# 《多模态机器学习》

## 第二章 基础概念

黄文炳

中国人民大学高瓴人工智能学院

hwenbing@126.com

2024年秋季

# 内容提纲

① 单模态基础表示

② 经典机器学习方法

③ 神经网络

④ 优化方法简介

# 内容提纲

① 单模态基础表示

② 经典机器学习方法
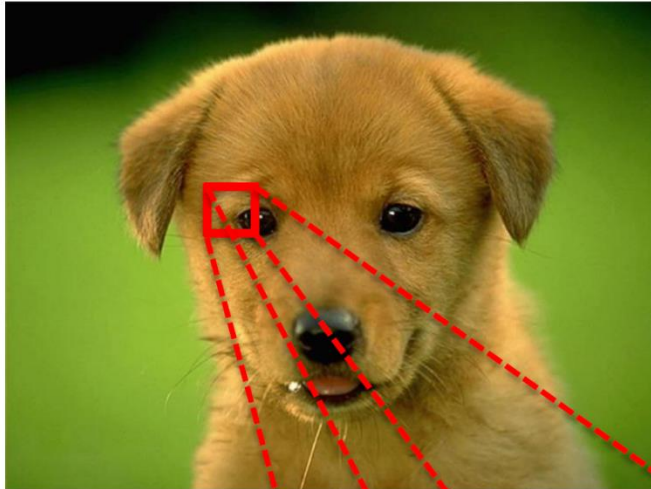
③ 神经网络

④ 优化方法简介

## Color image

**Binary classification problem**

Input observation $x_i$

Each pixel is represented in $\mathbb{R}^d$, $d$ is the number of colors ($d$=3 for RGB)

| 88 | 82 | 84 | 88 | 85 | 83 | 80 | 93 | 102 |
|----|----|----|----|----|----|----|----|-----|
| 88 | 80 | 78 | 80 | 80 | 78 | 73 | 94 | 100 |
| 85 | 79 | 80 | 78 | 77 | 74 | 65 | 91 | 99 |
| 38 | 35 | 40 | 35 | 39 | 74 | 77 | 70 | 65 |
| 20 | 25 | 23 | 28 | 37 | 69 | 64 | 60 | 57 |
| 22 | 26 | 22 | 28 | 40 | 65 | 64 | 59 | 34 |
| 24 | 28 | 24 | 30 | 37 | 60 | 58 | 56 | 66 |
| 21 | 22 | 23 | 27 | 38 | 60 | 67 | 65 | 67 |
| 23 | 22 | 22 | 25 | 38 | 59 | 64 | 67 | 66 |

88
88
85
38
20
22
24
21
23
82
80
79
35
25
26
28
22
22
84
78
80
⋮

# Dog ?

label $y_i \in \mathcal{Y} = \{0,1\}$

# Unimodal Representation: Visual Modality

## Color image



Each pixel is represented in $\mathbb{R}^d$, $d$ is the number of colors ($d$=3 for RGB)

Input observation $x_i$

| 88 | 82 | 84 | 88 | 85 | 83 | 80 | 93 | 102 |
| 88 | 80 | 78 | 80 | 80 | 78 | 73 | 94 | 100 |
| 85 | 79 | 80 | 78 | 77 | 74 | 65 | 91 | 99 |
| 38 | 35 | 40 | 35 | 39 | 74 | 77 | 70 | 65 |
| 20 | 25 | 23 | 28 | 37 | 69 | 64 | 60 | 57 |
| 22 | 26 | 22 | 28 | 40 | 65 | 64 | 59 | 34 |
| 24 | 28 | 24 | 30 | 37 | 60 | 58 | 56 | 66 |
| 21 | 22 | 23 | 27 | 38 | 60 | 67 | 65 | 67 |
| 23 | 22 | 22 | 25 | 38 | 59 | 64 | 67 | 66 |

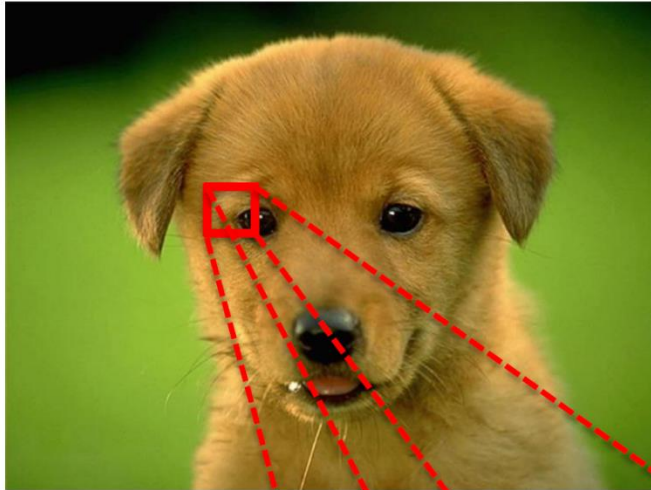## Multi-class classification problem

Duck
-or-
Cat
-or-
Dog
-or-
Pig
-or-
Bird ?

label $y_i \in \mathcal{Y} = \{0,1,2,3, \dots\}$

# Unimodal Representation: Visual Modality

## Color image

Each pixel is represented in $\mathbb{R}^d$, $d$ is the number of colors ($d$=3 for RGB)

Input observation $x_i$

## Multi-label (or multi-task) classification problem

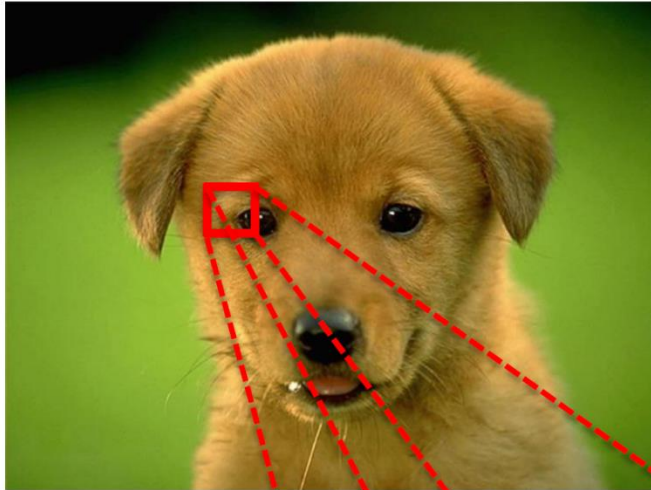Duck?

Cat ?

Dog ?

Pig ?

Bird ?

Puppy ?

label vector $y_i \in \mathcal{Y}^m = \{0,1\}^m$

# Unimodal Representation: Visual Modality

## Color image



Each pixel is represented in $\mathbb{R}^d$, $d$ is the number of colors ($d$=3 for RGB)

Input observation $x_i$

## Multi-label (or multi-task) regression problem

Age ?

Height ?

Weight ?

Distance ?

Happy ?

label vector $y_i \in \mathcal{Y}^m = \mathbb{R}^m$

# Unimodal Representation: Language Modality

**Written language**

★★★★★ Masterful!

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful
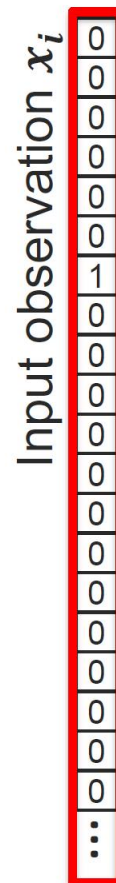
**Spoken language**

MARTHA(CON'T)
Look around you. Look at all the great things you've done and the people you've helped.

CLARK
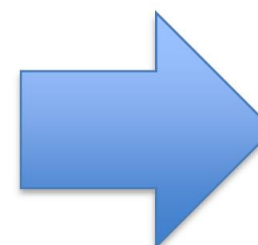But you've only put up the good things they say about me.

MARTHA
Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation $x_i$

| 0 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ⋮ |

**"one-hot" vector**

$|x_i|$ = number of words in dictionary

**Word-level classification**

Part-of-speech ?
(noun, verb,…)

Sentiment ?
(positive or negative)

Named entity ?
(names of person,…)

# Unimodal Representation: Language Modality

**Written language**

⭐⭐⭐⭐⭐ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.
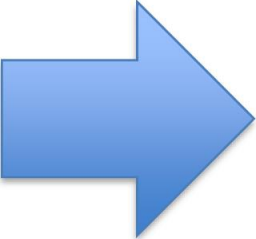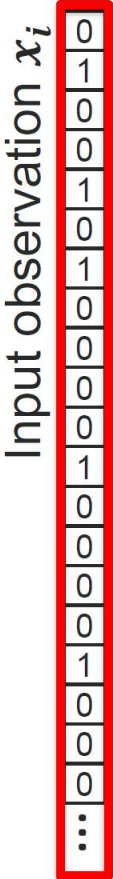
0 of 4 people found this review helpful

**Spoken language**

MARTHA(CON'T)
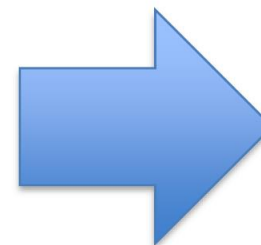Look around you. Look at all the great things you've done and the people you've helped.

CLARK
But you've only put up the good things they say about me.

MARTHA
Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation $x_i$

0
1
0
0
0
1
0
1
0
0
0
0
1
0
0
0
0
1
0
0
0
0
⋮

**"bag-of-word" vector**

$|x_i|$ = number of words in dictionary

**Document-level classification**

Sentiment ?
(positive or negative)

# Unimodal Representation: Language Modality

**Written language**

★★★★★ **Masterful!**

By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful

**Spoken language**

MARTHA(CON'T)
Look around you. Look at all the great things you've done and the people you've helped.

CLARK
But you've only put up the good things they say about me.

MARTHA
Clark, honey. If I were to use the bad things they say I could cover the barn, the house and the outhouse.

Input observation $x_i$

```
0
1
0
0
1
0
1
0
0
0
0
1
0
0
0
0
1
0
0
0
⋮
```

**"bag-of-word" vector**

$|x_i|$ = number of words in dictionary

**Utterance-level classification**

Sentiment ?
(positive or negative)

**Digitalized acoustic signal**

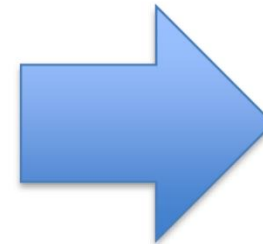- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
  - Offset: 10ms

Input observation $x_i$

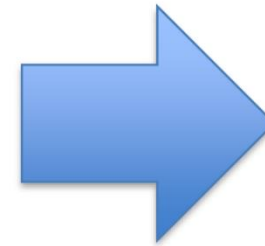| 0.21 |
| 0.14 |
| 0.56 |
| 0.45 |
| 0.9 |
| 0.98 |
| 0.75 |
| 0.34 |
| 0.24 |
| 0.11 |
| 0.02 |

Spoken word ?

