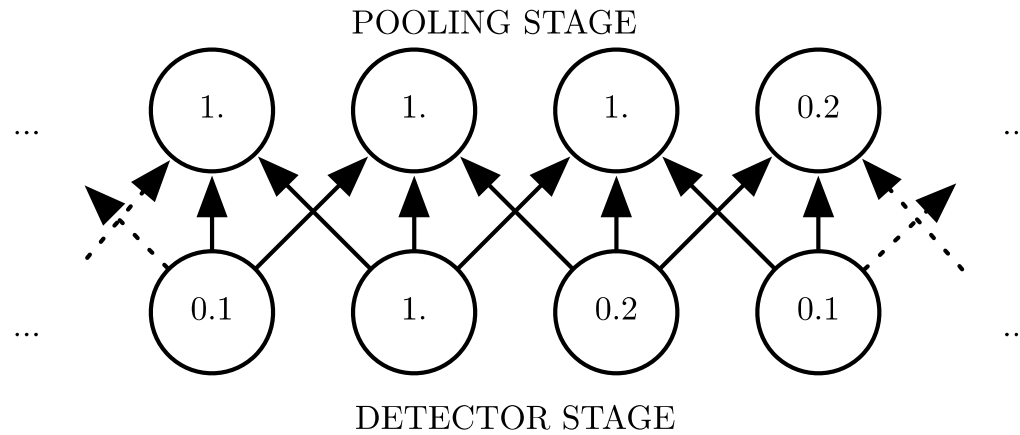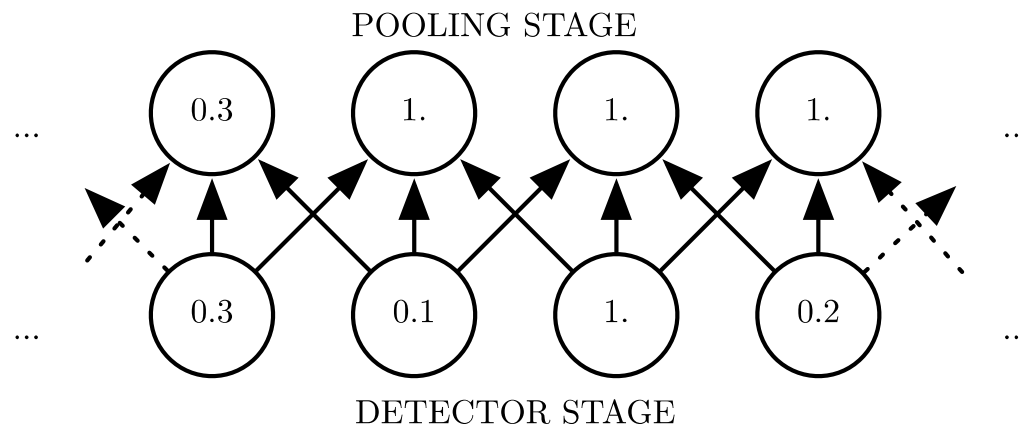
# Gabor-like learned kernels



- Features extractors provided by pretrained networks

[Goodfellow, Bengio, Courville 2016]

# Max pooling translation invariance

POOLING STAGE

... 1. 1. 1. 0.2 ...

... 0.1 1. 0.2 0.1 ...

DETECTOR STAGE

POOLING STAGE

... 0.3 1. 1. 1. ...

... 0.3 0.1 1. 0.2 ...

