

# Special Topics in Deep Learning

COMP 6211D & ELEC 6910T

# Course website

<https://course.cse.ust.hk/comp6211d>

# Who we are

Instructors: Qifeng Chen (cqf@ust.hk)

TA: Hyukryul Yang (hyangbd@connect.ust.hk)  
Nayeon Lee (nayeon.lee@connect.ust.hk)

# Syllabus

## Instructor:

Week 1-2: Overview of Deep Learning: Architecture, Losses, and Optimization

## Student presentation:

Week 3-4: Convolutional Neural Networks: Dilated Convolutions, ResNet, Perceptual losses

Week 5: Deep 3D Vision: PointNet++, OctNet, Tangent convolutions

Week 6-7: Graph Convolutional Networks for Graph Processing and Optimization

Week 8-9: Sequential Modelling and Signal Processing: RNN, LSTM, TCN, and WaveNet

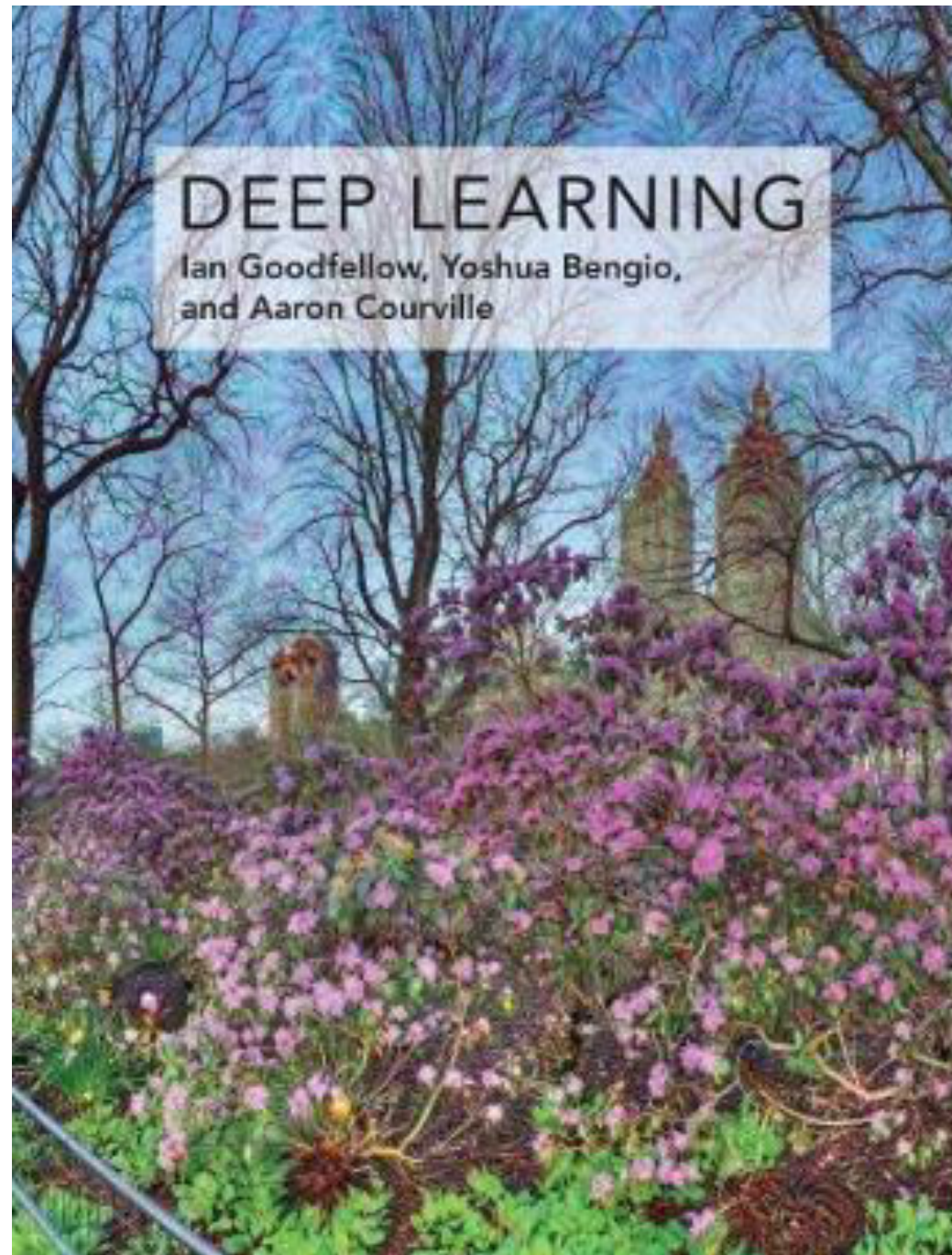
Week 10-11: Generative Models: GAN, Pix2pix, CycleGAN, CRN, VAE

Week 12-13: Final project presentation and project report due

# Grading

Class participation:	10%
In-class presentation:	15%
Homework:	30%
Final project:	45%