**Spectogram**

## Digitalized acoustic signal

- Sampling rates: 8~96kHz
- Bit depth: 8, 16 or 24 bits
- Time window size: 20ms
  - Offset: 10ms

**Spectogram**

Input observation $x_i$

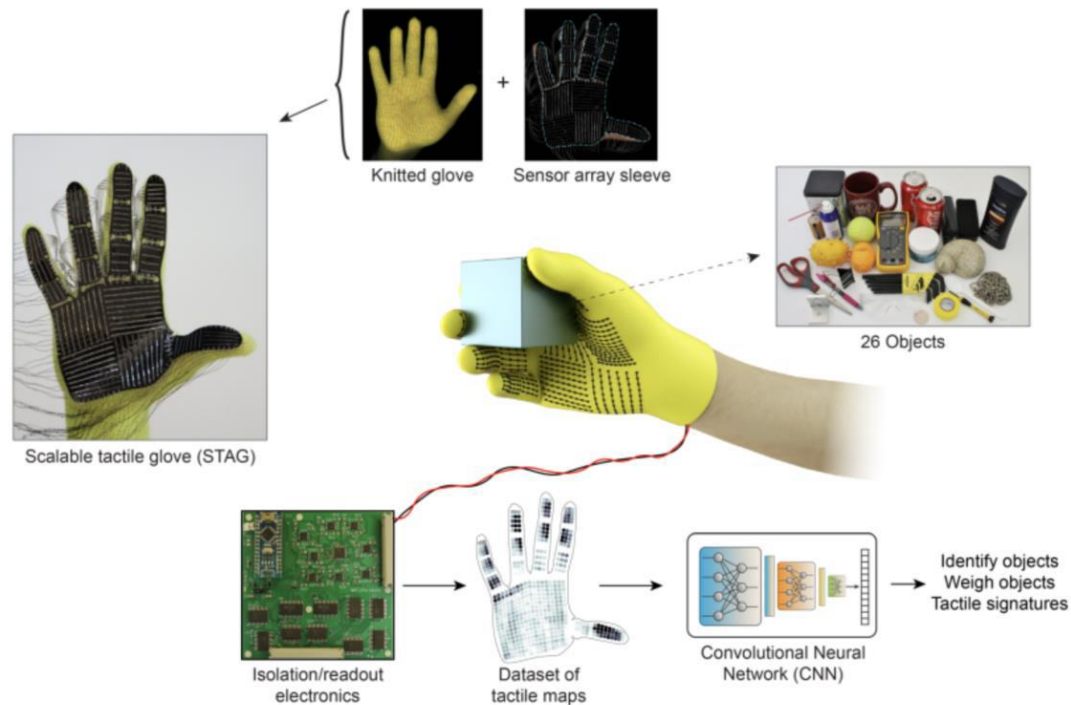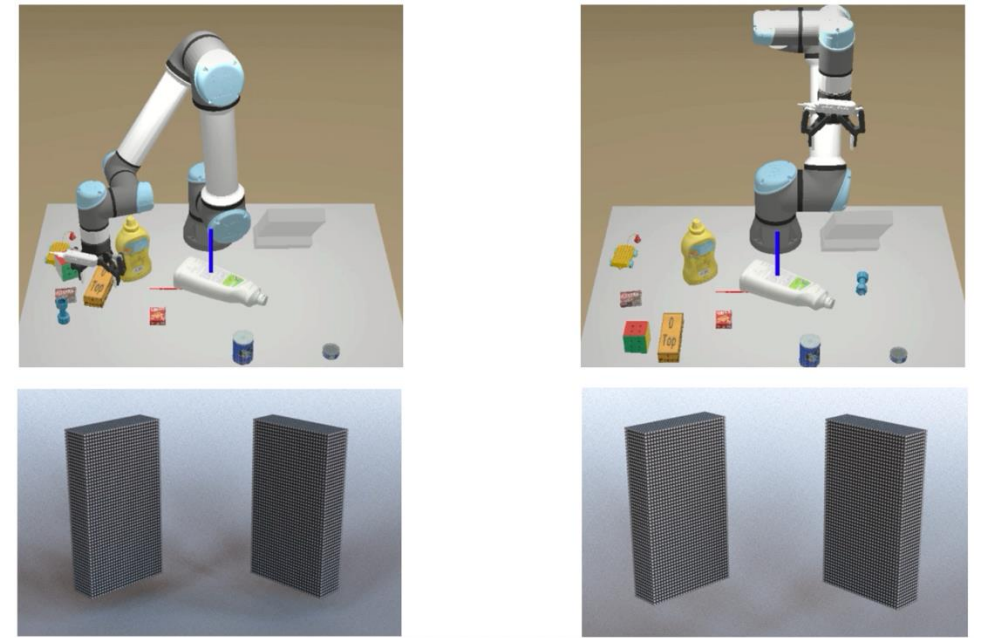| 0.21 |
|------|
| 0.14 |
| 0.56 |
| 0.45 |
| 0.9 |
| 0.98 |
| 0.75 |
| 0.34 |
| 0.24 |
| 0.11 |
| 0.02 |
| 0.24 |
| 0.26 |
| 0.58 |
| 0.9 |
| 0.99 |
| 0.79 |
| 0.45 |
| 0.34 |
| 0.24 |

Emotion ?

Spoken word ?

Voice quality ?

# Unimodal Representation: Tactile Modality

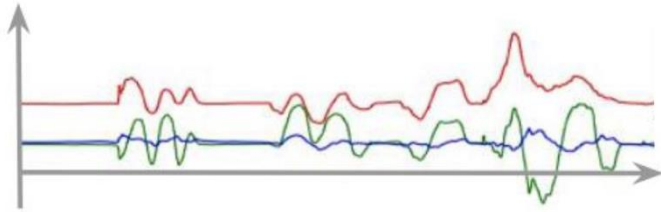The tactile sensor array (548 sensors) is assembled on a knitted glove uniformly distributed over the hand.

Model the elastic property of the tactile sensor



Sundaram et al., Learning the signatures of the human grasp using a scalable tactile glove. Nature 2019

Elastic Interaction of Particles for Robotic Tactile Simulation, ACMMM 2021

# Unimodal Representation: Tactile Modality

Time series data across six-axis Force-Torque sensor: **T × 6 signal.**

Force-Torque Sensor

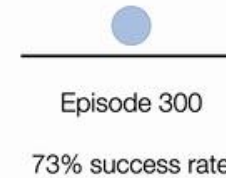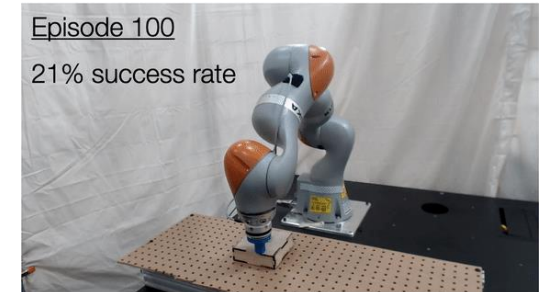## Next action

Proprioception

Time series data across current position and velocity of the end-effector: **T × 2d signal.**

Measure values internal to the system (robot); e.g. motor speed, wheel load, **robot arm joint angles**, battery voltage.

Episode 100

21% success rate

Episode 300

73% success rate

Episode 300

71% success rate

Episode 300

92% success rate

x1

Lee et al., *Making Sense of Vision and Touch: Self-Supervised Learning of Multimodal Representations for Contact Rich Tasks.* ICR A 2019

# Unimodal Representation: Tables



| Caption | Singapore Cup | | |
|---------|------|-----------|-----------|
| **Attribute** | **Year** | **Champions** | **Runners-up** |
| **Cell** | 1996 | Geylang United | Singapore Armed Forces |
| | 1997 | Singapore Armed Forces | Woodlands Wellington |
| | 1998 | Tanjong Pagar United | Sembawang Rangers |

Text — *Singapore Armed forces was the champion of Singapore Cup in 1997.*

Table-to-text generation

Bao et al., Table to Text: Describing Table Region with Natural Language. AAAI 2018
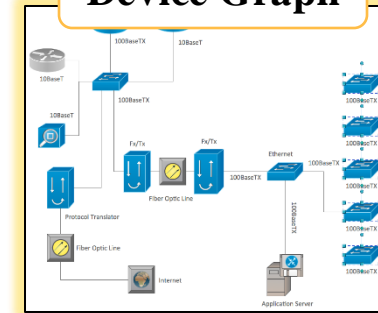
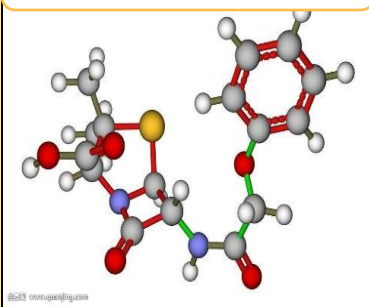# Unimodal Representation: Graphs
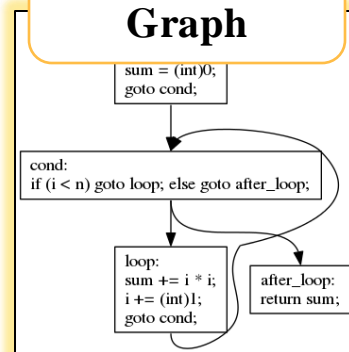
**Web Graph**