DETECTOR STAGE

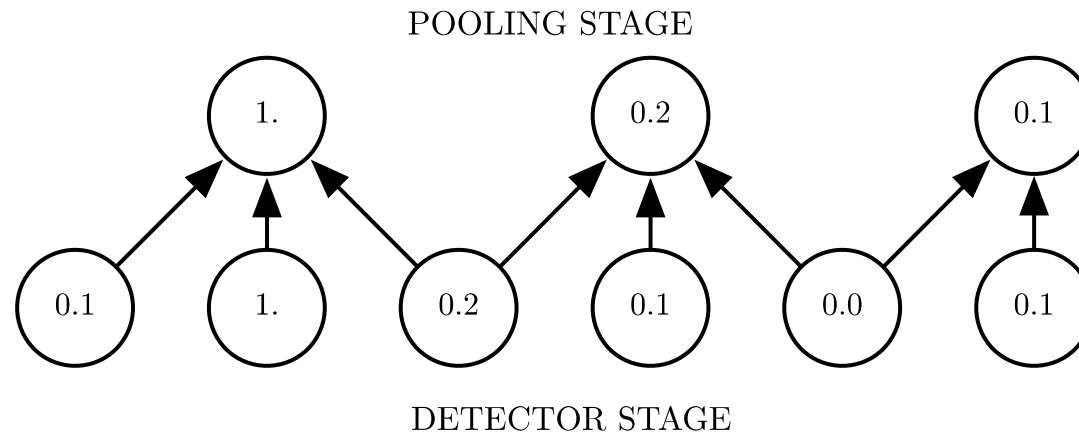- Take max of certain neighbourhood

[Goodfellow, Bengio, Courville 2016]
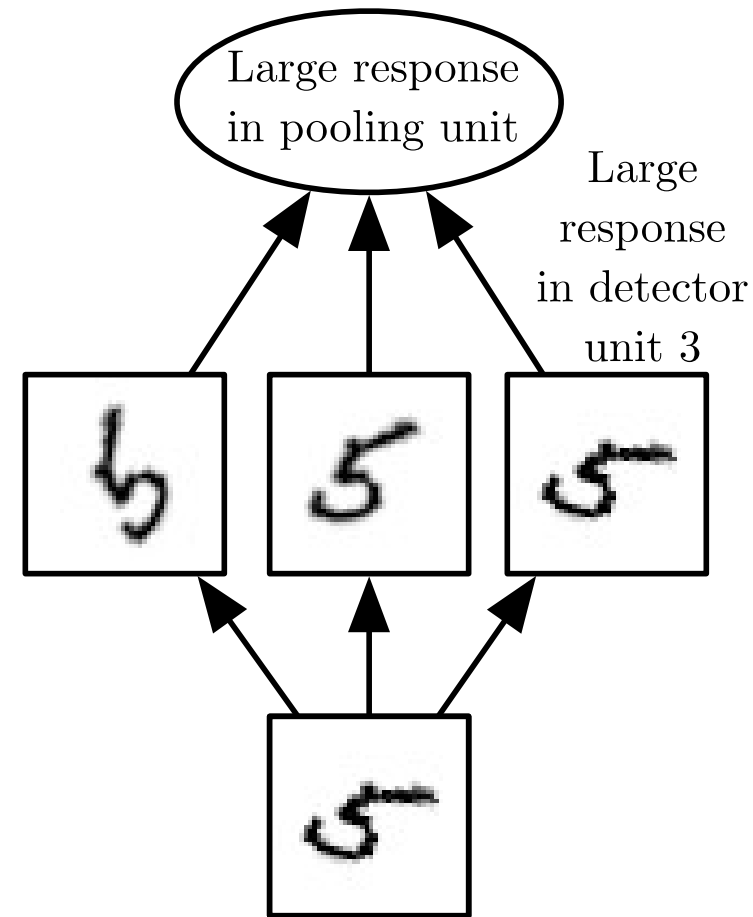
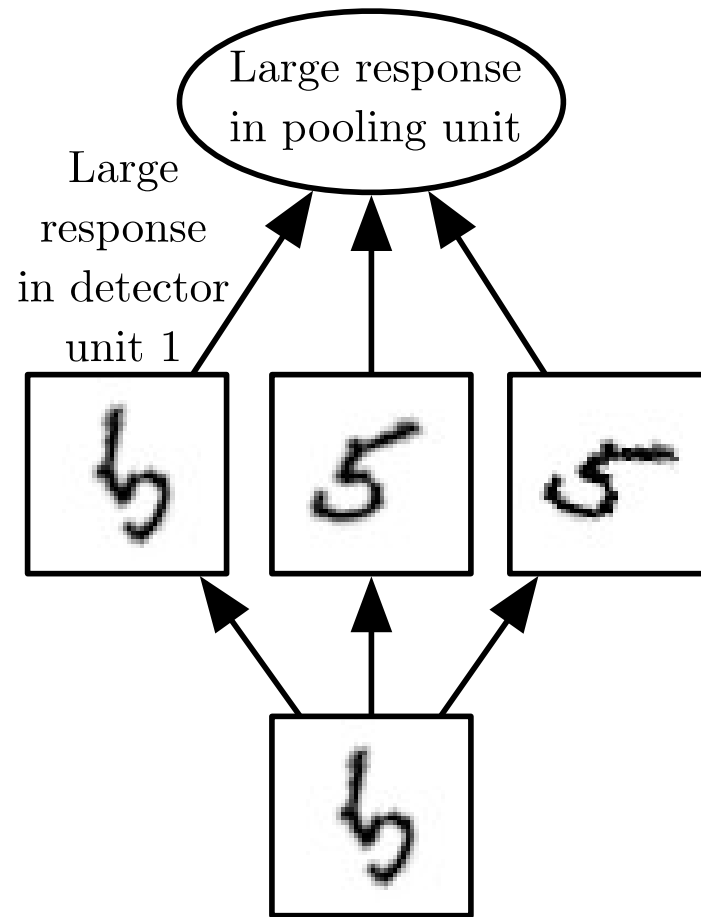# Max pooling translation invariance



- Take max of certain neighbourhood
- Often combined followed by downsampling

[Goodfellow, Bengio, Courville 2016]

# Max pooling transform invariance



[Goodfellow, Bengio, Courville 2016]

42

# Types of connectivity

$s_1$  $s_2$  $s_3$  $s_4$  $s_5$

a  b  c  d  e  f  g  h  i

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

Local connection:
like convolution,
but no sharing

[Goodfellow, Bengio, Courville 2016]

# Types of connectivity



Local connection:
 like convolution,
  but no sharing

Convolution

[Goodfellow, Bengio, Courville 2016]

# Types of connectivity



Local connection:
like convolution,
but no sharing

Convolution

Fully connected

[Goodfellow, Bengio, Courville 2016]

# Choosing architecture family

# Choosing architecture family

- No structure $\rightarrow$ fully connected

# Choosing architecture family

- No structure $\rightarrow$ fully connected

- Spatial structure $\rightarrow$ convolutional

# Choosing architecture family

- No structure $\rightarrow$ fully connected

- Spatial structure $\rightarrow$ convolutional

- Sequential structure $\rightarrow$ recurrent

# Optimization Algorithm

- Lots of variants address choice of learning rate

- See Visualization of Algorithms

- AdaDelta and RMSprop often work well

# Software for Deep Learning

# Current Frameworks

- Tensorflow / Keras

- Pytorch

- DL4J

- Caffe

- And many more

- Most have CPU-only mode but much faster on NVIDIA GPU

# Development strategy

- Identify needs: High accuracy or low accuracy?

- Choose metric
  - Accuracy (% of examples correct), Coverage (% examples processed)
  - Precision TP/(TP+FP), Recall TP/(TP+FN)
  - Amount of error in case of regression

- Build end-to-end system
  - Start from baseline, e.g. initialize with pre-trained network

- Refine driven by data

# Sources

- I. Goodfellow, Y. Bengio, A. Courville "Deep Learning" MIT Press 2016 [link]