# Deep Learning Book



# Representations Matter

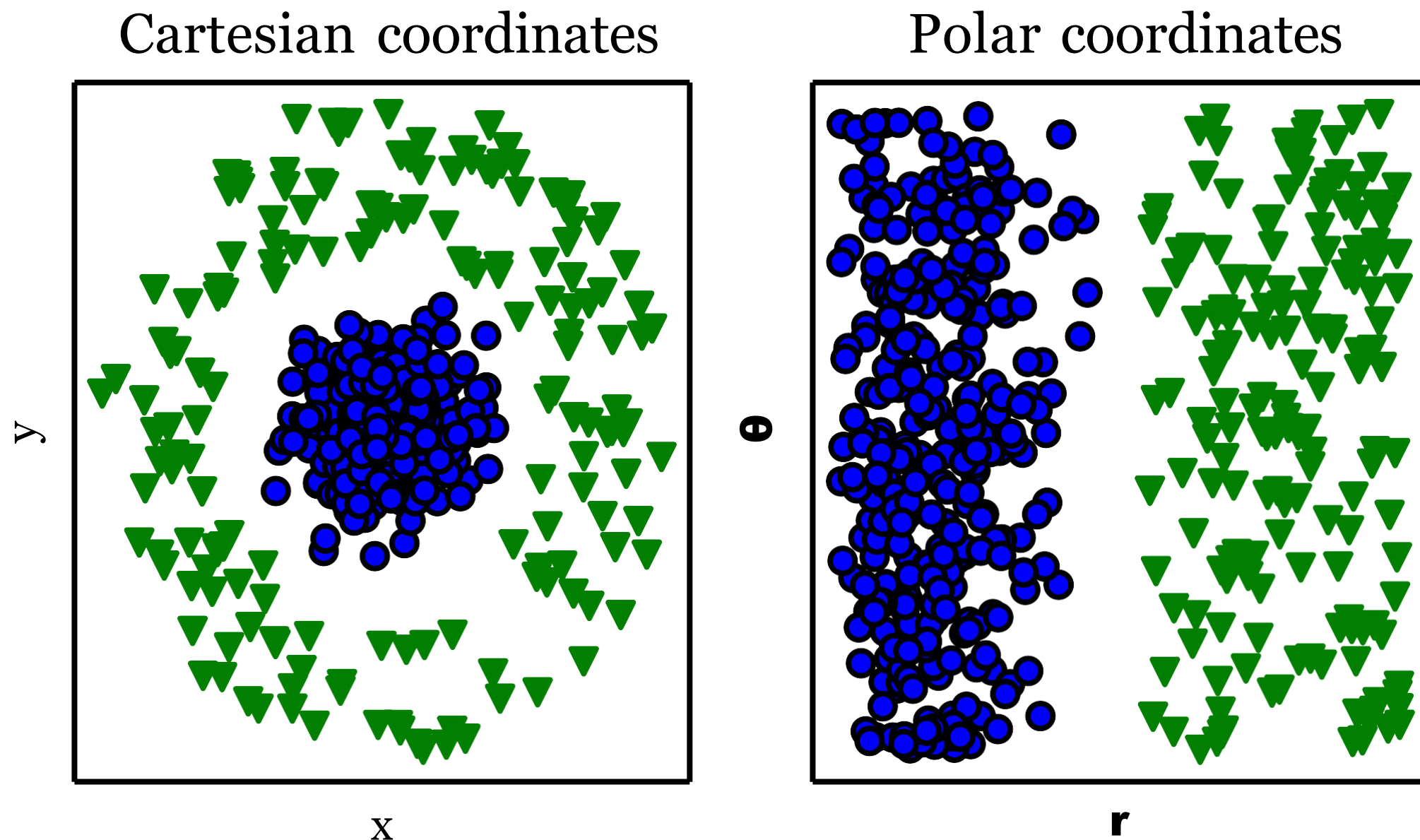


Figure 1.1



# Depth: Repeated Composition

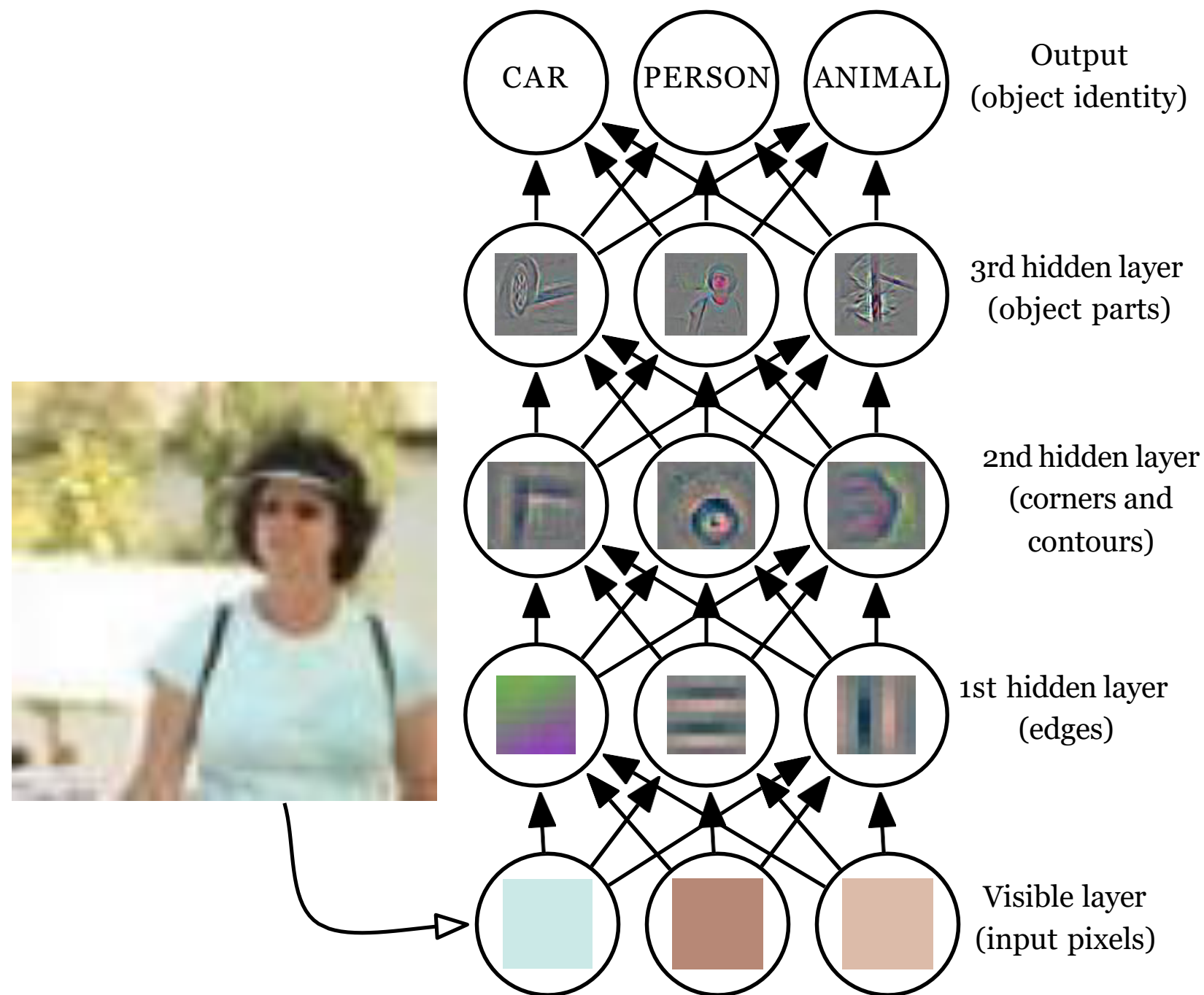


Figure 1.2



# Computational Graphs

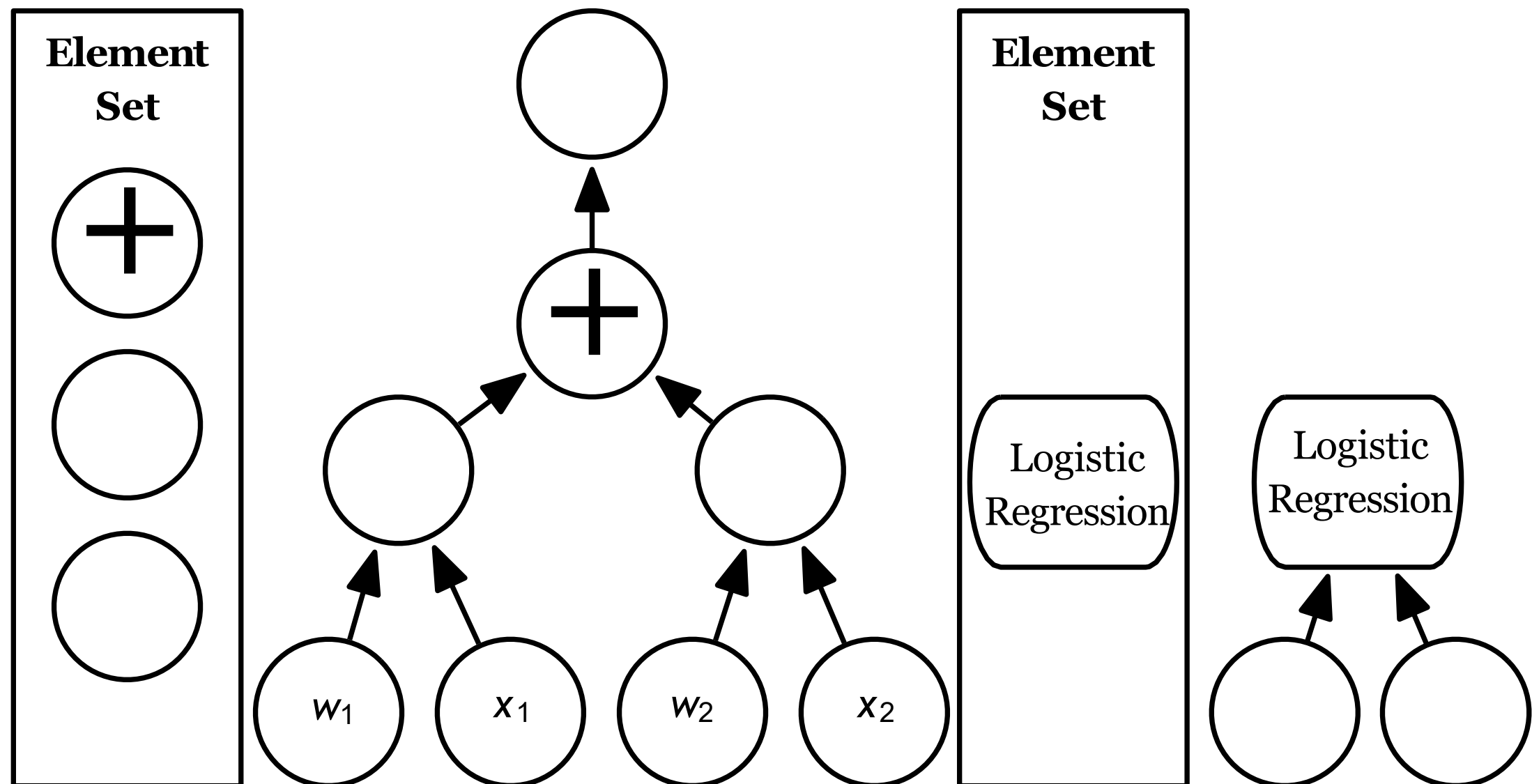


Figure 1.3

# Machine Learning and AI

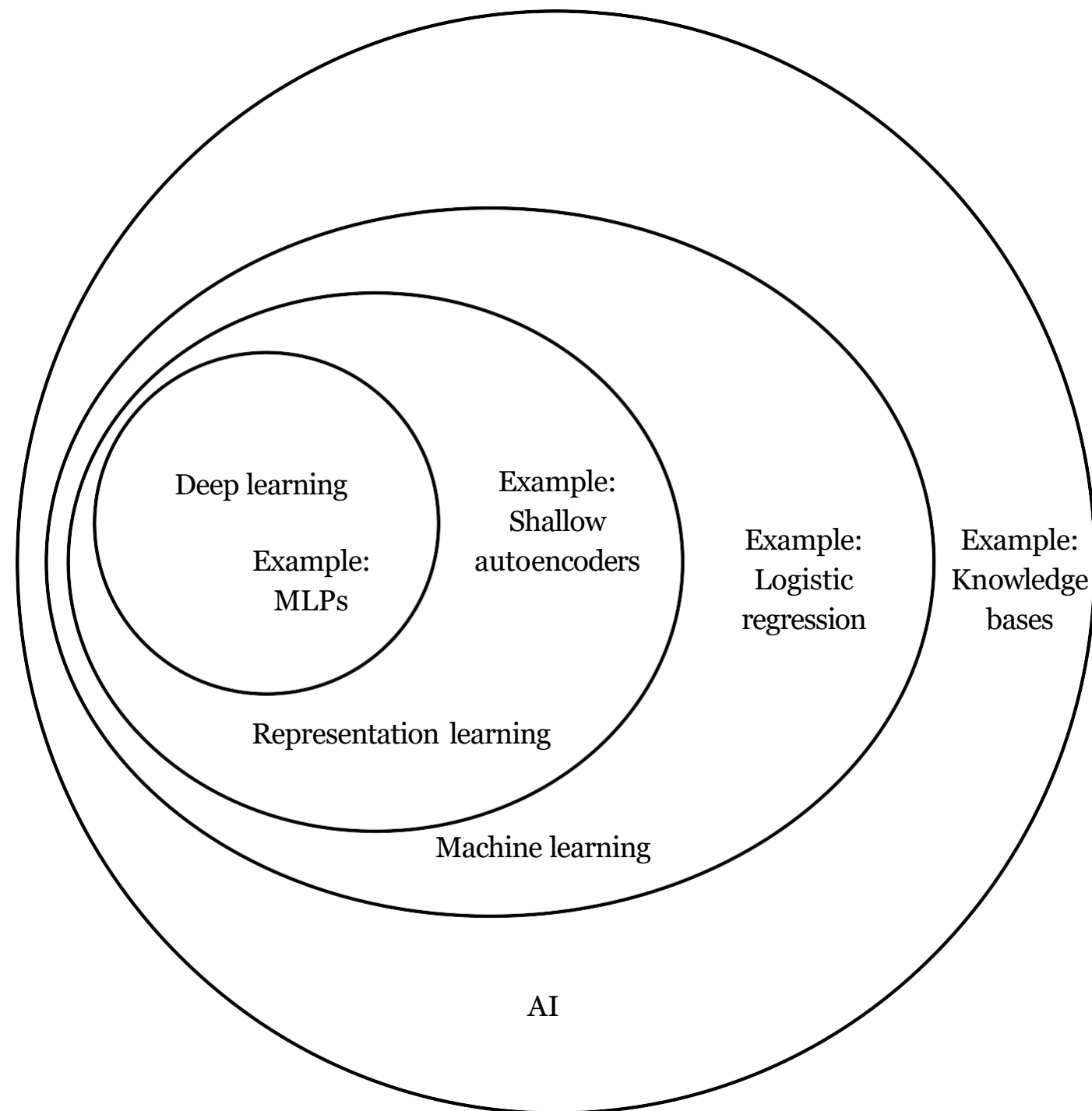
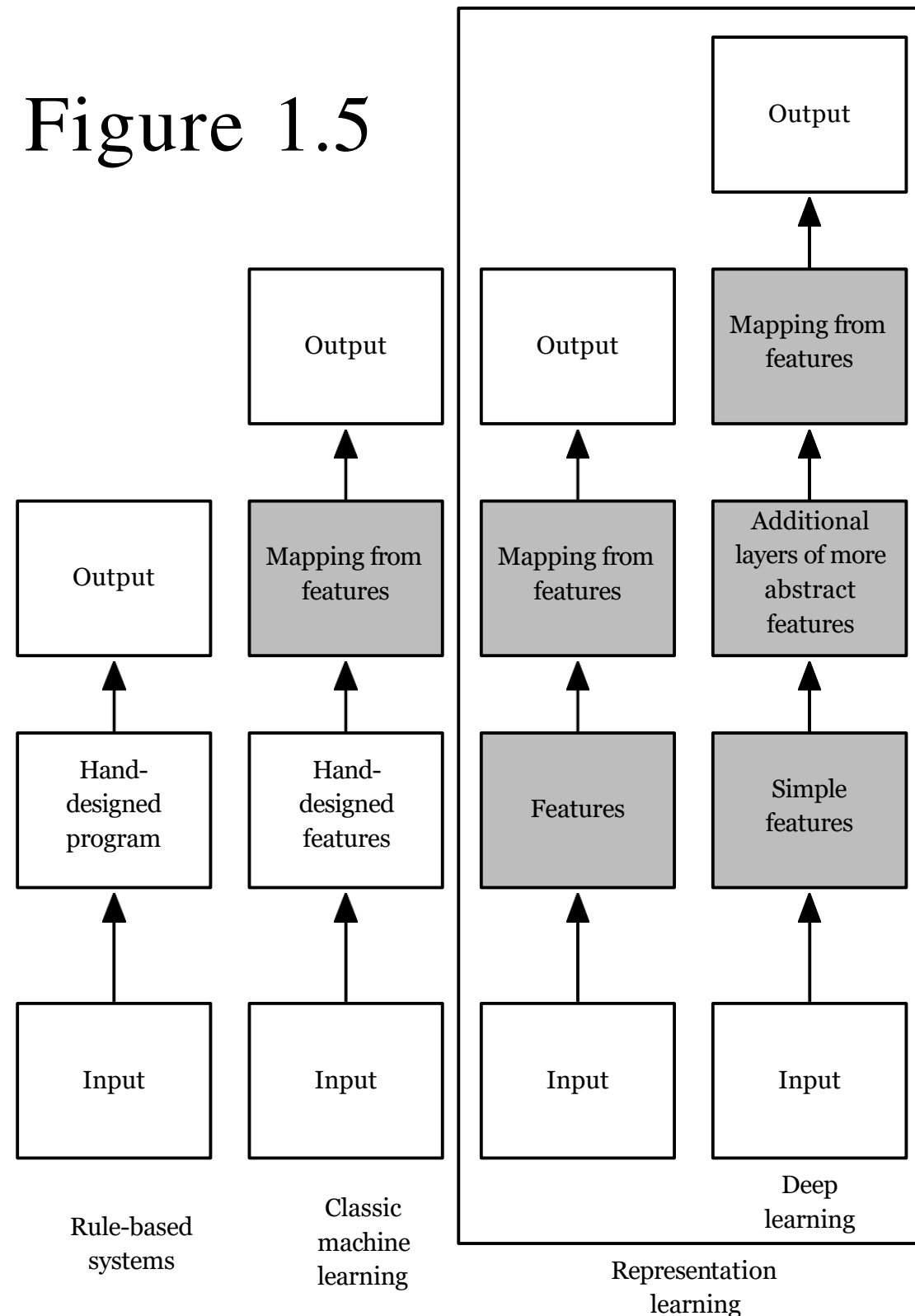