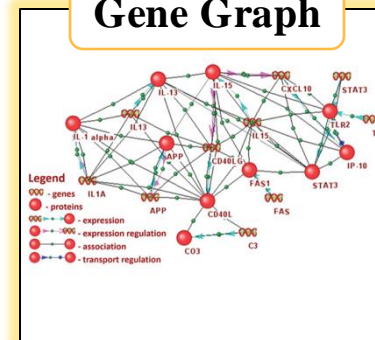**Knowledge Graph**

**Social Graph**

**Device Graph**

## Graphs are Everywhere

**Molecular Graph**

**Control Flow Graph**

**Gene Graph**

**Brain Graph**

# Unimodal Representation: Sets



Sets

black hair & rosey cheeks

attractive & heavy makeup

double-chin & wavy hair

black hair & brown hair

attractive & mouth slightly open
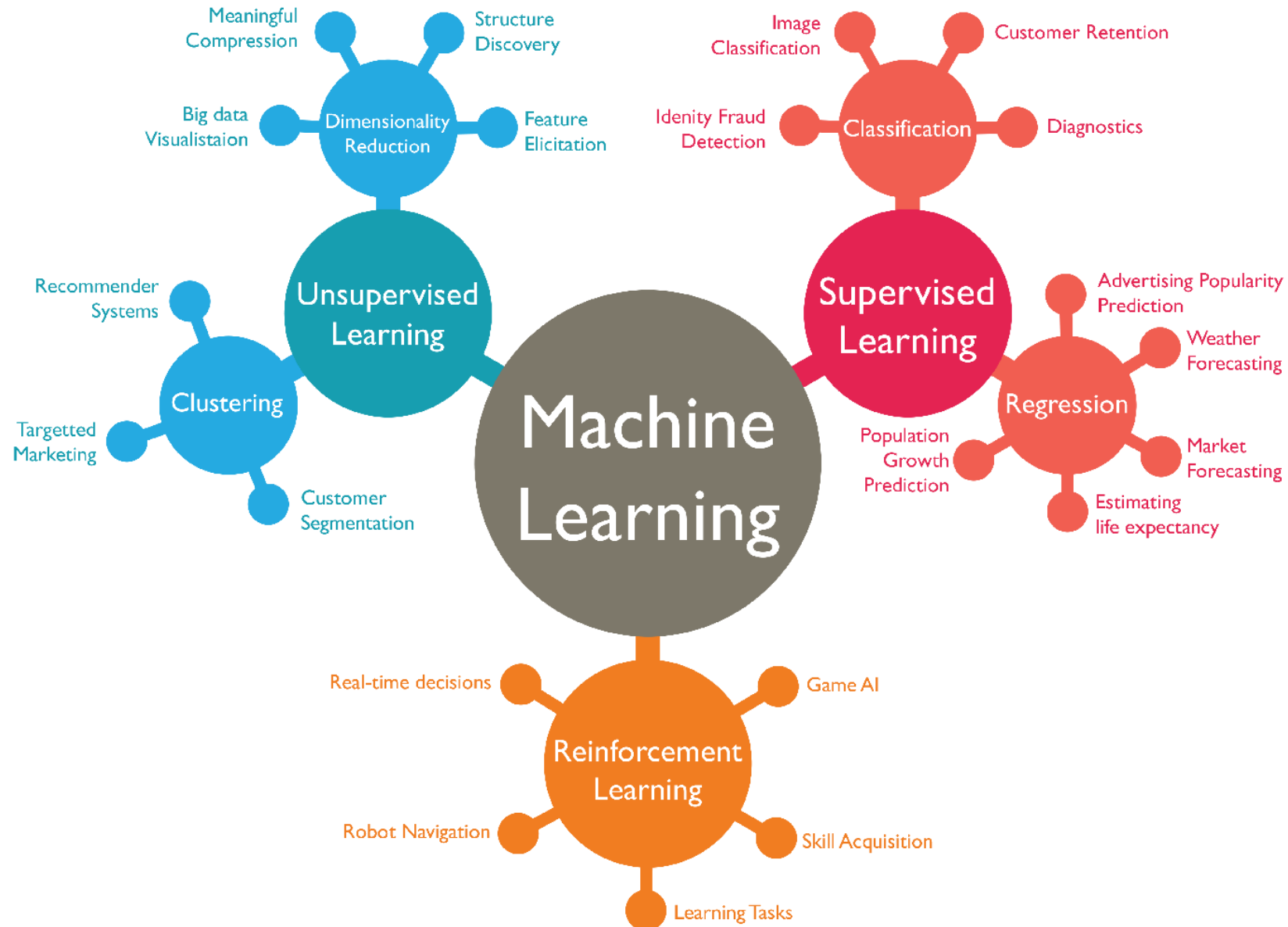
Point clouds

Set anomaly detection
Set expansion
Set completion
Point cloud classification
Point cloud generation

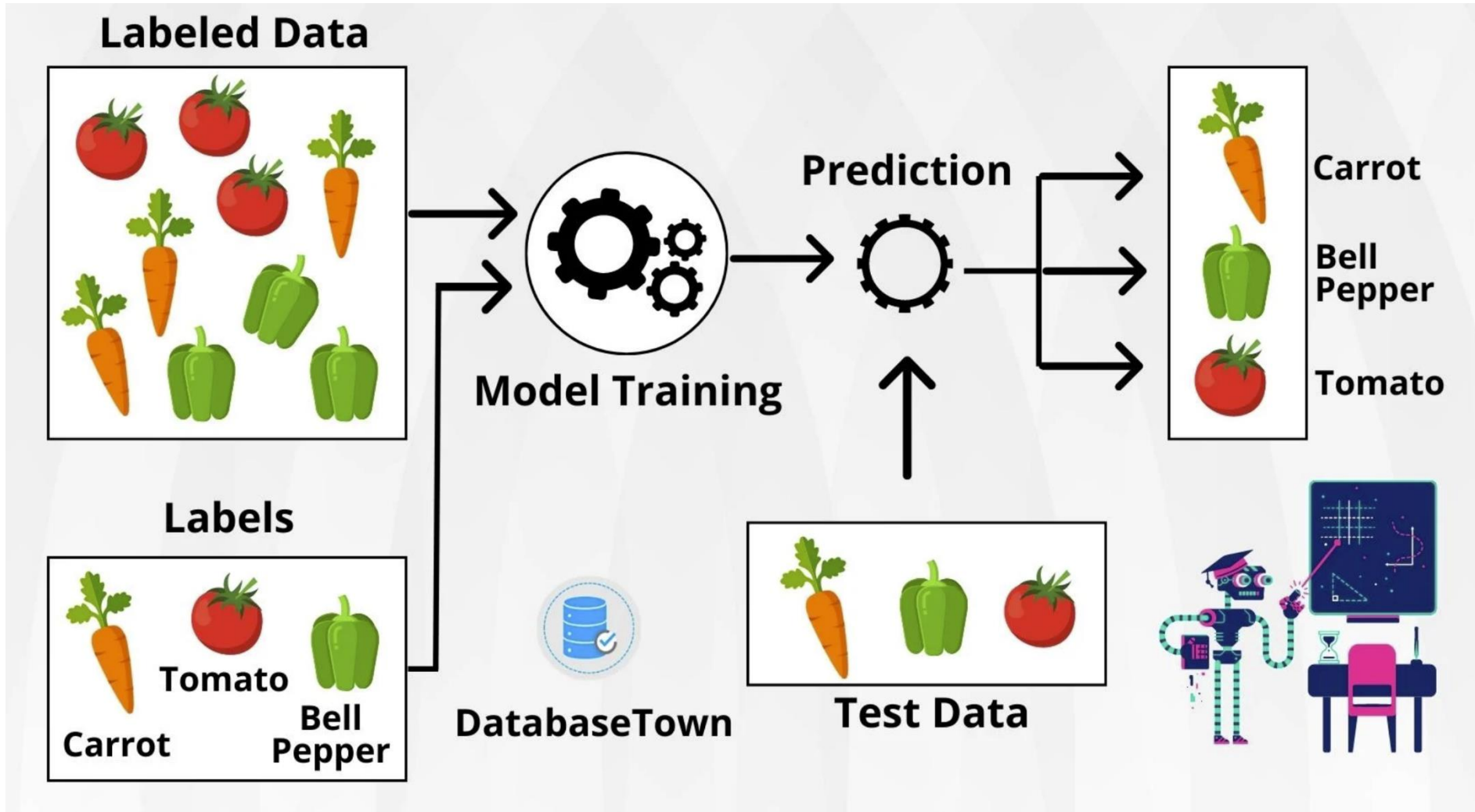Zaheer et al., DeepSets . NeurIPS 2017, Li et al., Point Cloud GAN. arxiv 2018

# 内容提纲

# Supervised Learning

# Simplest Classifier ?



Traffic light
-or-
Dog
-or-
Basket
-or-
Kayak ?

Dataset

# Simple Classifier: Nearest Neighbor

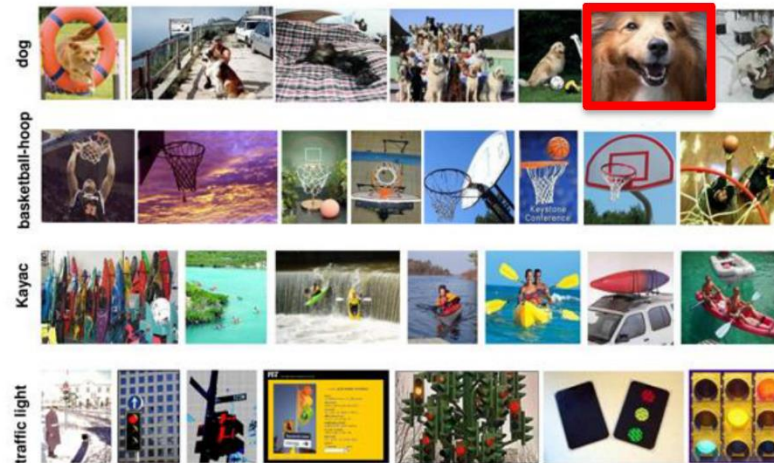- Non-parametric approaches—key ideas:
  - *"Let the data speak for themselves"*
  - *"Predict new cases based on similar cases"*
  - *"Use multiple local models instead of a single global model"*

## Distance metrics

L1 (Manhattan) distance:

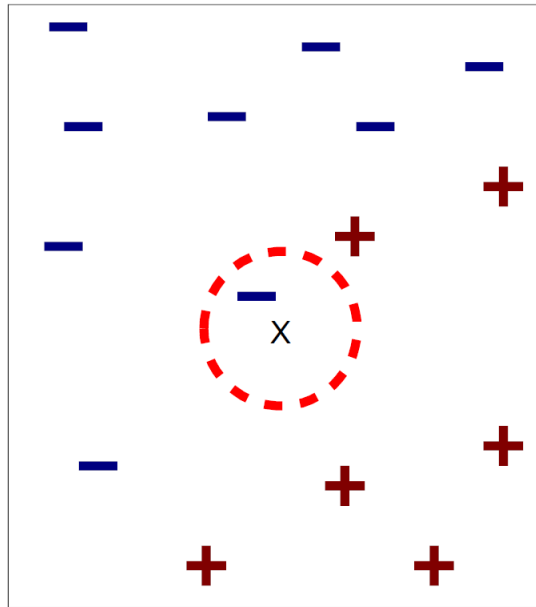$$d_1(x_1, x_2) = \sum_j \left| x_1^j - x_2^j \right|$$

L2 (Eucledian) distance:

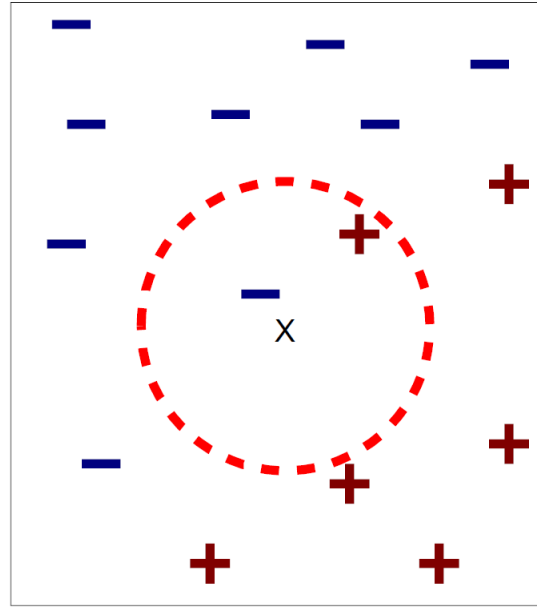$$d_2(x_1, x_2) = \sqrt{\sum_j \left( x_1^j - x_2^j \right)^2}$$

**Which distance metric to use?**

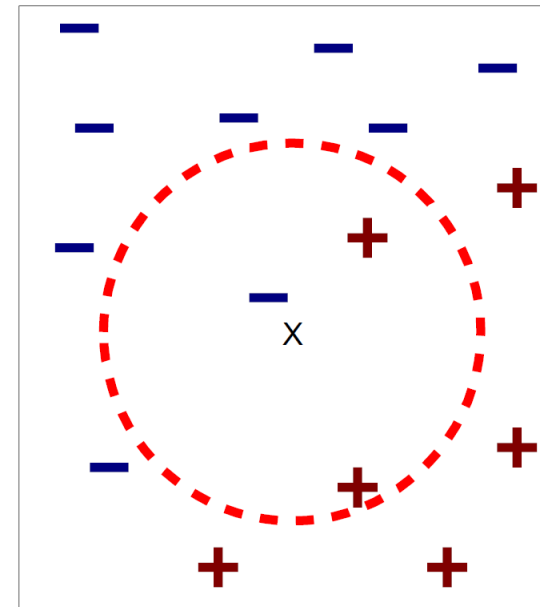First hyper-parameter!