Figure 1.4

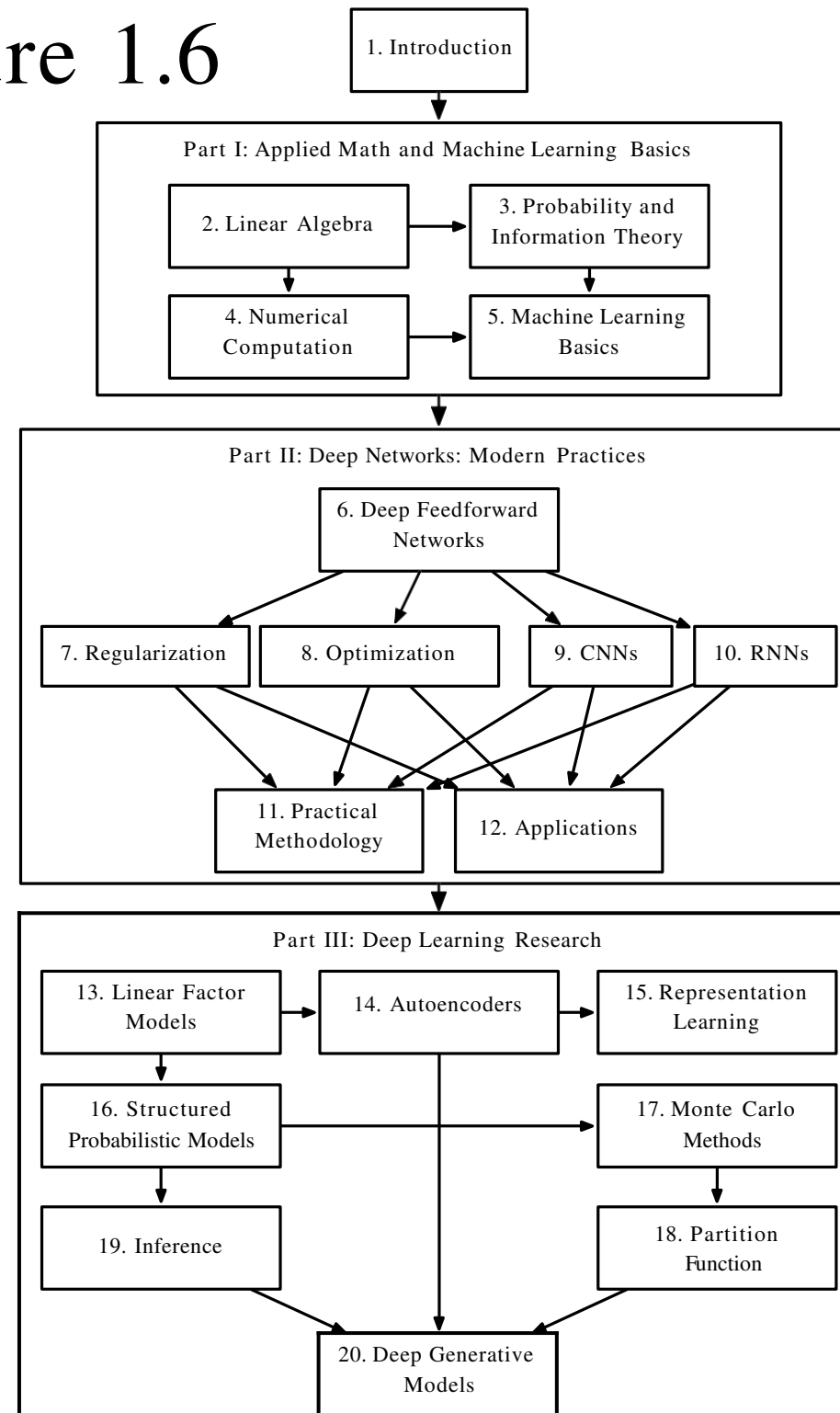
# Learning Multiple Components

Figure 1.5



# Organization of the Book

Figure 1.6



# Historical Waves

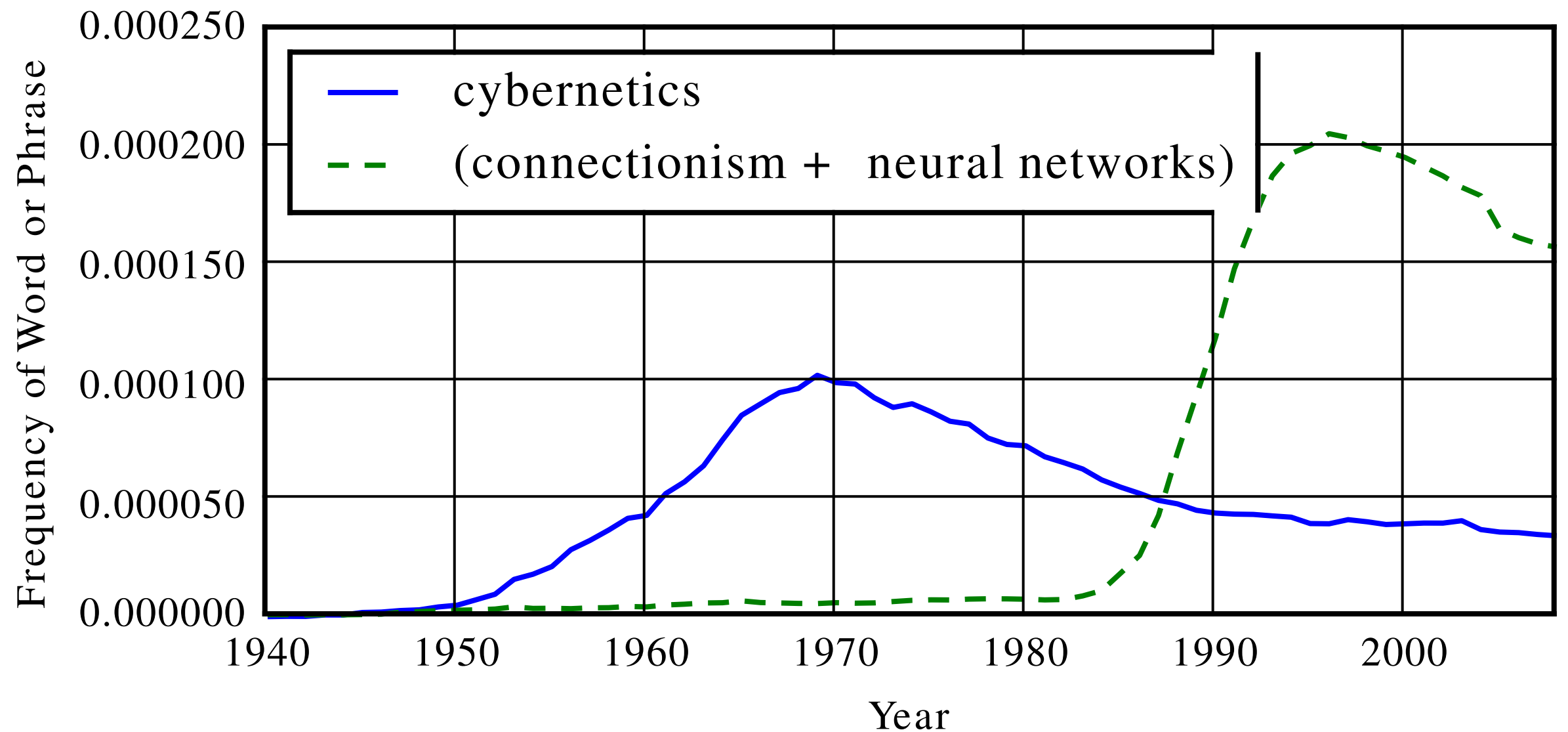


Figure 1.7

# Historical Trends: Growing Datasets

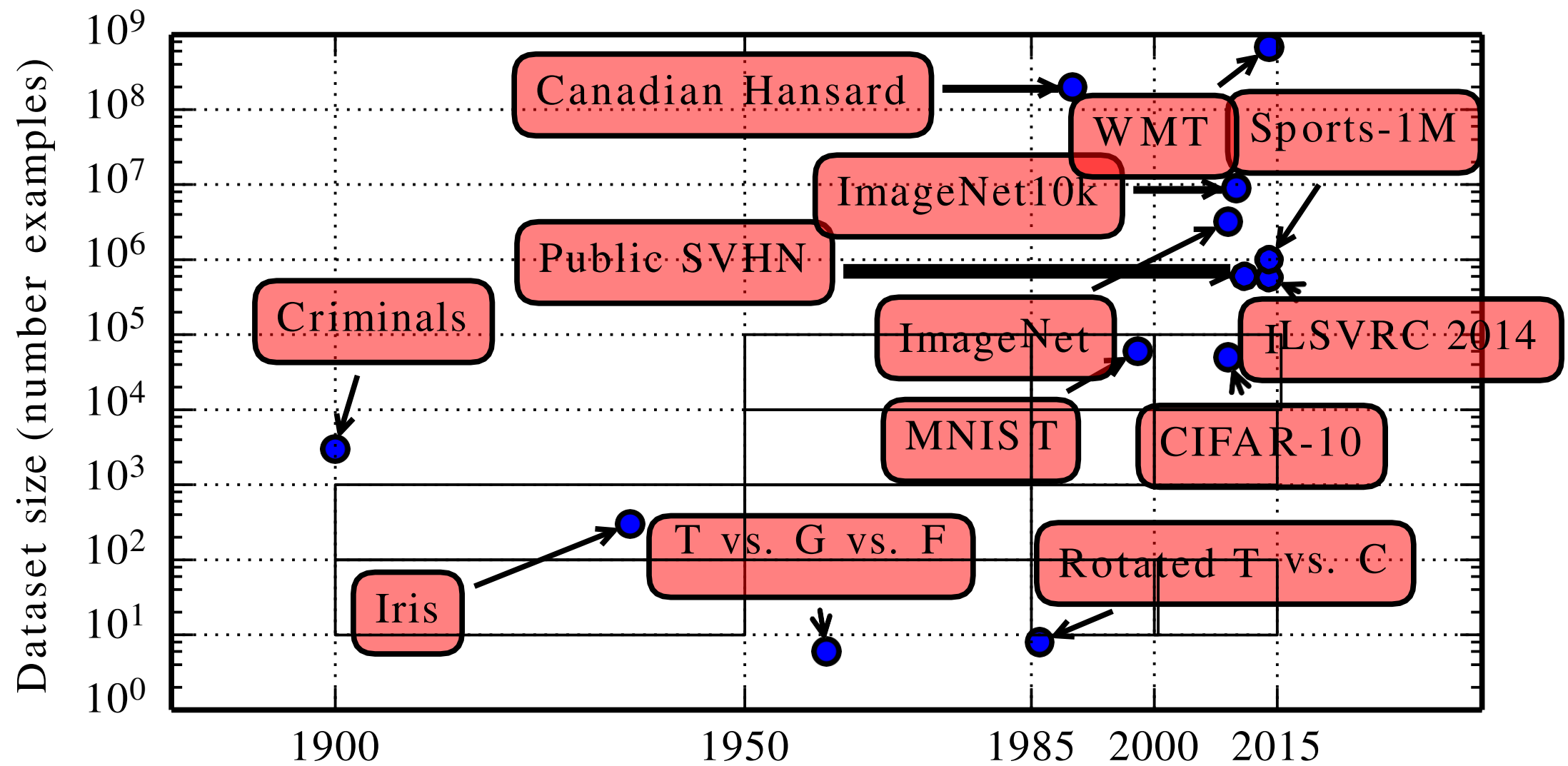


Figure 1.8

# The MNIST Dataset

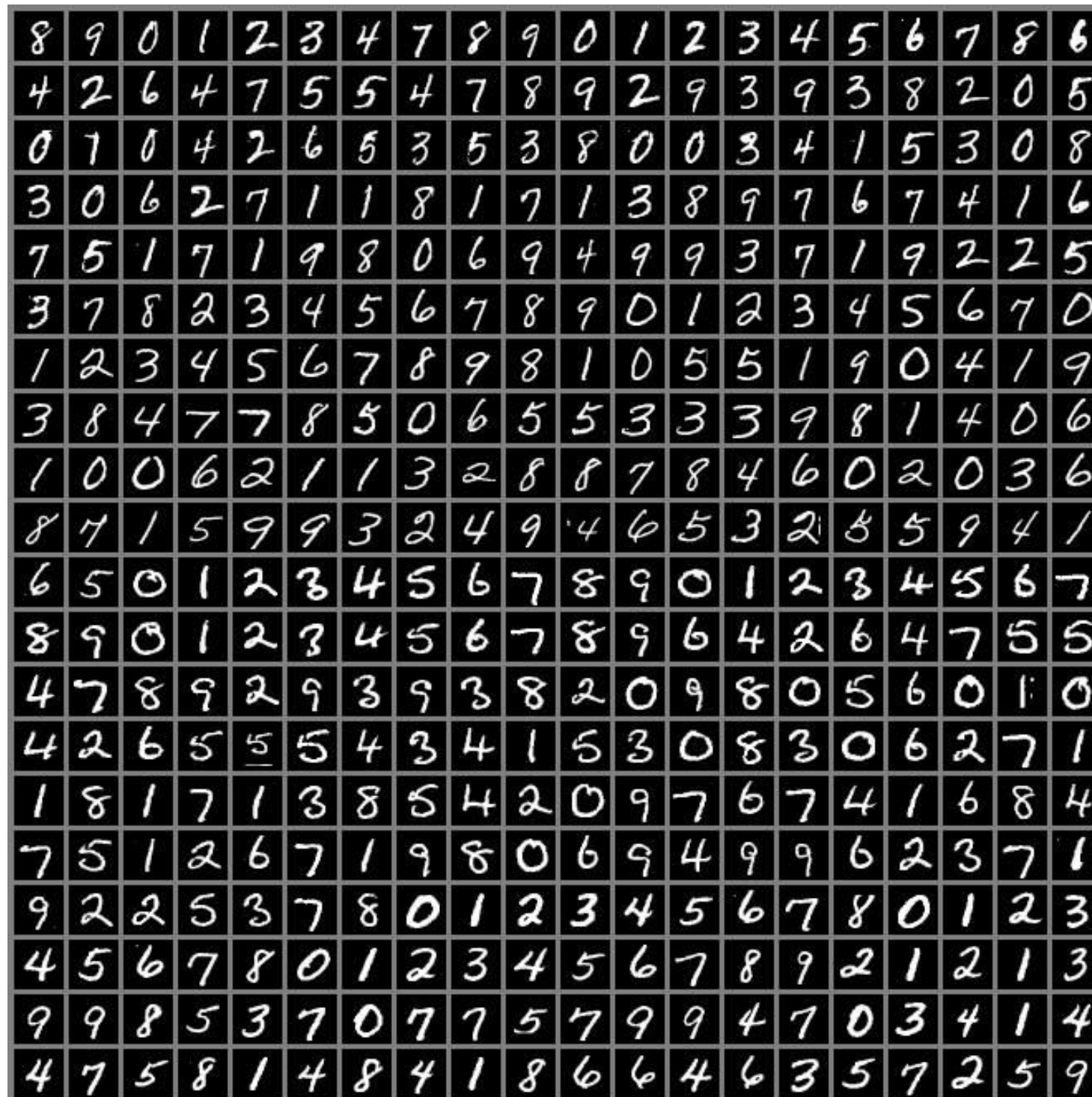


Figure 1.9



# Connections per Neuron

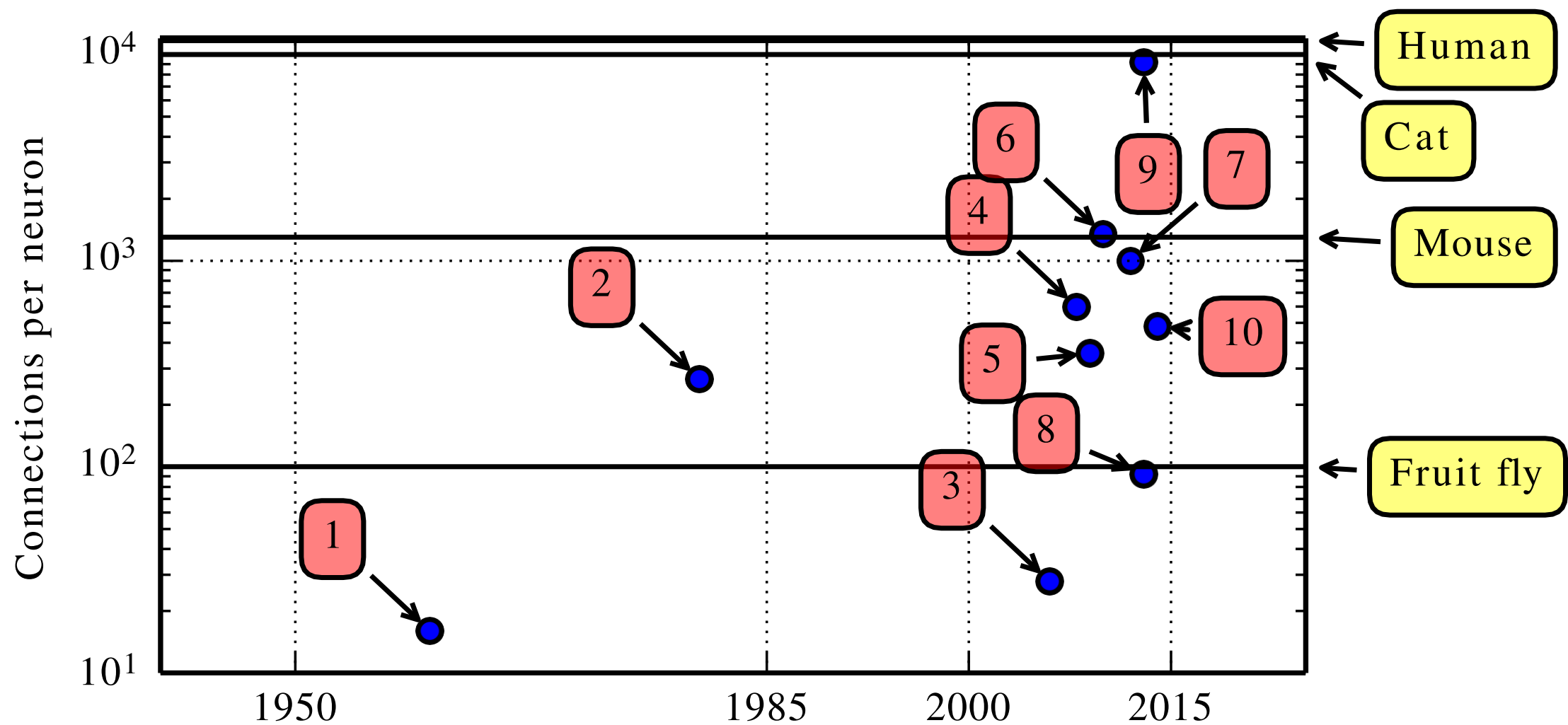


Figure 1.10

# Number of Neurons

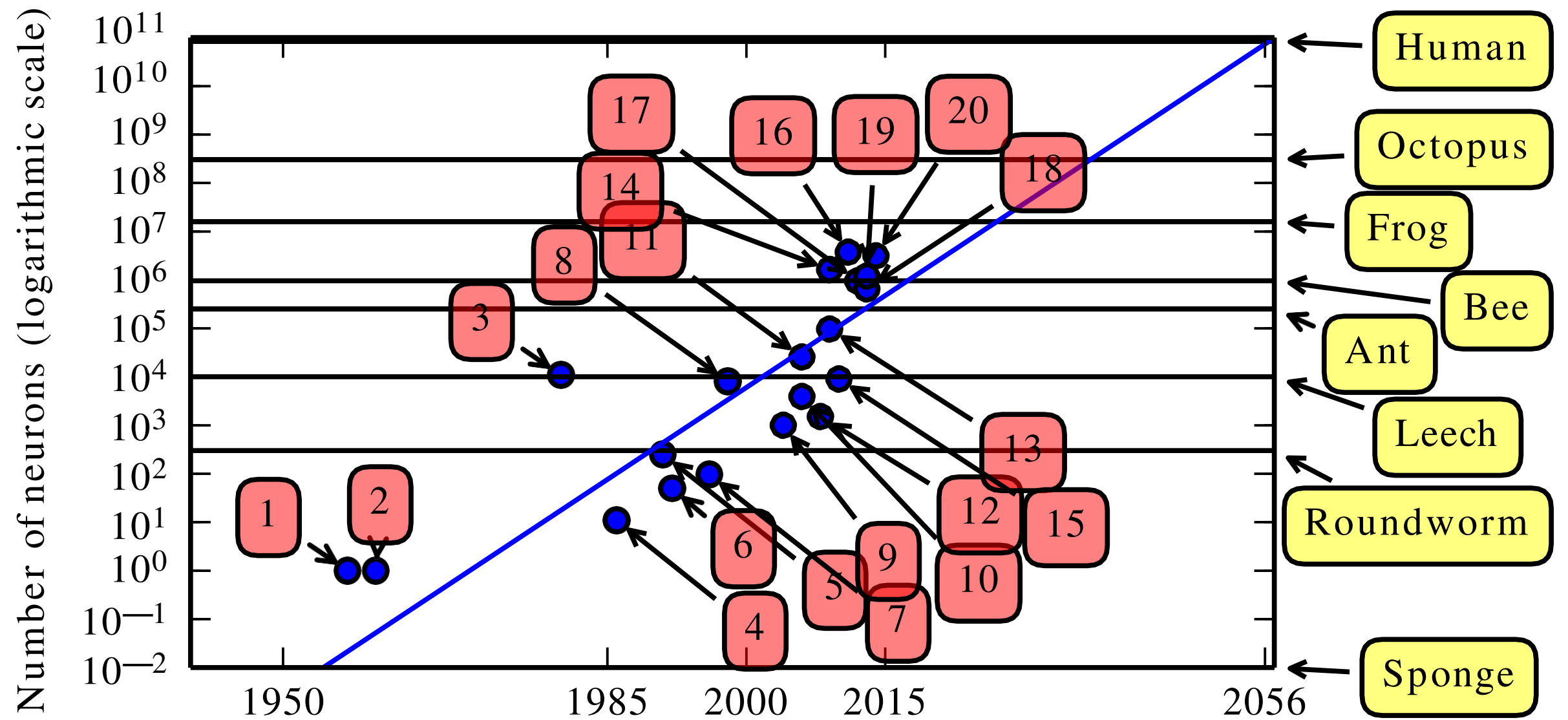


Figure 1.11

# Solving Object Recognition

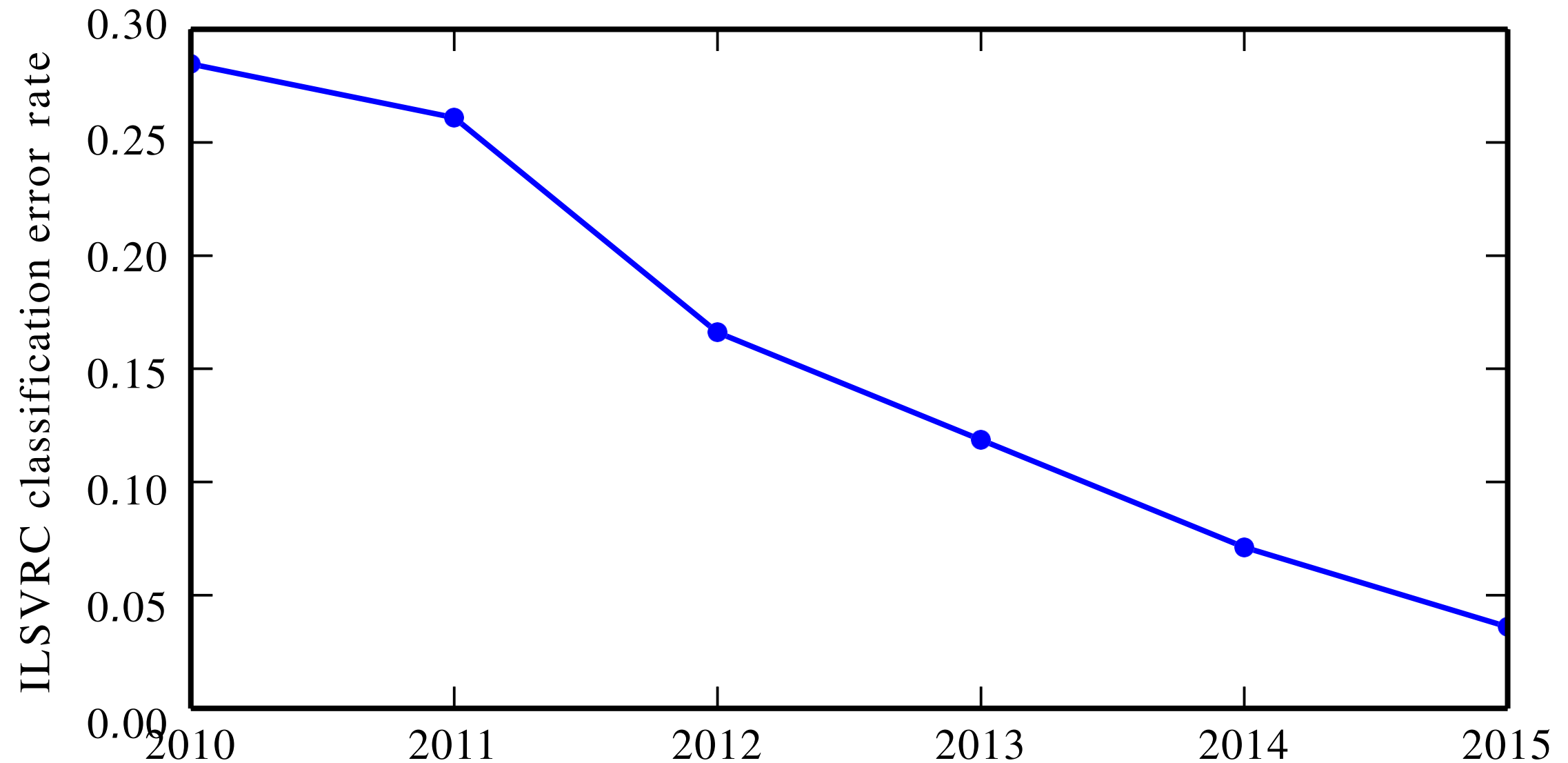


Figure 1.12

# Linear Algebra

Lecture slides for Chapter 2 of *Deep Learning*

Ian Goodfellow

2016-06-24

# About this chapter

- Not a comprehensive survey of all of linear algebra
- Focused on the subset most relevant to deep learning
- Larger subset: e.g., *Linear Algebra* by Georgi Shilov