(a) 1-nearest neighbor     (b) 2-nearest neighbor     (c) 3-nearest neighbor

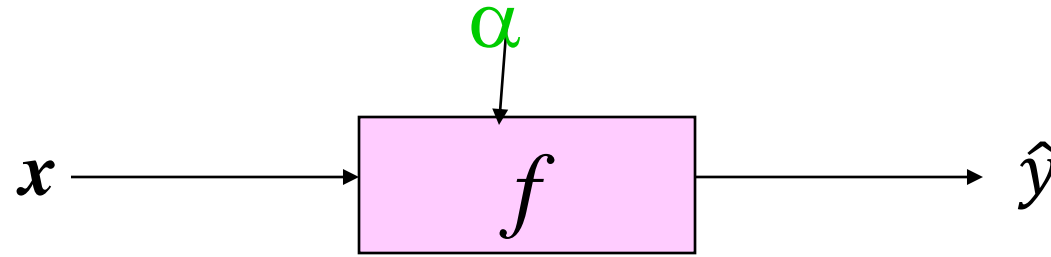**What value should we set K?**

Second hyper-parameter!

# Linear Classifier

$$\alpha$$

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$f(x,w,b) = \text{sign}(wx - b)$$

● denotes +1
○ denotes -1

How would you classify this data?

# Linear Classifier

$\alpha$

$x \longrightarrow f \longrightarrow \hat{y}$
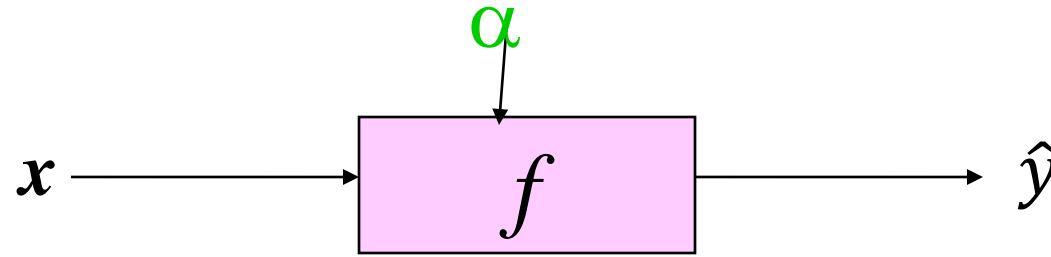
$f(x,w,b) = \text{sign}(wx - b)$

●   denotes +1

  denotes -1

○

How would you
classify this data?

# Linear Classifier

α

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$f(x,w,b) = \text{sign}(wx - b)$$

- • denotes +1
- ○ denotes -1

Any of these would be fine..

..but which is best?

# Classifier Margin



$$\alpha$$

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$f(x,w,b) = \text{sign}(wx - b)$$

• denotes +1
○ denotes -1

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

$$\alpha$$

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$f(\boldsymbol{x},\boldsymbol{w},b) = \text{sign}(\boldsymbol{w}\boldsymbol{x} - b)$$

• denotes +1
○ denotes -1

The maximum margin linear classifier is the linear classifier with the maximum margin.
This is the simplest kind of SVM (Called an LSVM)

Linear SVM

# Support Vector Machine: Maximum Margin

$$\alpha$$

$$x \longrightarrow \boxed{f} \longrightarrow \hat{y}$$

$$f(x,w,b) = \text{sign}(wx - b)$$

• denotes +1
○ denotes -1

Support Vectors are those datapoints that the margin pushes up against

The maximum margin linear classifier is the linear classifier with the maximum margin.
This is the simplest kind of SVM (Called an LSVM)

Linear SVM

# Support Vector Machine: Maximum Margin



$M = $ Margin Width $= \dfrac{2}{\sqrt{\mathbf{w.w}}}$

"Predict Class = +1" zone
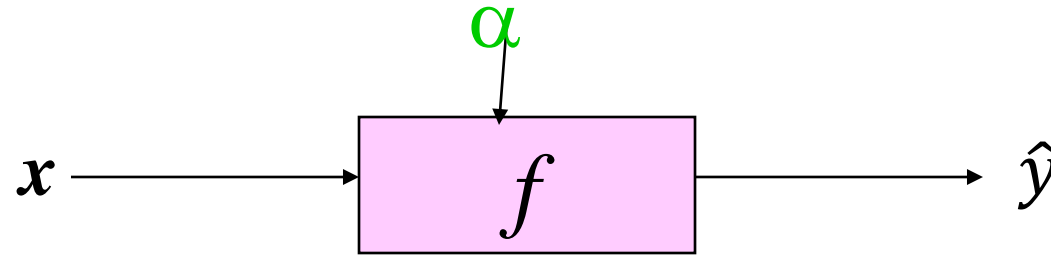
"Predict Class = -1" zone

$wx+b=1$

$wx+b=0$

$wx+b=-1$

$x^+$

$x^-$

- Plus-plane $= \{\, x : w \cdot x + b = +1 \,\}$
- Minus-plane $= \{\, x : w \cdot x + b = -1 \,\}$

Classify as..   +1      if    $w \cdot x + b >= 1$

-1      if    $w \cdot x + b <= -1$

# Representation Learning: A Review and New Perspectives

Yoshua Bengio[†], Aaron Courville, and Pascal Vincent[†]
Department of computer science and operations research, U. Montreal
† also, Canadian Institute for Advanced Research (CIFAR)

◆

**Data-driven feature representation learning: Deep neural networks**

# 内容提纲

- ## Made up of artificial neurons

# Neural Networks: score function

- Made up of artificial neurons
  - Linear function (dot product) followed by a nonlinear activation function
- Example a Multi Layer Perceptron



input layer

hidden layer 1    hidden layer 2

output layer

# Basic Neural Network Building Block

**Input**

Weighted sum
$Wx + b$

Activation function $f(\ )$

**Output**

Linear classifier

Nonlinear activation function

$$\hat{y} = f(Wx + b)$$

$$f(x_i; W, b) = W x_i + b$$

The linear classifier defines a decision plane:

$$W x_i + b > 0$$

$f(x) > 0$

$f(x) = 0$

$f(x) < 0$

$x_2$

$\mathcal{R}_1$

$\mathcal{R}_2$

$\mathbf{x}$

$\mathbf{w}$

$\dfrac{f(x)}{\|w\|}$

$\mathbf{x}_\perp$

$x_1$

$\dfrac{-b}{\|w\|}$

Output can be seen as distance to decision plane

# Neural Networks: activation function

- $f(x) = \tanh(x)$

- Sigmoid - $f(x) = (1 + e^{-x})^{-1}$

- Linear $- f(x) = ax + b$

- ReLU $\qquad f(x) = \max(0, x) \sim \log(1 + \exp(x))$
  - Rectifier Linear Units
  - Faster training - no gradient vanishing
  - Induces sparsity

# Multi-Layer Perceptron (MLP)

## Activation functions (individual layers)

$$f_{1;W_1}(x) = \sigma(W_1 x + b_1)$$

$$f_{2;W_2}(x) = \sigma(W_2 x + b_2)$$

$$f_{3;W_3}(x) = \sigma(W_3 x + b_3)$$



## Score function

$$y_i = f(x_i) = f_{3;W_3}(f_{2;W_2}(f_{1;W_1}(x_i)))$$

How to integrate all the output scores?

# Neural Network: loss function

(or cost function or objective)

**Scores**  **Label**  ⟶  **Loss**

$f(x_i; W)$   $y_i = 2 \ (dog)$   $L_i = ?$

Image $x_i$



(Size: 32*32*3)

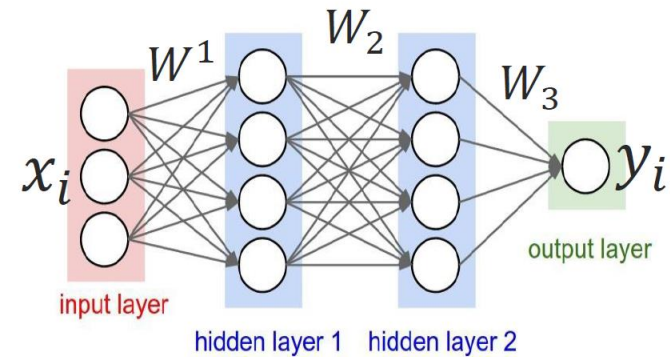| | |
|---|---|
| 0 (duck) ? | -12.3 |
| 1 (cat) ? | 45.6 |
| 2 (dog) ? | 98.7 |
| 3 (pig) ? | 12.2 |
| 4 (bird) ? | -45.3 |

**Multi-class problem**

How to assign only one number representing how "unhappy" we are about these scores?

**The loss function quantifies the amount by which the prediction scores deviate from the actual values.**

# First Loss Function: Cross Entropy Loss

(or logistic loss)

**Logistic function:**

$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

**Logistic regression:**
(two classes)

$$p(y_i = \text{"dog"} | x_i; w) = \sigma(w^T x_i)$$

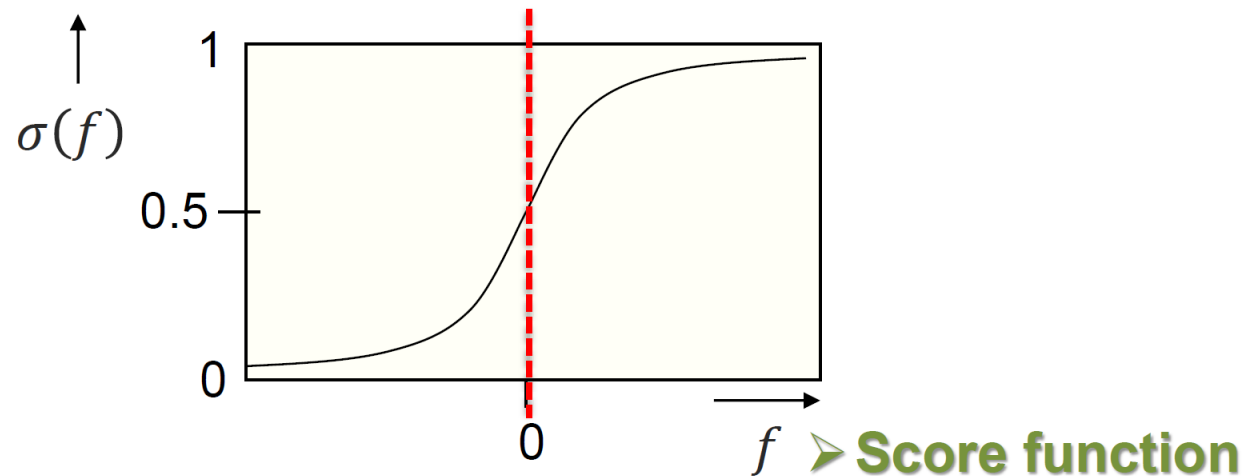**= true**
for two-class problem



➢ **Score function**

# First Loss Function: Cross Entropy Loss

(or logistic loss)

**Logistic function:**
$$\sigma(f) = \frac{1}{1 + e^{-f}}$$

**Logistic regression:**
(two classes)
$$p(y_i = "dog"|x_i; w) = \sigma(w^T x_i)$$

**= true**
for two-class problem

**Softmax function:**
(multiple classes)
$$p(y_i|x_i; W) = \frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$

**Cross-entropy loss:**
$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right)$$

Softmax function

Minimizing the negative log likelihood.

# Second Loss Function: Hinge Loss

(or max-margin loss or Multi-class SVM loss)

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

loss due to example i

sum over all incorrect labels

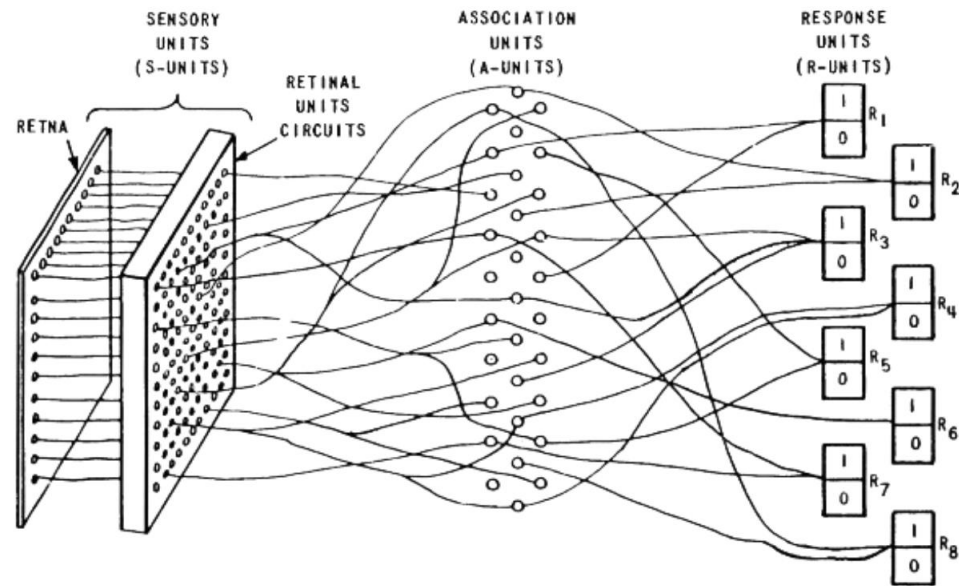difference between the correct class score and incorrect class score

delta

scores for other classes

score for correct class

score

# Different concepts

| Sigmoid | Relu | Logistic Function | Softmax Function |
|---|---|---|---|
| $\sigma(x) := \dfrac{1}{(1 + e^{-x})}$ | $\sigma(x) := \max(0, x)$ | $\sigma(\boldsymbol{x}) := \dfrac{1}{\left(1 + e^{-(\boldsymbol{w}^\top \boldsymbol{x} + b)}\right)}$ (Or Linear Function + Sigmoid) | $\sigma(\boldsymbol{x})_y := \dfrac{e^{x_y}}{\sum_j e^{x_j}}$ |

| Cross-Entropy Loss | Hinge Loss |
|---|---|
| $L(\boldsymbol{p}, y) := -\log \hat{p}_y$ | $L(\boldsymbol{x}, y) := \sum_{j \neq y} \max(0, x_j - x_y + \Delta)$ |

| Linear Classifier (Perceptron) | Logistic Regression | MLP |
|---|---|---|
| $\text{sign}(\boldsymbol{w}^\top \boldsymbol{x} + b)$ | Logistic Function + Cross-Entropy Loss | (Linear Function + Activation) $\times L$ + Softmax + Cross-Entropy Loss |

*Early neural networks*



Perceptron, one of the first neural network architectures

F. Rosenblatt

1957

Rosenblatt 1957
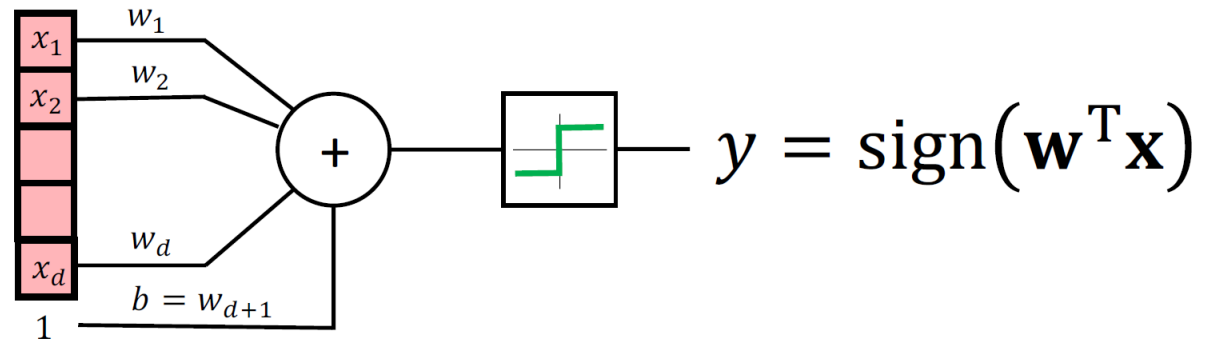
*"Simple perceptron"*


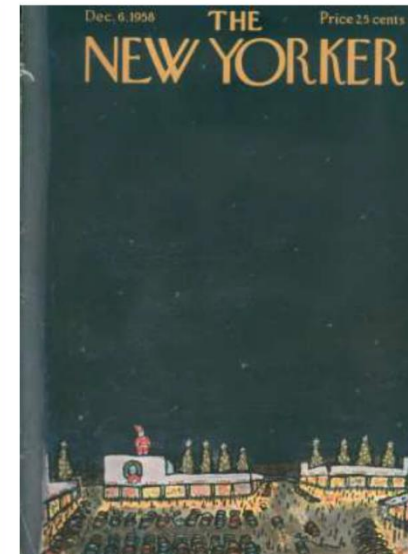
$$y = \text{sign}(\mathbf{w}^{\mathrm{T}}\mathbf{x})$$

# Brief History of Deep Learning

*Early hype*

"First serious rival to the human brain even devised."

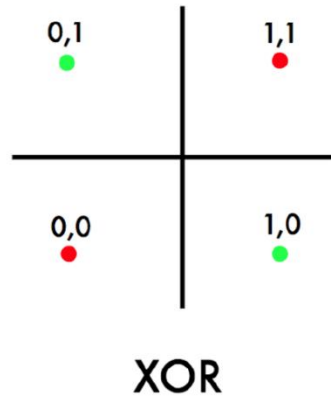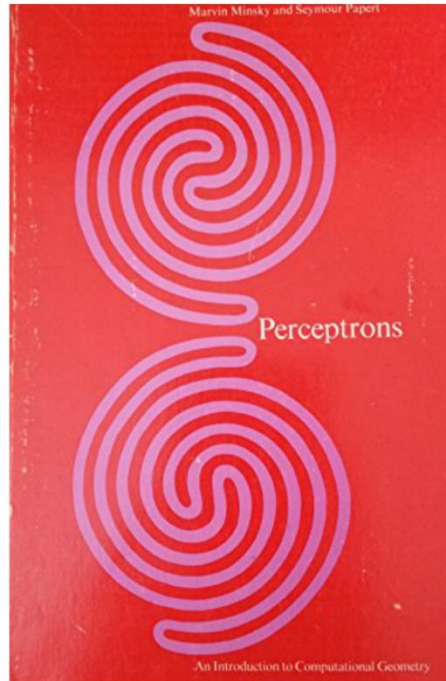"Remarkable machine capable of what amounts to thought"

— The New Yorker



1958

Manson, Stewart, Gill 1958

*The "XOR Affair"*



0,1    1,1

0,0    1,0

**XOR**

"[simple] perceptron cannot represent even the XOR function"

*Perceptrons*

Marvin Minsky and Seymour Papert

An Introduction to Computational Geometry

M. Minsky    S. Papert

1969

Minsky, Papert 1969

"AI WINTER"

*Universal approximation*



| D. Hilbert | A. Kolmogorov | V. Arnold | G. Cybenko | K. Hornik |

13th Problem

$$f(x_1, \ldots, x_n) = \sum_{q=0}^{2n} \Phi_q \left( \sum_{p=1}^{n} \phi_{q,p}(x_p) \right)$$

Results specific to multilayer
neural networks

Hilbert 1900; Arnold 1956; Kolmogorov 1957; Cybenko 1989; Hornik 1991

# Brief History of Deep Learning

*Universal approximation*



"A 2-layer perceptron can approximate a continuous function to any desired accuracy"

Cybenko 1989; Hornik 1991; Barron 1993; Leshno et al 1993; Maiorov 1999; Pinkus 1999

# 内容提纲

We have our training data
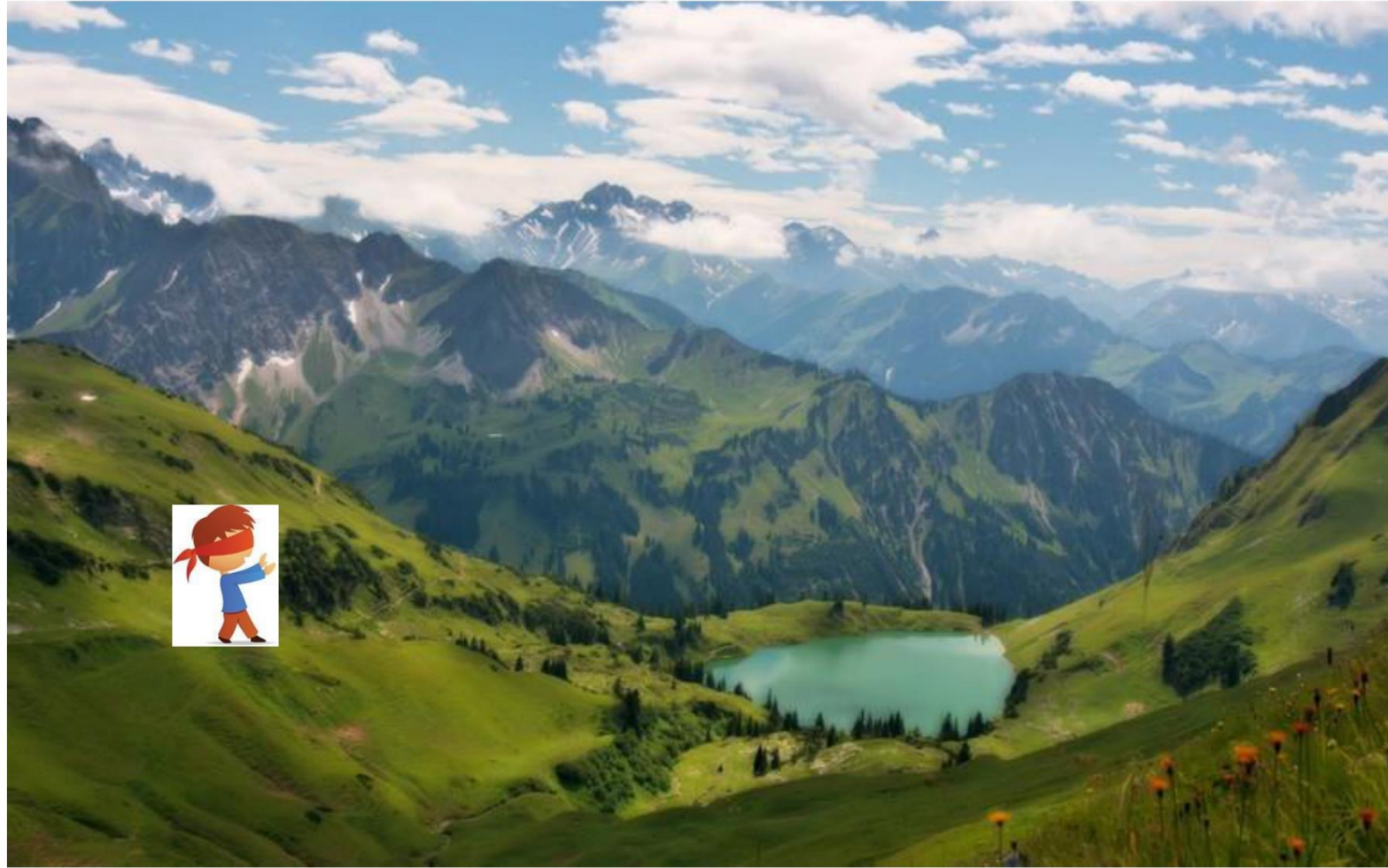
- $\mathbb{X} = \{x_1, x_2, \dots, x_n\}$ (e.g. images, videos, text etc.)
- $\mathbb{Y} = \{y_1, y_2, \dots, y_n\}$ (labels)