# Scalars

- A scalar is a single number
- Integers, real numbers, rational numbers, etc.
- We denote it with italic font:

*a, n, x*

# Vectors

- A vector is a 1-D array of numbers:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}. \quad (2.1)$$

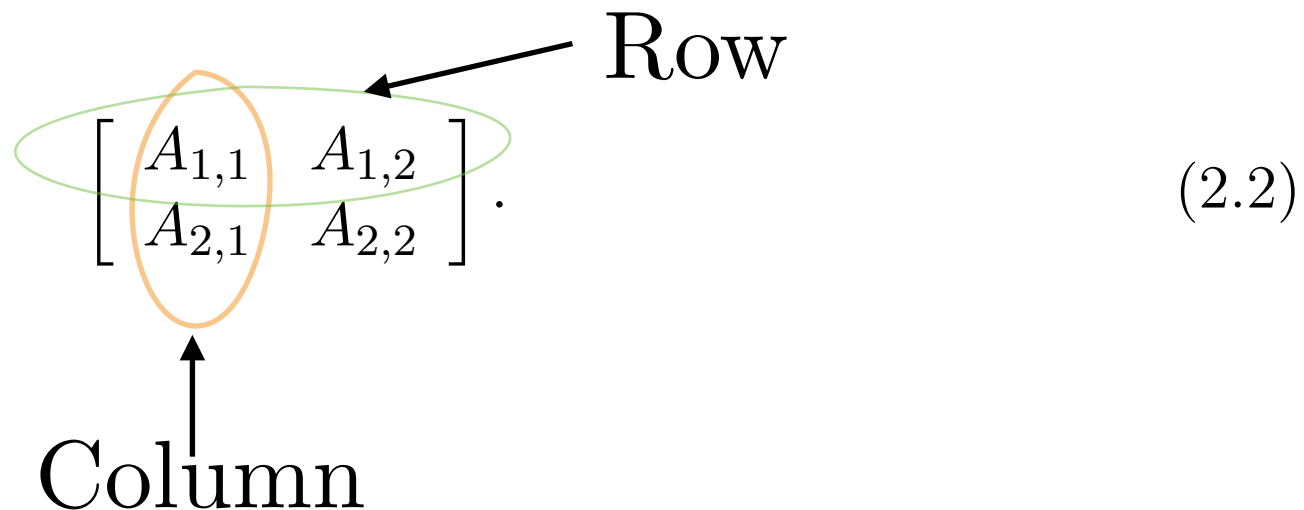
- Can be real, binary, integer, etc.
- Example notation for type and size:

$$\mathbb{R}^n$$



# Matrices

- A matrix is a 2-D array of numbers:



The diagram shows a 2x2 matrix  $\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$ . A green oval encircles the top row, with an arrow pointing to it from the word "Row". An orange oval encircles the left column, with an arrow pointing to it from the word "Column". The matrix is followed by the label (2.2).

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}. \quad (2.2)$$

- Example notation for type and shape:

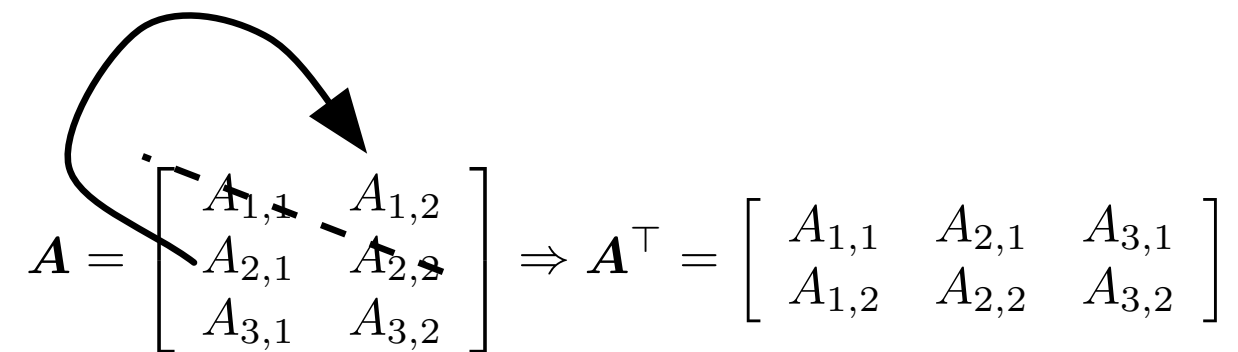
$$\mathbf{A} \in \mathbb{R}^{m \times n}$$

# Tensors

- A tensor is an array of numbers, that may have
  - zero dimensions, and be a scalar
  - one dimension, and be a vector
  - two dimensions, and be a matrix
  - or more dimensions.

# Matrix Transpose

$$(\mathbf{A}^\top)_{i,j} = A_{j,i}. \quad (2.3)$$



$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow \mathbf{A}^\top = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

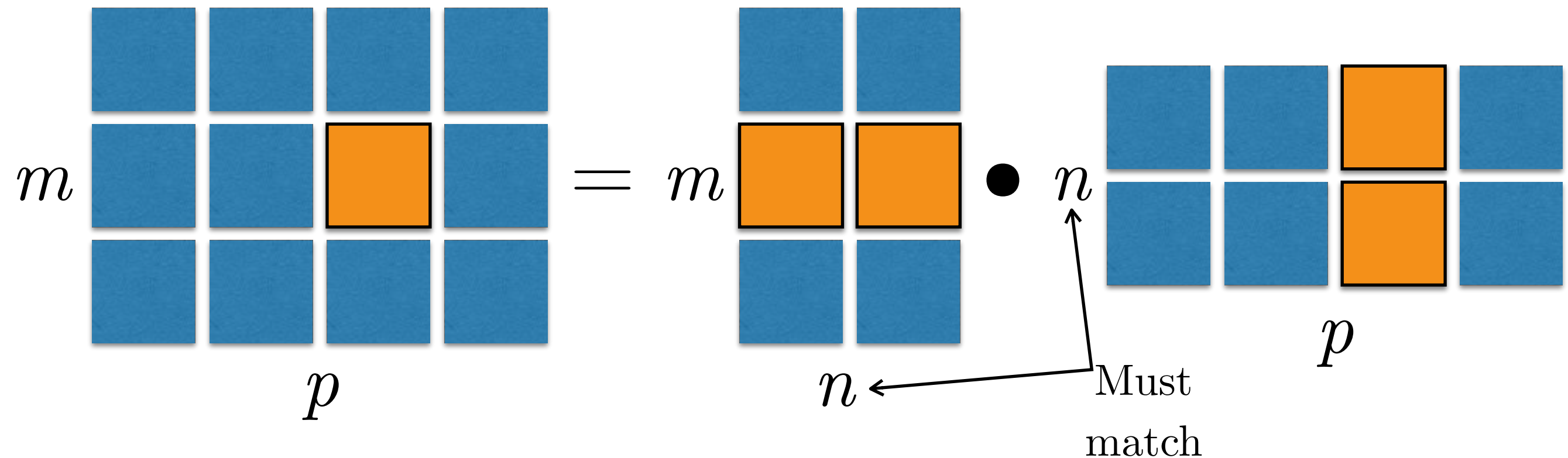
Figure 2.1: The transpose of the matrix can be thought of as a mirror image across the main diagonal.

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top. \quad (2.9)$$

# Matrix (Dot) Product

$$\mathbf{C} = \mathbf{A}\mathbf{B}. \quad (2.4)$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}. \quad (2.5)$$



# Identity Matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Figure 2.2: *Example identity matrix:* This is  $\mathbf{I}_3$ .

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{I}_n \mathbf{x} = \mathbf{x}. \tag{2.20}$$

# Systems of Equations

$$\mathbf{A}\mathbf{x} = \mathbf{b} \tag{2.11}$$

expands to

$$\mathbf{A}_{1,:}\mathbf{x} = b_1 \tag{2.12}$$

$$\mathbf{A}_{2,:}\mathbf{x} = b_2 \tag{2.13}$$

$$\dots \tag{2.14}$$

$$\mathbf{A}_{m,:}\mathbf{x} = b_m \tag{2.15}$$

# Solving Systems of Equations

- A linear system of equations can have:
  - No solution
  - Many solutions
  - Exactly one solution: this means multiplication by the matrix is an invertible function



# Matrix Inversion

- Matrix inverse:

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}_n. \quad (2.21)$$

- Solving a system using an inverse:

$$\mathbf{A} \mathbf{x} = \mathbf{b} \quad (2.22)$$

$$\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.23)$$

$$\mathbf{I}_n \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.24)$$

- Numerically unstable, but useful for abstract analysis

# Invertibility

- Matrix can't be inverted if...
  - More rows than columns
  - More columns than rows
  - Redundant rows/columns (“linearly dependent”, “low rank”)

# Norms

- Functions that measure how “large” a vector is
- Similar to a distance between zero and the point represented by the vector
  - $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
  - $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$  (the *triangle inequality*)
  - $\forall \alpha \in \mathbb{R}, f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$

# Norms

- $L^p$  norm

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

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

- L1 norm,  $p=1$ :  $||\mathbf{x}||_1 = \sum_i |x_i|.$  (2.31)

- Max norm, infinite  $p$ :  $||\mathbf{x}||_\infty = \max_i |x_i|.$  (2.32)

# Special Matrices and Vectors

- Unit vector:

$$||\boldsymbol{x}||_2 = 1. \quad (2.36)$$

- Symmetric Matrix:

$$\boldsymbol{A} = \boldsymbol{A}^\top. \quad (2.35)$$

- Orthogonal matrix:

$$\begin{aligned} \boldsymbol{A}^\top \boldsymbol{A} &= \boldsymbol{A} \boldsymbol{A}^\top = \boldsymbol{I}. \\ \boldsymbol{A}^{-1} &= \boldsymbol{A}^\top \end{aligned} \quad (2.37)$$

# Eigendecomposition

- Eigenvector and eigenvalue:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}. \quad (2.39)$$