We want to learn the W (weights and biases) that leads to best loss

$$\operatorname*{argmin}_{W}[L(\mathbb{X}, \mathbb{Y}, W)]$$

The notation means find $W$ for which $L(\mathbb{X}, \mathbb{Y}, W)$ has the lowest value

# Optimization

If we know the function and it is **differentiable**

- Derivative/gradient is defined at every point in *f*
- Sometimes use differentiable approximations
- Some are locally differentiable

Examples:

$$f(x) = \frac{1}{1 + e^{-x}}; \frac{df}{dx} = (1 - f(x))f(x)$$

$$f(x) = (x - y)^2; \frac{df}{dx} = 2(x - y)$$

Many methods for optimization

- **Gradient Descent (actually the "simplest" one)**
- Newton methods (use Hessian – second derivative)
- Quasi-Newton (use approximate Hessian)
    - BFGS
    - LBFGS
    - Don't require learning rates (fewer hyperparameters)
    - But, do not work with stochastic and batch methods so rarely used to train modern Neural Networks

**All of them look at the gradient**

- Very few non gradient based optimization methods

# Parameter Update Strategies

## Gradient descent:

$$\theta^{(t+1)} = \theta^t - \epsilon_k \boxed{\nabla_\theta L} \rightarrow \text{Gradient of our loss function}$$

New model parameters

Previous parameters

**Learning rate** at iteration $k$

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha \boxed{\epsilon_\tau} \rightarrow \text{Decay learning rate linearly until iteration } \tau$$

**Learning rate** at iteration $k$

**Decay**

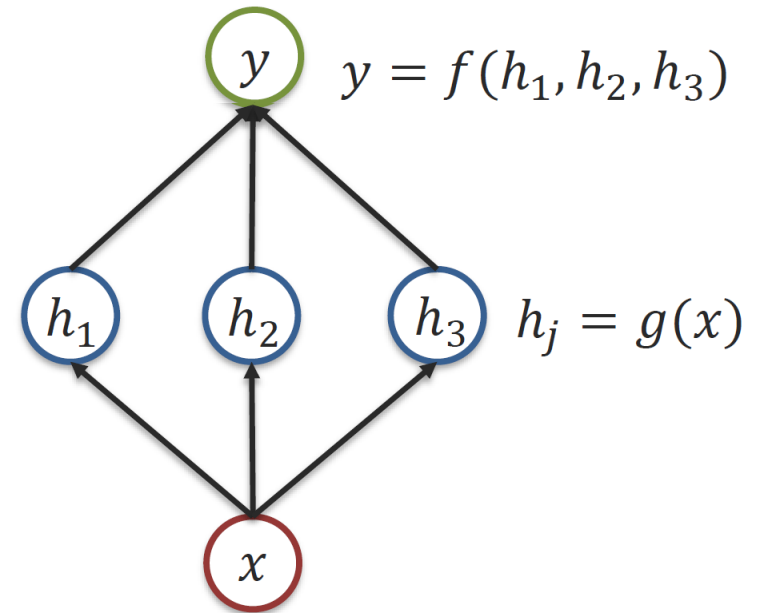Initial learning rate

# Gradient Computation

Chain rule:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial h}\frac{\partial h}{\partial x}$$

$y = f(h)$

$h = g(x)$

Multiple-path chain rule:

$$\frac{\partial y}{\partial x} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x}$$



$y = f(h_1, h_2, h_3)$

$h_j = g(x)$

Multiple-path chain rule:

$$\frac{\partial y}{\partial x_1} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_1}$$

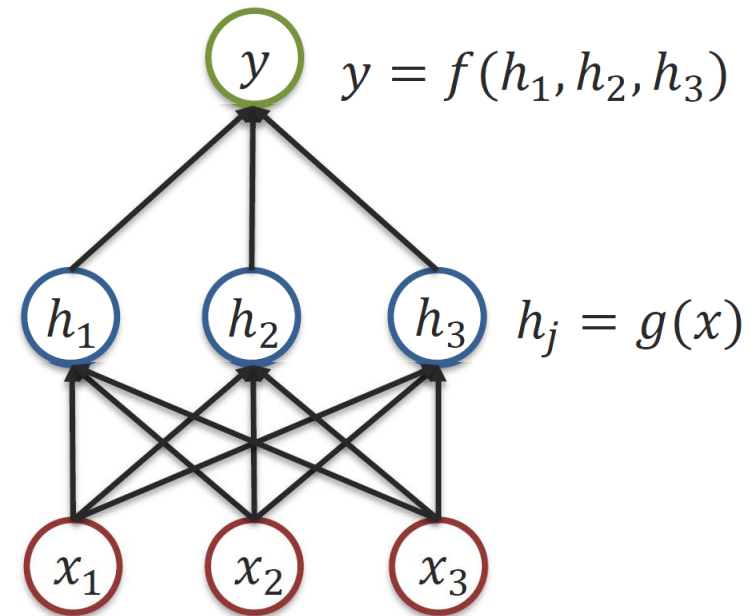$$\frac{\partial y}{\partial x_2} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_2}$$

$$\frac{\partial y}{\partial x_3} = \sum_j \frac{\partial y}{\partial h_j} \frac{\partial h_j}{\partial x_3}$$
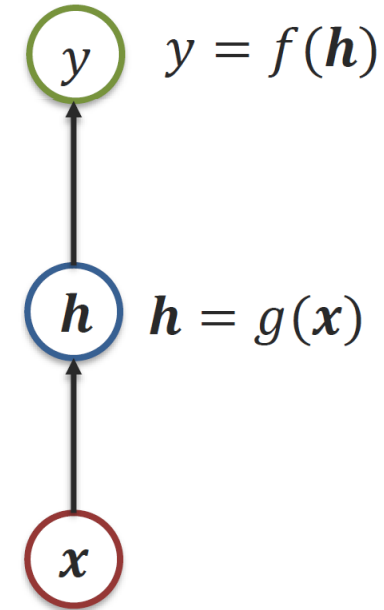


$y = f(h_1, h_2, h_3)$

$h_j = g(x)$

# Gradient Computation

Vector representation:

$$\nabla_x y = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \frac{\partial y}{\partial x_3} \right]$$

Gradient

$$\nabla_x y = \left( \frac{\partial h}{\partial x} \right)^T \nabla_h y$$

"backprop" Gradient

"local" Jacobian
(matrix of size $|h| \times |x|$ computed using partial derivatives)

$y$   $y = f(h)$

$h$   $h = g(x)$

$x$

# Back-Propagation Algorithm (efficient gradient)
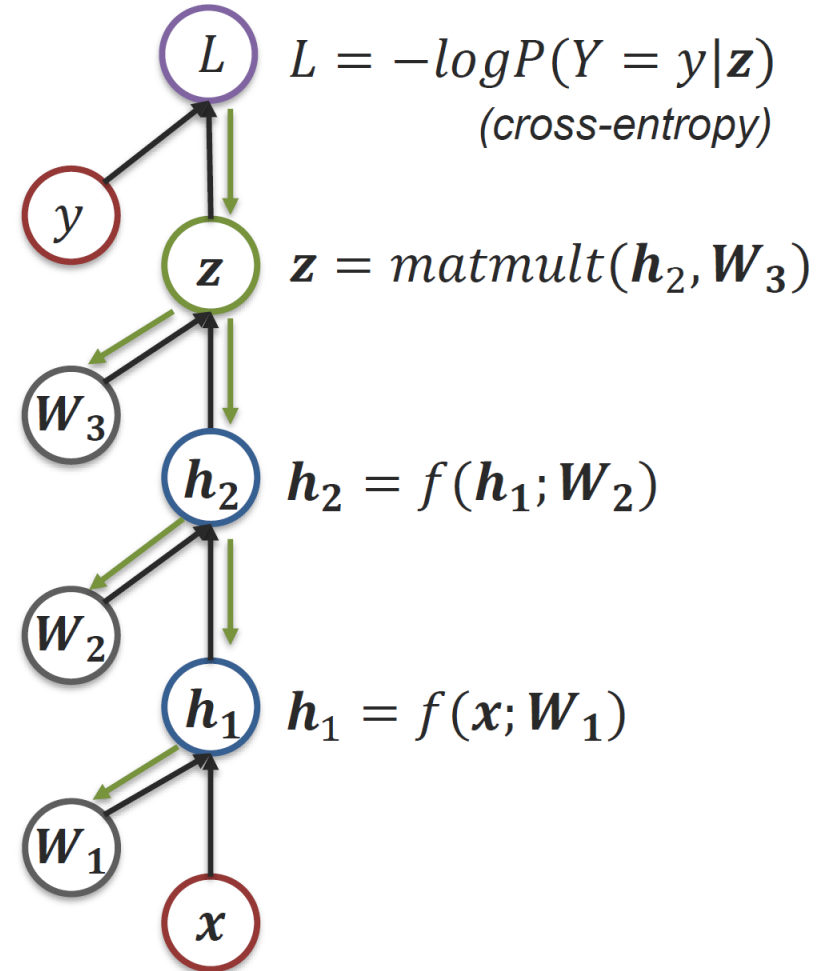
## Forward pass

- Following the graph topology, compute value of each unit

## Backpropagation pass

- Initialize output gradient = 1

- Compute "local" Jacobian matrix using values from forward pass
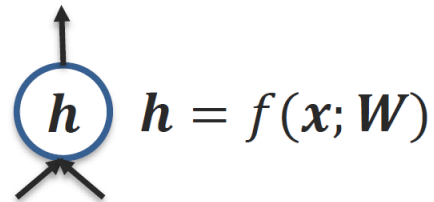
- Use the chain rule:

  Gradient = "local" Jacobian x
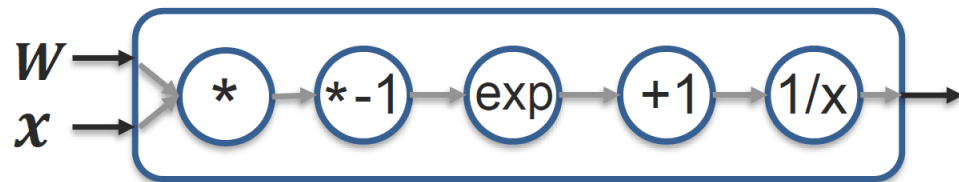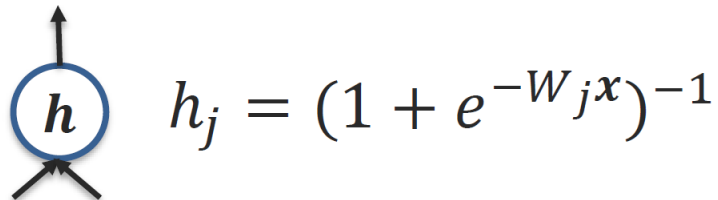
  "backprop" gradient

- Why is this rule important?



$L = -logP(Y = y|z)$

*(cross-entropy)*

$z = matmult(h_2, W_3)$

$h_2 = f(h_1; W_2)$

$h_1 = f(x; W_1)$

# Computational Graph: Multilayer Feedforward Network

Computational unit:

$$h \quad \boldsymbol{h} = f(\boldsymbol{x}; \boldsymbol{W})$$

- Multiple input
- One output
- Vector/tensor

- Sigmoid unit:

$$h \quad h_j = (1 + e^{-W_j \boldsymbol{x}})^{-1}$$

$W \rightarrow$
$x \rightarrow$   $* \rightarrow *{-}1 \rightarrow \exp \rightarrow +1 \rightarrow 1/x \rightarrow$
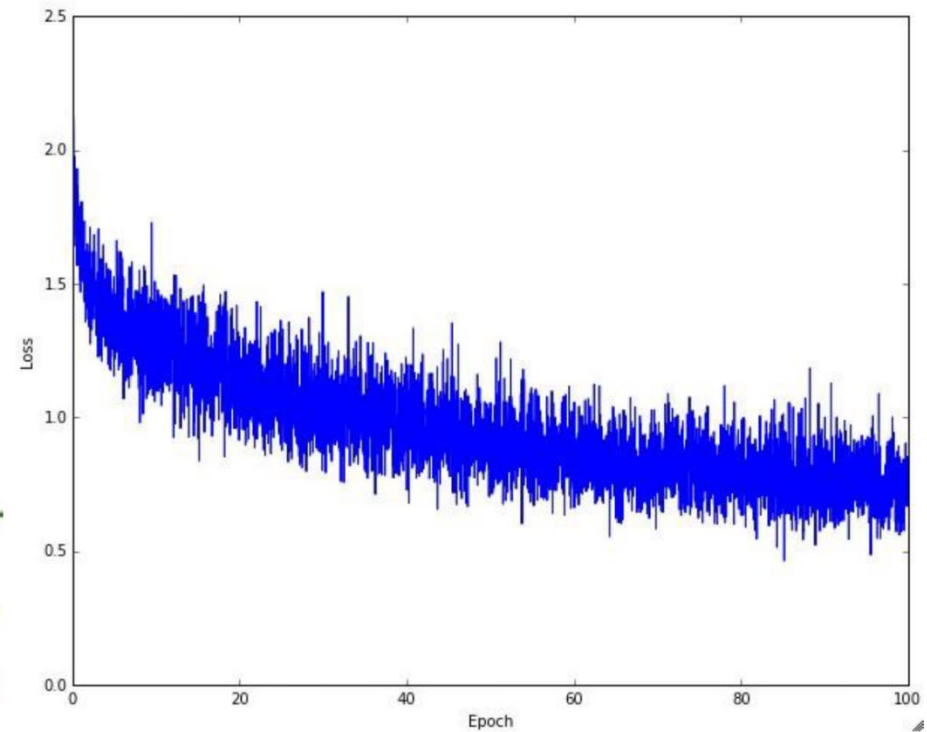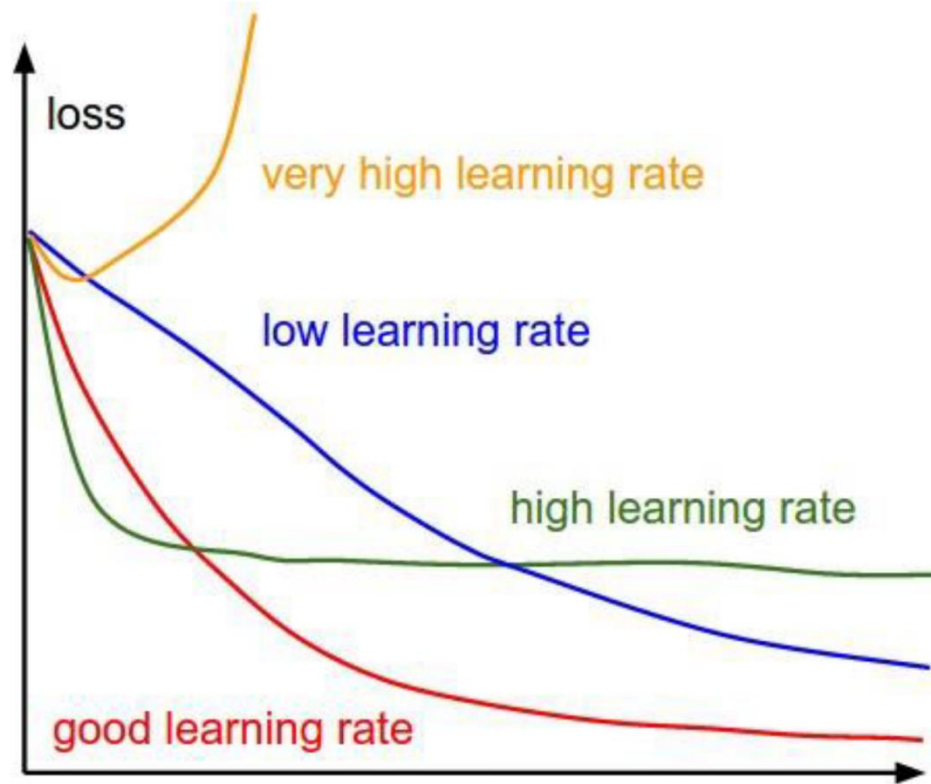
**Differentiable "unit" function!**
(or close approximation to compute "local Jacobian)

$L \qquad L = -logP(Y = y|\boldsymbol{z})$
*(cross-entropy)*

$y$

$z \qquad \boldsymbol{z} = matmult(\boldsymbol{h}_2, \boldsymbol{W}_3)$

$W_3$

$h_2 \qquad \boldsymbol{h}_2 = f(\boldsymbol{h}_1; \boldsymbol{W}_2)$

$W_2$

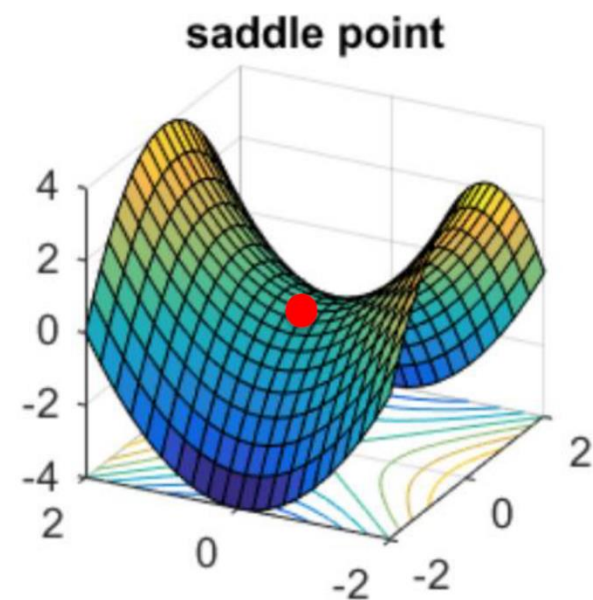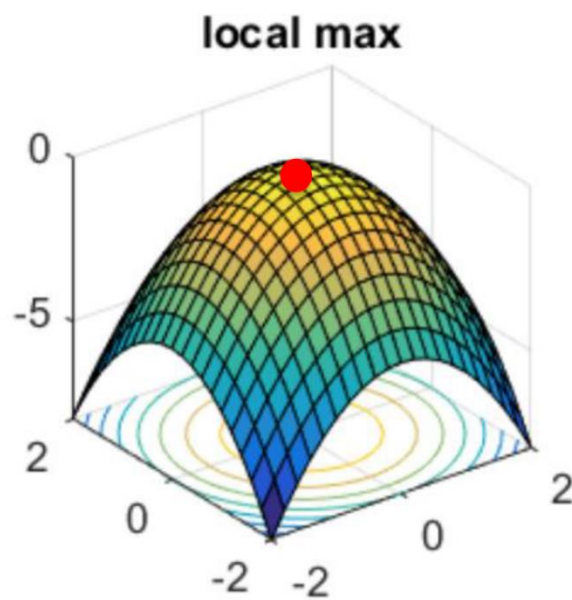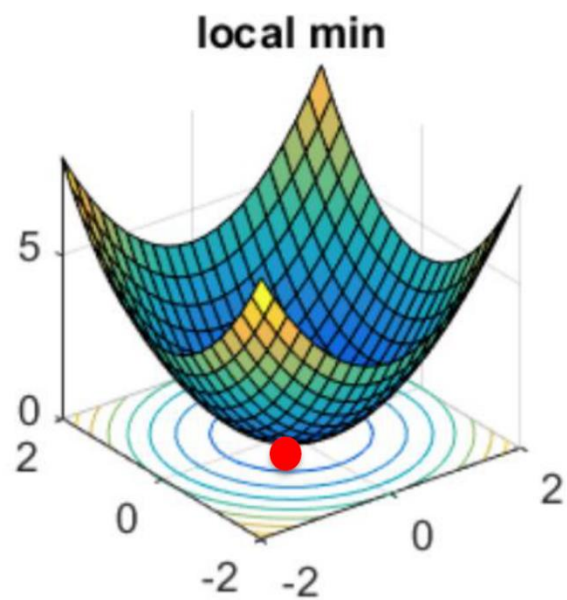$h_1 \qquad \boldsymbol{h}_1 = f(\boldsymbol{x}; \boldsymbol{W}_1)$

$W_1$

$x$

# Interpreting learning rates

# Critical Points

# Detecting Saddles

One way to detect saddles:

- Calculate Hessian at point $x$
- If Hessian is indefinite you have a saddle for sure.
- If Hessian is not indefinite you really can't tell.

"My loss isn't changing"

- You are definitely close to a critical point
  - You may be in a saddle point
  - You may be in the local minima/maxima
- One trick: quickly check the surrounding
  - Best practical trick if Hessian is not indefinite.