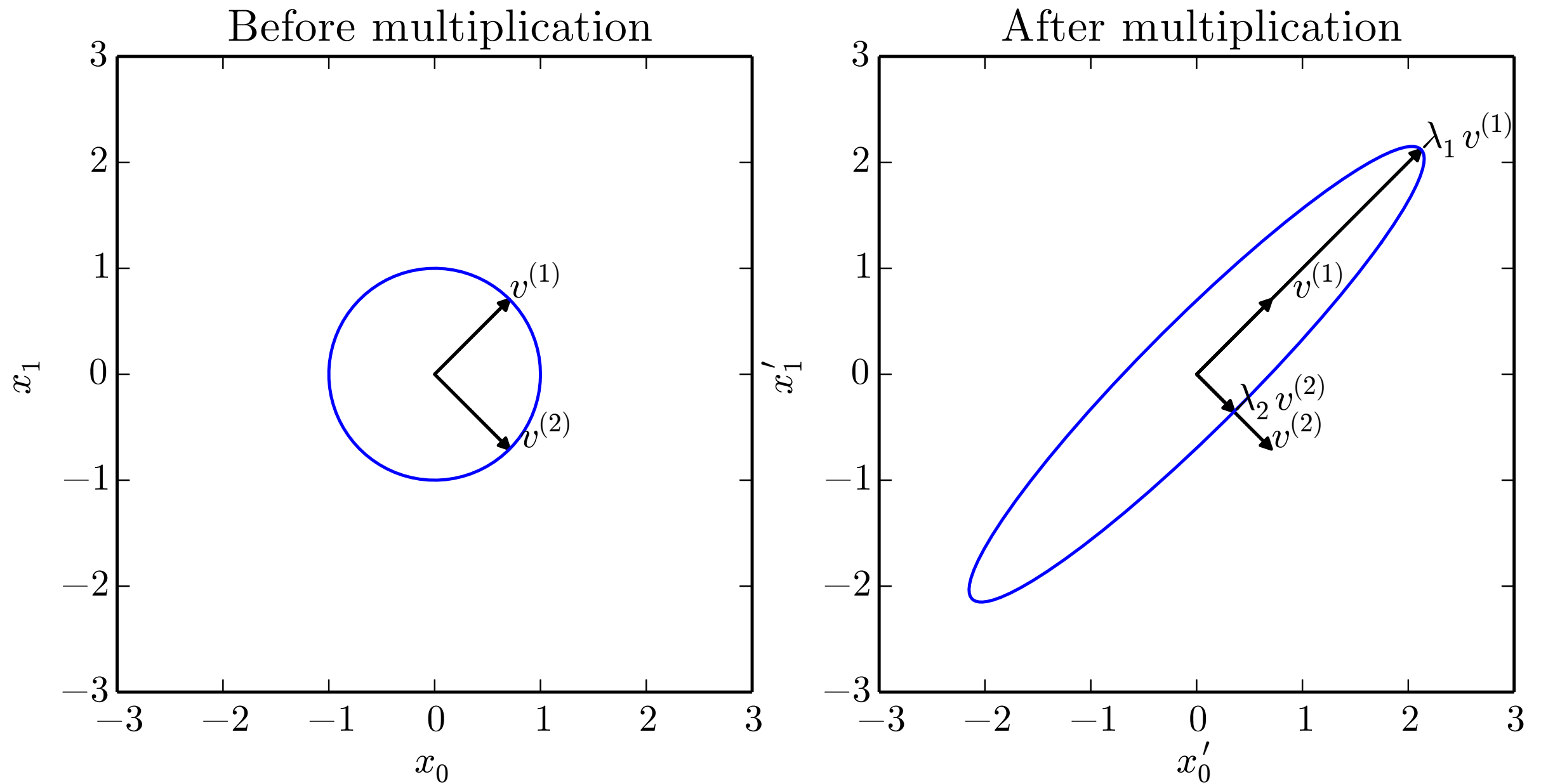
- Eigendecomposition of a diagonalizable matrix:

$$\mathbf{A} = \mathbf{V}\text{diag}(\boldsymbol{\lambda})\mathbf{V}^{-1}. \quad (2.40)$$

- Every real symmetric matrix has a real, orthogonal eigendecomposition:

$$\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^{\top} \quad (2.41)$$

# Effect of Eigenvalues





# Singular Value Decomposition

- Similar to eigendecomposition
- More general; matrix need not be square

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^{\top}. \tag{2.43}$$

# Moore-Penrose Pseudoinverse

$$\mathbf{x} = \mathbf{A}^+ \mathbf{y}$$

- If the equation has:
  - Exactly one solution: this is the same as the inverse.
  - No solution: this gives us the solution with the smallest error  $\|\mathbf{Ax} - \mathbf{y}\|_2$ .
  - Many solutions: this gives us the solution with the smallest norm of  $\mathbf{x}$ .

# Computing the Pseudoinverse

The SVD allows the computation of the pseudoinverse:

$$\mathbf{A}^+ = \mathbf{V} \mathbf{D}^+ \mathbf{U}^\top, \quad (2.47)$$



Take reciprocal of non-zero entries

# Trace

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}. \quad (2.48)$$

$$\text{Tr}(\mathbf{A}\mathbf{B}\mathbf{C}) = \text{Tr}(\mathbf{C}\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{C}\mathbf{A}) \quad (2.51)$$

# Learning linear algebra

- Do a lot of practice problems
- Start out with lots of summation signs and indexing into individual entries
- Eventually you will be able to mostly use matrix and vector product notation quickly and easily

# Probability and Information Theory

Lecture slides for Chapter 3 of *Deep Learning*

[www.deeplearningbook.org](http://www.deeplearningbook.org)

Ian Goodfellow

2016-09-26

# Probability Mass Function

- The domain of  $P$  must be the set of all possible states of  $\mathbf{x}$ .
- $\forall x \in \mathbf{x}, 0 \leq P(x) \leq 1$ . An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring.
- $\sum_{x \in \mathbf{x}} P(x) = 1$ . We refer to this property as being **normalized**. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring.

Example: uniform distribution:  $P(\mathbf{x} = x_i) = \frac{1}{k}$

# Probability Density Function

- The domain of  $p$  must be the set of all possible states of  $\mathbf{x}$ .
- $\forall x \in \mathbf{x}, p(x) \geq 0$ . Note that we do not require  $p(x) \leq 1$ .
- $\int p(x)dx = 1$ .

Example: uniform distribution:  $u(x; a, b) = \frac{1}{b-a}$ .



# Computing Marginal Probability with the Sum Rule

$$\forall x \in \mathbf{x}, P(\mathbf{x} = x) = \sum_y P(\mathbf{x} = x, y = y). \quad (3.3)$$

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

# Conditional Probability

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}. \quad (3.5)$$

# Chain Rule of Probability

$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}) = P(\mathbf{x}^{(1)}) \prod_{i=2}^n P(\mathbf{x}^{(i)} \mid \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i-1)}). \quad (3.6)$$

# Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, \quad p(\mathbf{x} = x, \mathbf{y} = y) = p(\mathbf{x} = x)p(\mathbf{y} = y). \quad (3.7)$$

# Conditional Independence

$$\forall x \in \mathbf{x}, y \in \mathbf{y}, z \in \mathbf{z}, \quad p(\mathbf{x} = x, \mathbf{y} = y \mid \mathbf{z} = z) = p(\mathbf{x} = x \mid \mathbf{z} = z)p(\mathbf{y} = y \mid \mathbf{z} = z). \quad (3.8)$$

# Expectation

$$\mathbb{E}_{\mathbf{x} \sim P}[f(x)] = \sum_x P(x) f(x), \quad (3.9)$$

$$\mathbb{E}_{\mathbf{x} \sim p}[f(x)] = \int p(x) f(x) dx. \quad (3.10)$$

linearity of expectations:

$$\mathbb{E}_{\mathbf{x}}[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_{\mathbf{x}}[f(x)] + \beta \mathbb{E}_{\mathbf{x}}[g(x)], \quad (3.11)$$

# Variance and Covariance

$$\text{Var}(f(x)) = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right]. \quad (3.12)$$

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)]) (g(y) - \mathbb{E}[g(y)])]. \quad (3.13)$$

Covariance matrix:

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j). \quad (3.14)$$

# Bernoulli Distribution

$$P(\mathbf{x} = 1) = \phi \tag{3.16}$$

$$P(\mathbf{x} = 0) = 1 - \phi \tag{3.17}$$

$$P(\mathbf{x} = x) = \phi^x (1 - \phi)^{1-x} \tag{3.18}$$

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}] = \phi \tag{3.19}$$

$$\text{Var}_{\mathbf{x}}(\mathbf{x}) = \phi(1 - \phi) \tag{3.20}$$



# Gaussian Distribution

Parametrized by variance:

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

Parametrized by precision:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

# Gaussian Distribution

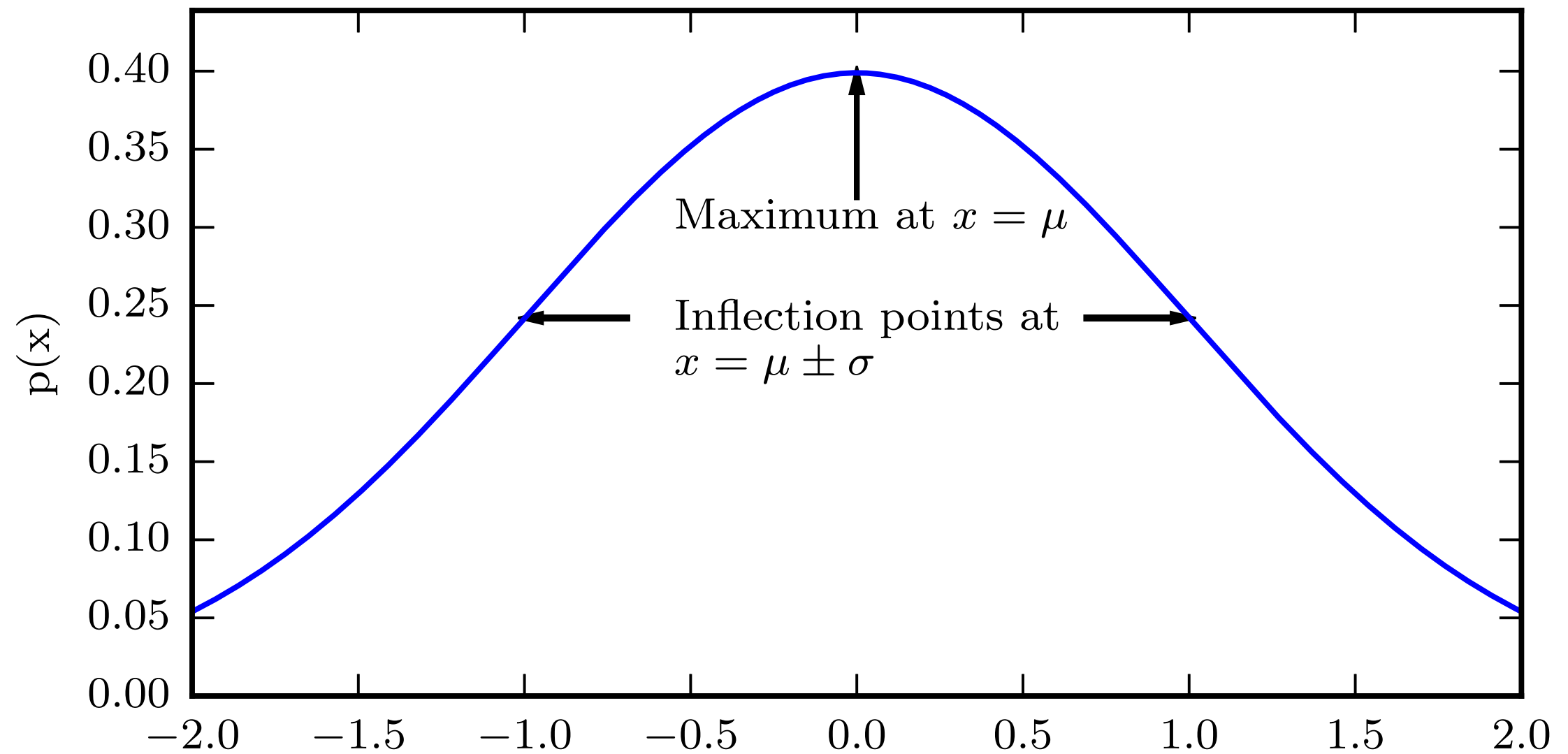


Figure 3.1

# Multivariate Gaussian

Parametrized by covariance matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right). \quad (3.23)$$

Parametrized by precision matrix:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\beta} (\mathbf{x} - \boldsymbol{\mu}) \right). \quad (3.24)$$

# More Distributions

Exponential:

$$p(x; \lambda) = \lambda \mathbf{1}_{x \geq 0} \exp(-\lambda x). \quad (3.25)$$

Laplace:

$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right). \quad (3.26)$$

Dirac:

$$p(x) = \delta(x - \mu). \quad (3.27)$$

# Empirical Distribution

$$\hat{p}(\boldsymbol{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\boldsymbol{x} - \boldsymbol{x}^{(i)}) \quad (3.28)$$