# Adaptive Learning Rate

**Key Idea:** Let neurons who just started learning have huge learning rate.
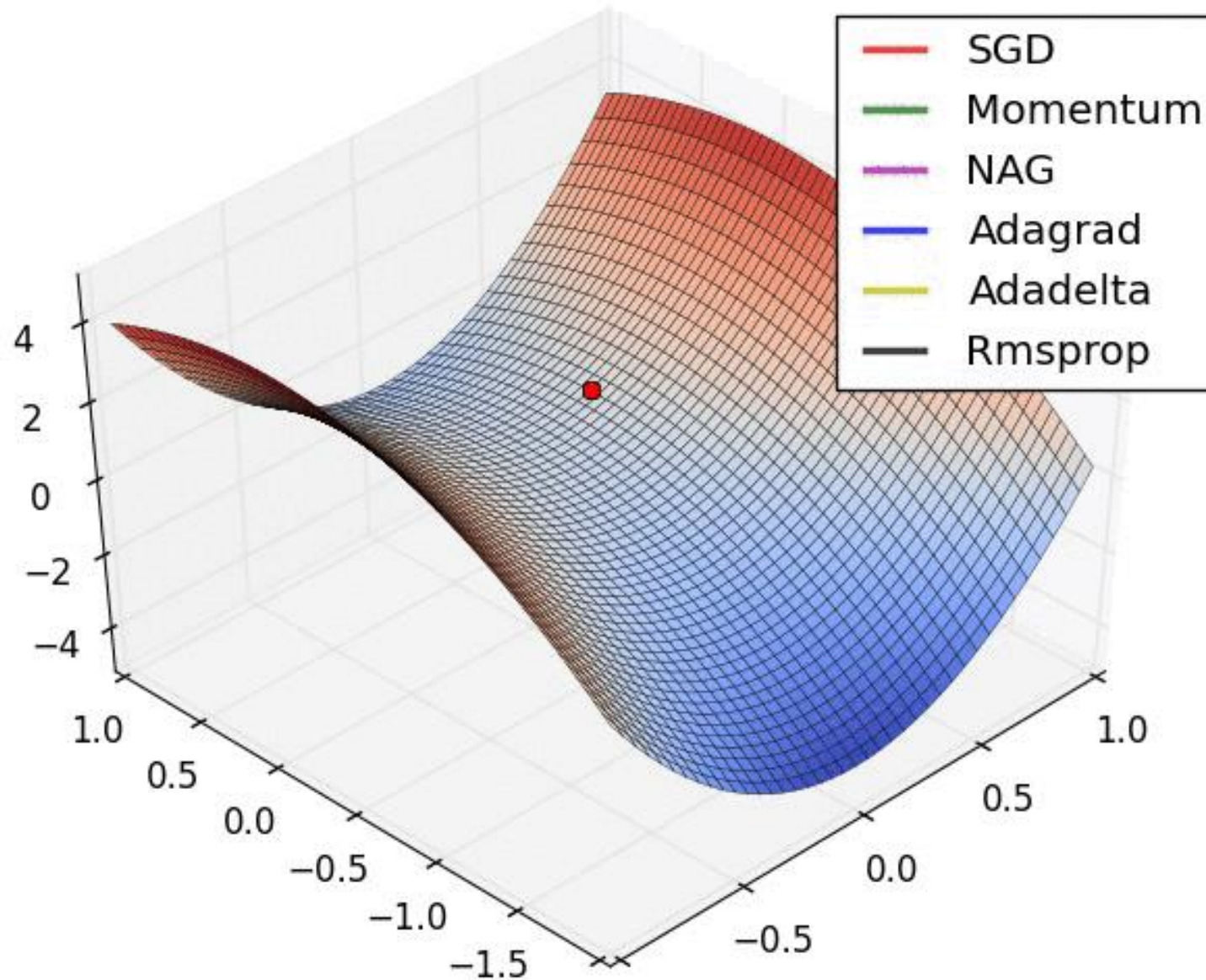
Adaptive Learning Rate is an active area of research:

- Adadelta
- RMSProp

  cache = decay_rate * cache + (1 - decay_rate) * dx**2

  x += - learning_rate * dx / (np.sqrt(cache) + eps)

- Adam

  m = beta1*m + (1-beta1)*dx

  v = beta2*v + (1-beta2)*(dx**2)

  x += - learning_rate * m / (np.sqrt(v) + eps)
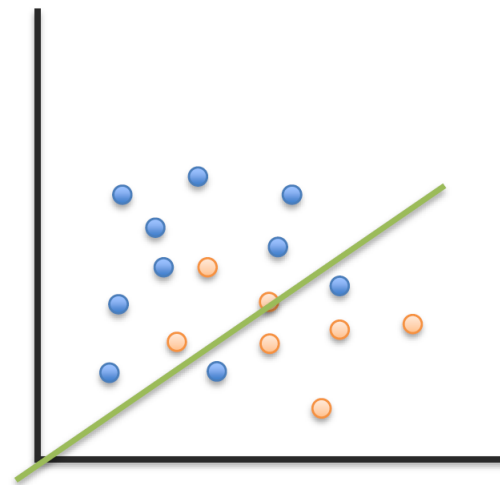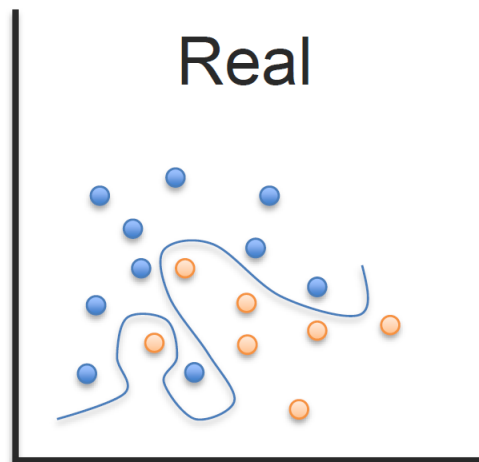
# Adaptive Learning Rate

# Bias Variance

## Problem of bias and variance

- Simple models are unlikely to find the solution to a hard problem, thus probability of finding the right model is low.

Not an issue these days!



Real

# Bias Variance

## Problem of bias and variance

- Simple models are unlikely to find the solution to a hard problem, thus probability of finding the right model is low.
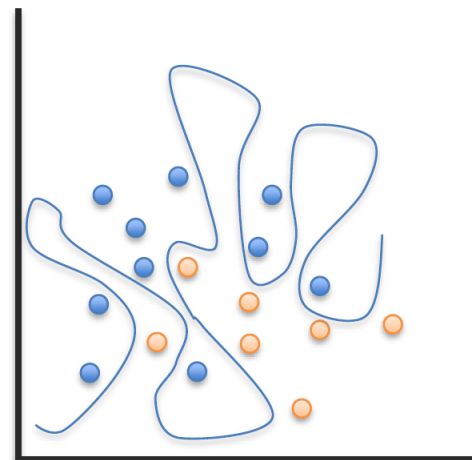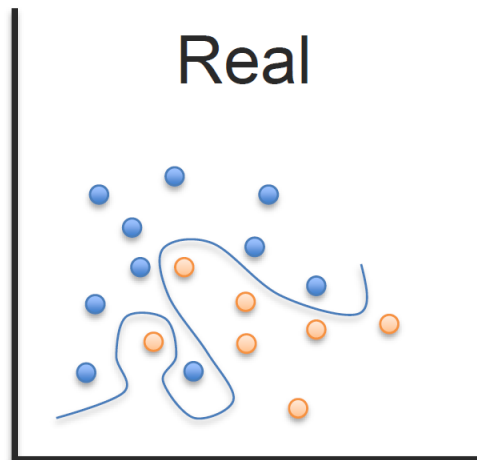
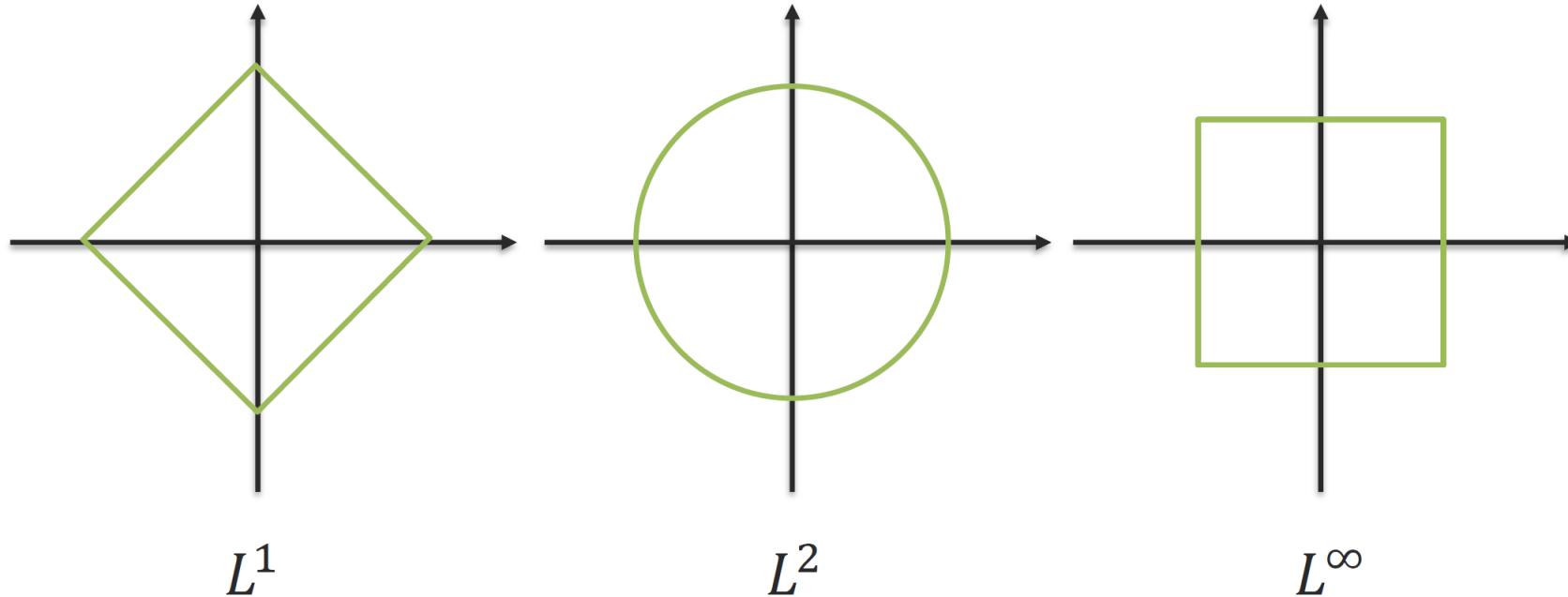- Complex models find many solutions to a problem, thus probability of finding the right model is again low.

A big issue with deep learning!

Real

## Adding prior to the network parameters

- $L^p$ Norms



$$L^1 \qquad\qquad L^2 \qquad\qquad L^\infty$$

Minimize: $Loss(x; \theta) + \propto \|\theta\|$

# Structural Regularization

Lots of models can learn everything.

- Go for simpler ones. ← Occam's razor

Take advantage of the structure and "invariances" present in each modality:

- CNNs: translation invariance
- LSTMs: sequential structure
- GRUs: sequential structure

# 总结

- 了解不同模态的基础表示

- 了解经典机器学习算法

- 掌握神经网络基本原理

- 掌握神经网络的优化实践