# Mixture Distributions

$$P(\mathbf{x}) = \sum_i P(c = i) P(\mathbf{x} \mid c = i) \quad (3.29)$$

Gaussian mixture  
with three  
components

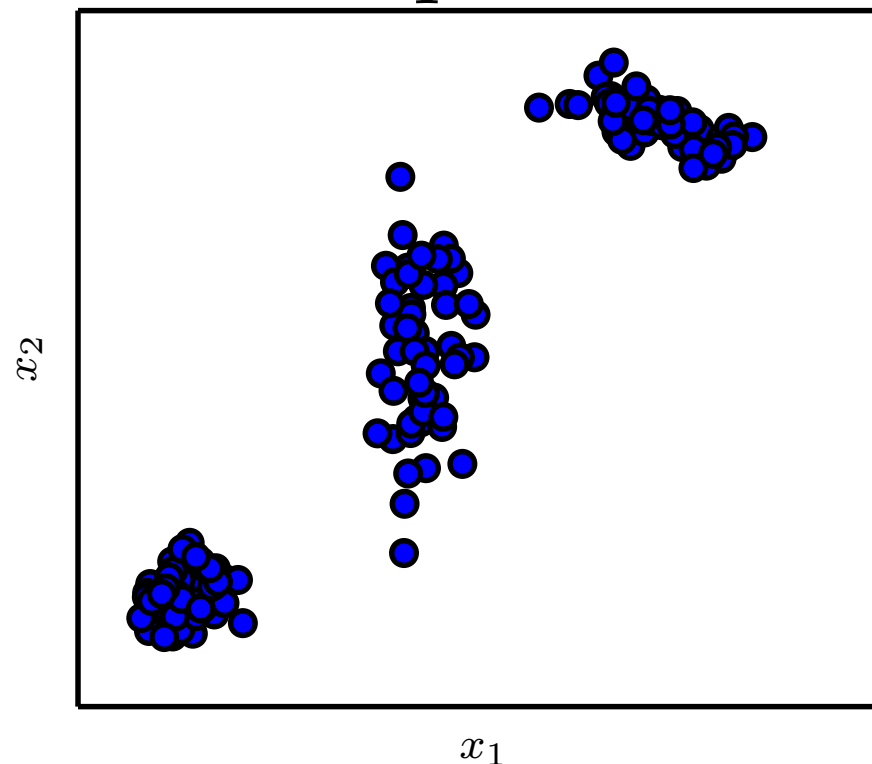


Figure 3.2

# Logistic Sigmoid

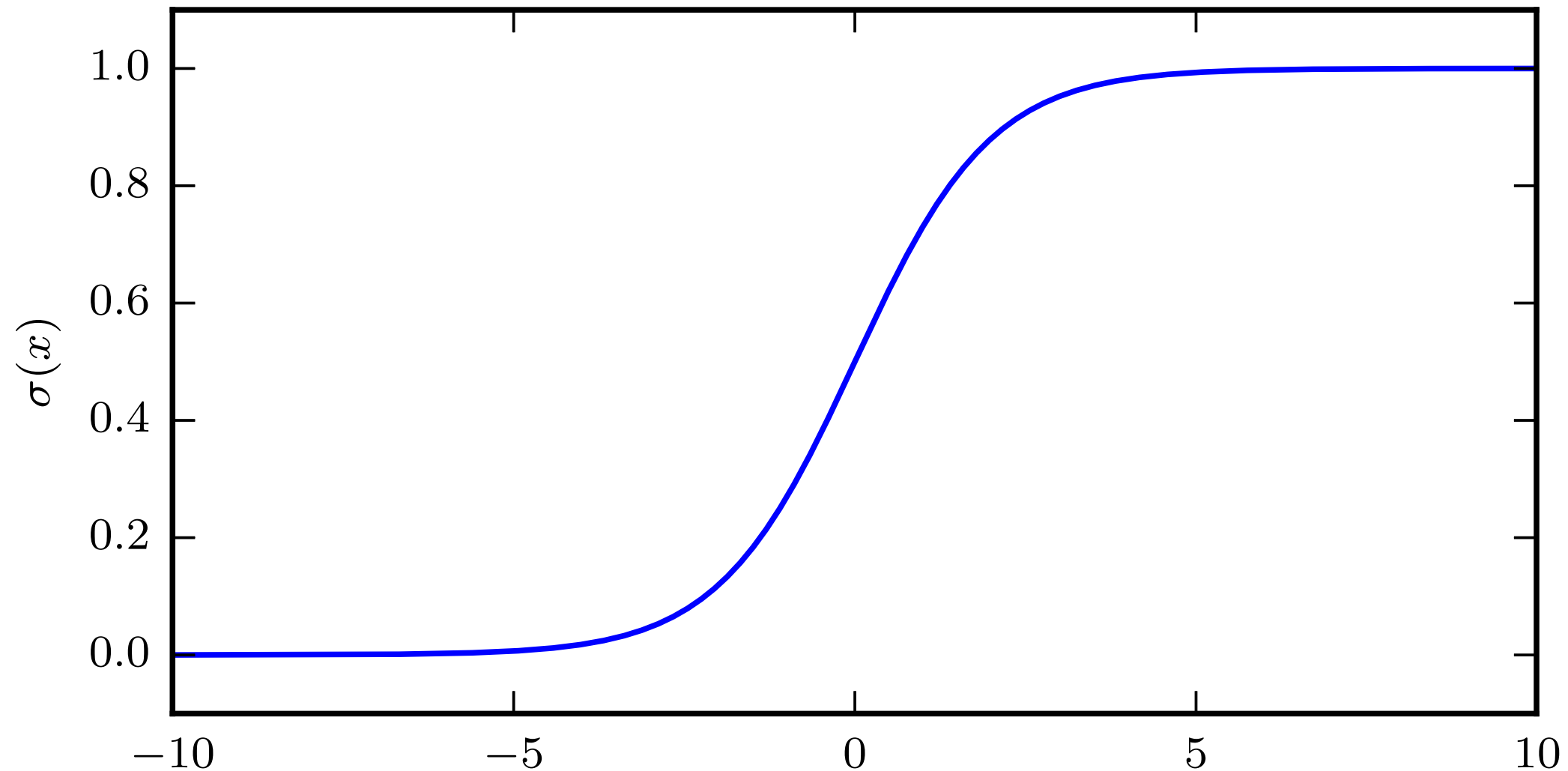


Figure 3.3: The logistic sigmoid function.

Commonly used to parametrize Bernoulli distributions

# Softplus Function

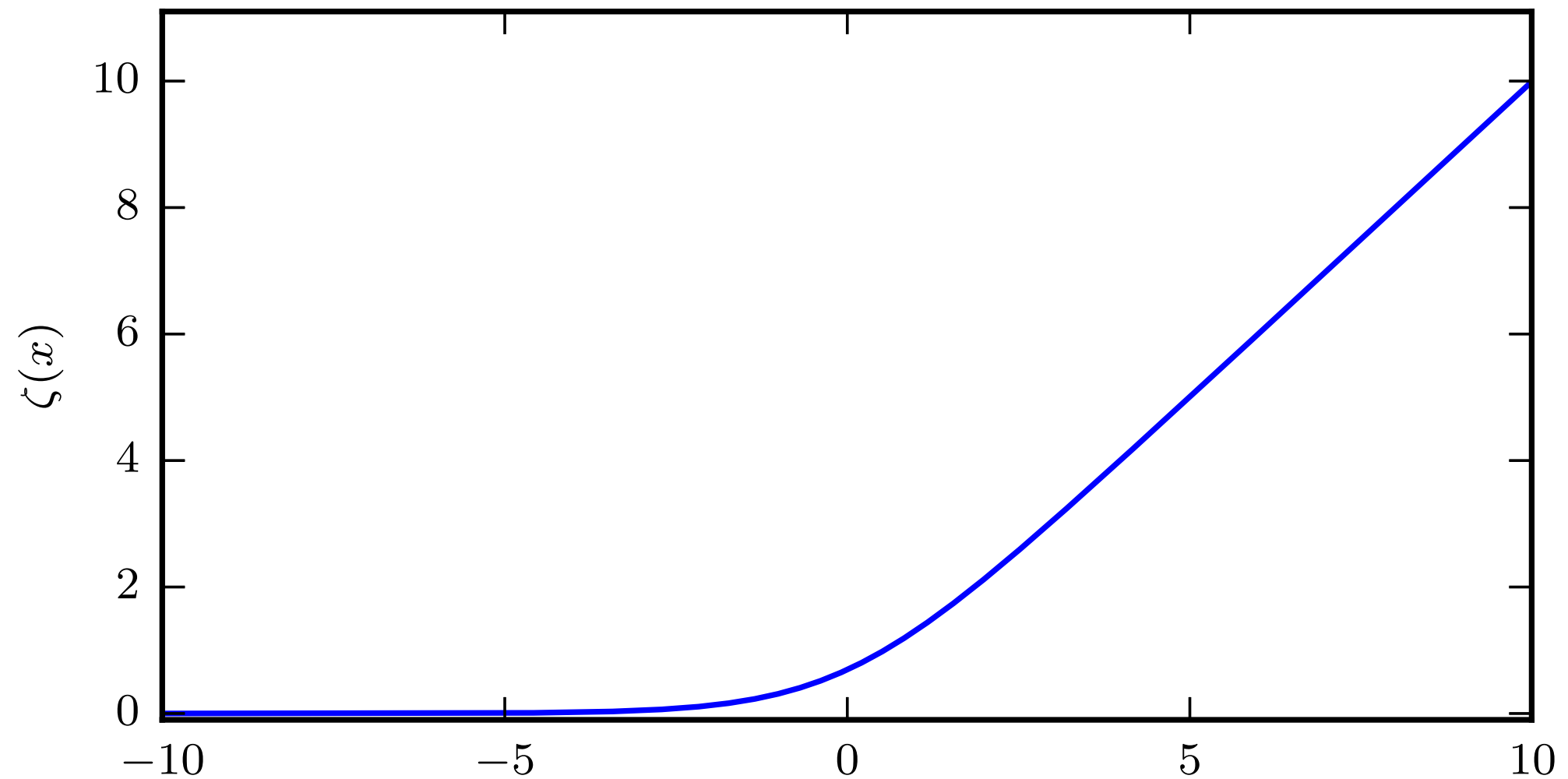


Figure 3.4: The softplus function.



# Bayes' Rule

$$P(\mathbf{x} \mid y) = \frac{P(\mathbf{x})P(y \mid \mathbf{x})}{P(y)}. \quad (3.42)$$

# Change of Variables

$$p_x(\boldsymbol{x}) = p_y(g(\boldsymbol{x})) \left| \det \left( \frac{\partial g(\boldsymbol{x})}{\partial \boldsymbol{x}} \right) \right|. \quad (3.47)$$

# Information Theory

Information:

$$I(x) = -\log P(x). \quad (3.48)$$

Entropy:

$$H(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim P}[I(x)] = -\mathbb{E}_{\mathbf{x} \sim P}[\log P(x)]. \quad (3.49)$$

KL divergence:

$$D_{\text{KL}}(P \| Q) = \mathbb{E}_{\mathbf{x} \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{\mathbf{x} \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

# Entropy of a Bernoulli Variable

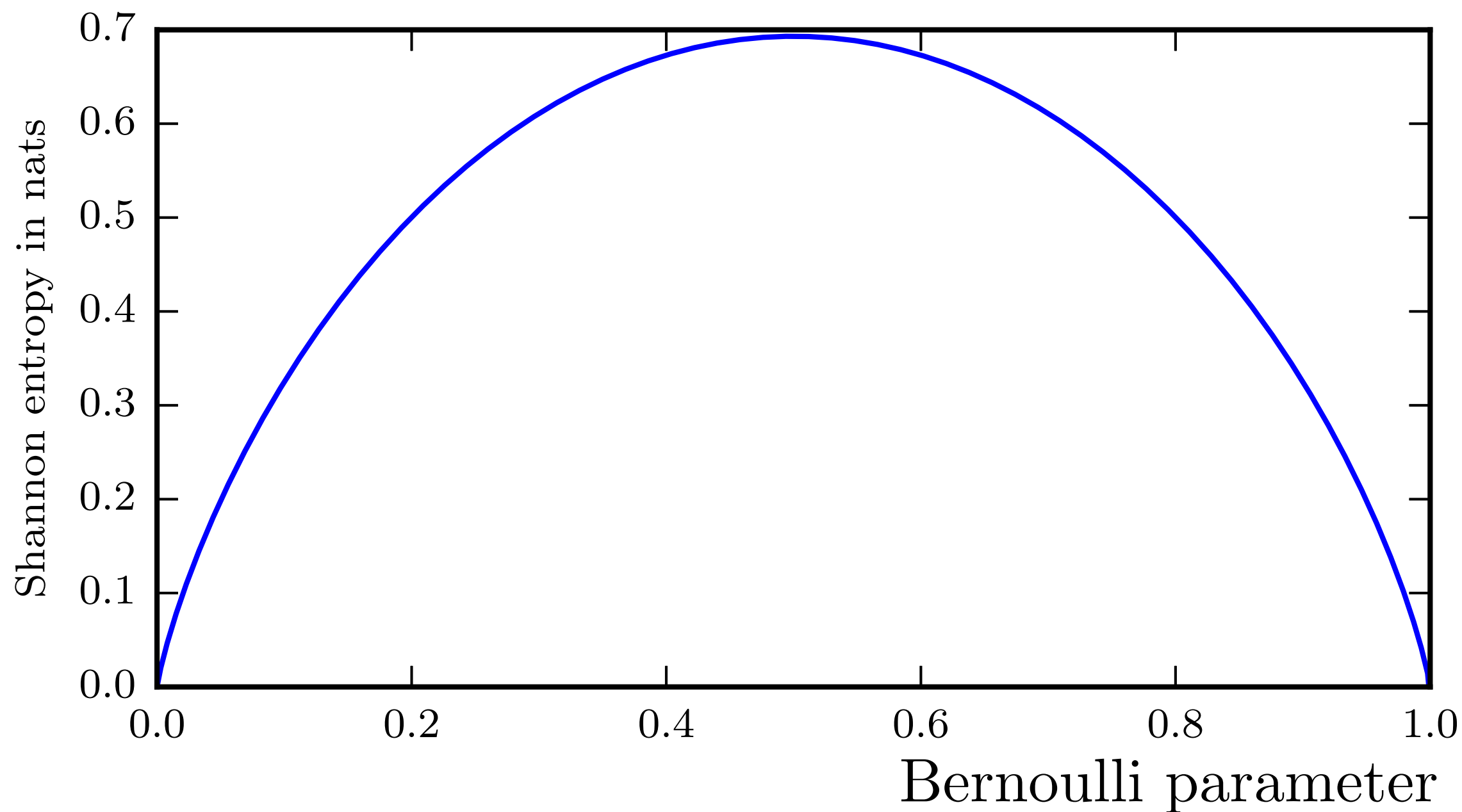
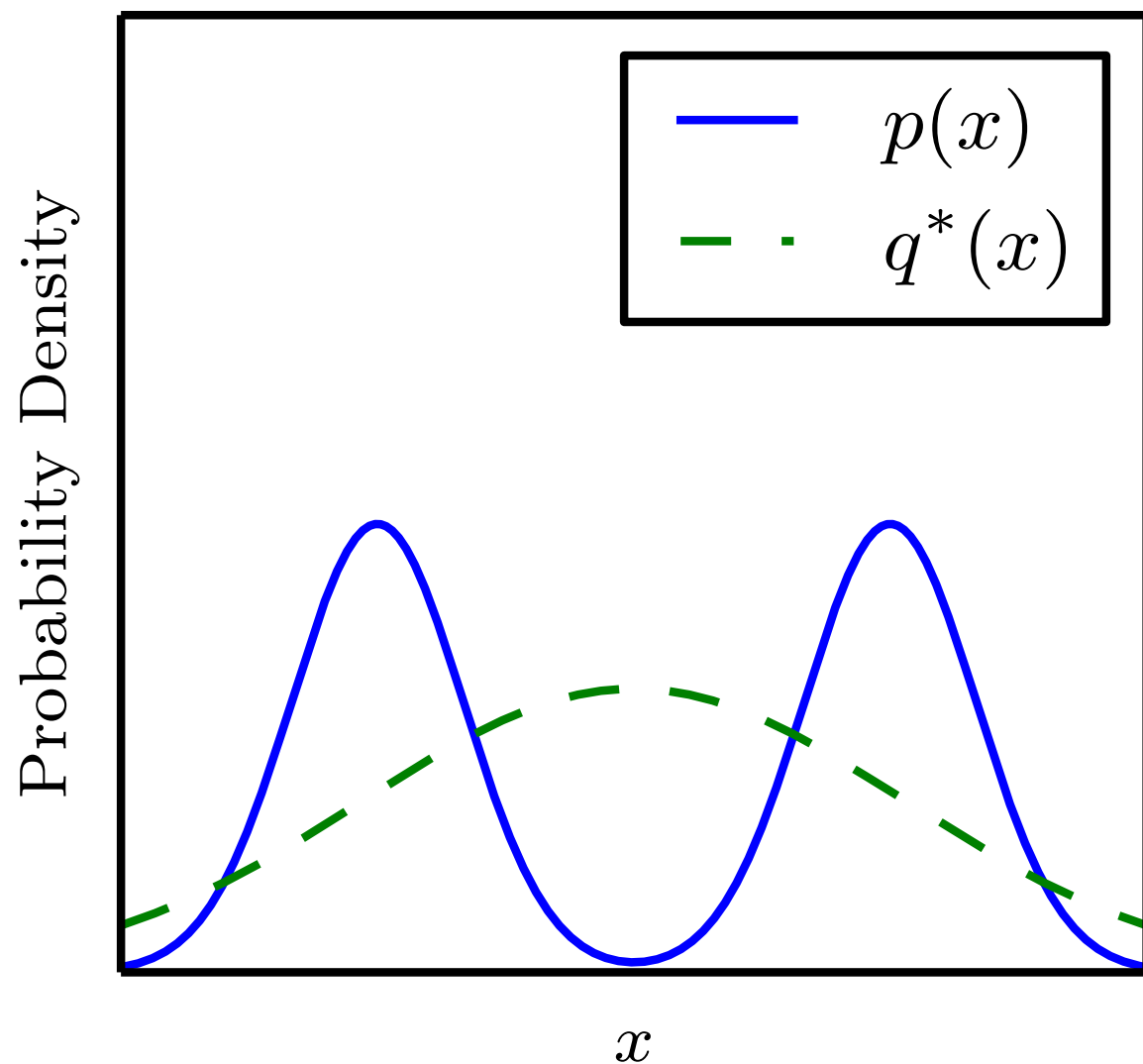


Figure 3.5

# The KL Divergence is Asymmetric

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p \| q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q \| p)$$

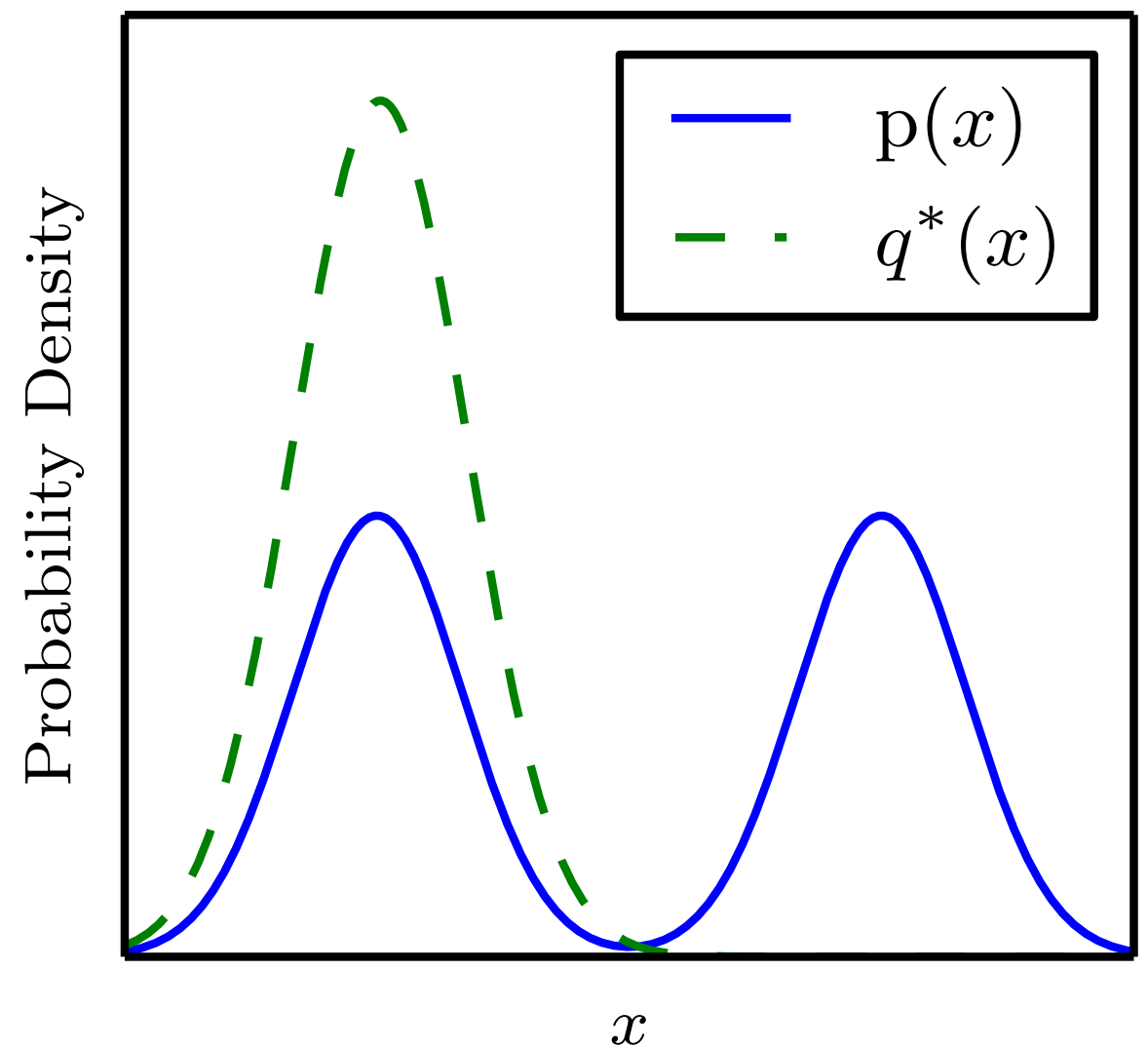
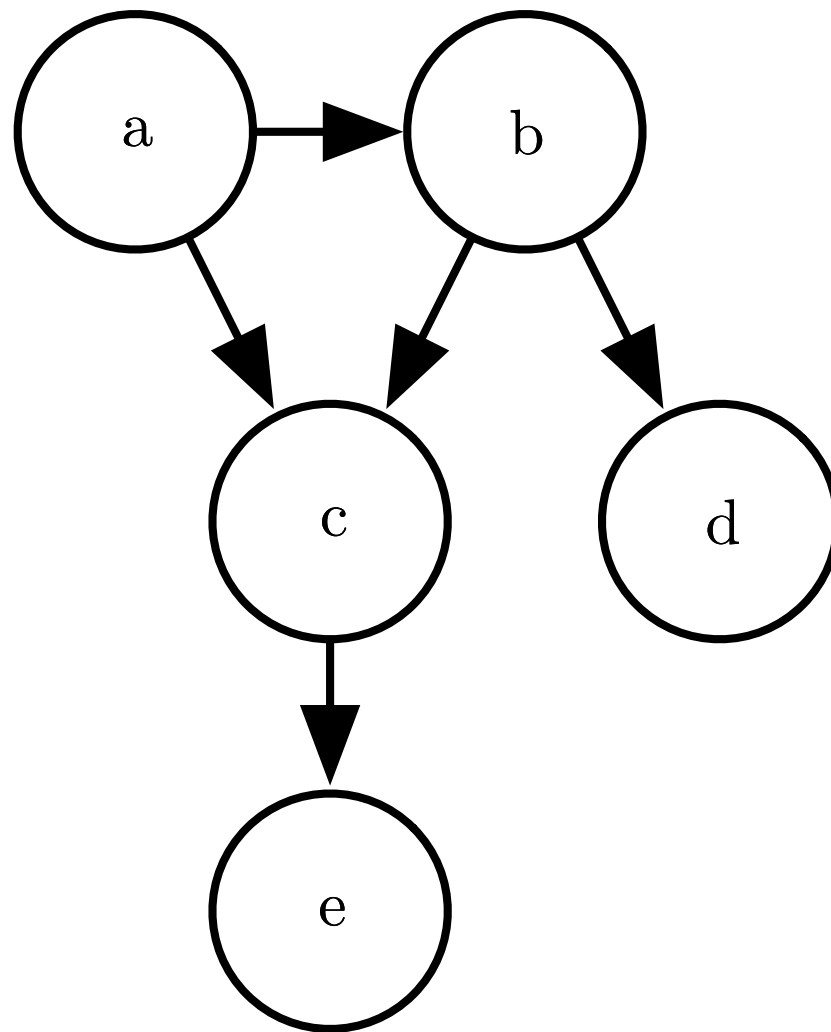


Figure 3.6

# Directed Model

Figure 3.7



$$p(a, b, c, d, e) = p(a)p(b \mid a)p(c \mid a, b)p(d \mid b)p(e \mid c). \quad (3.54)$$

# Undirected Model

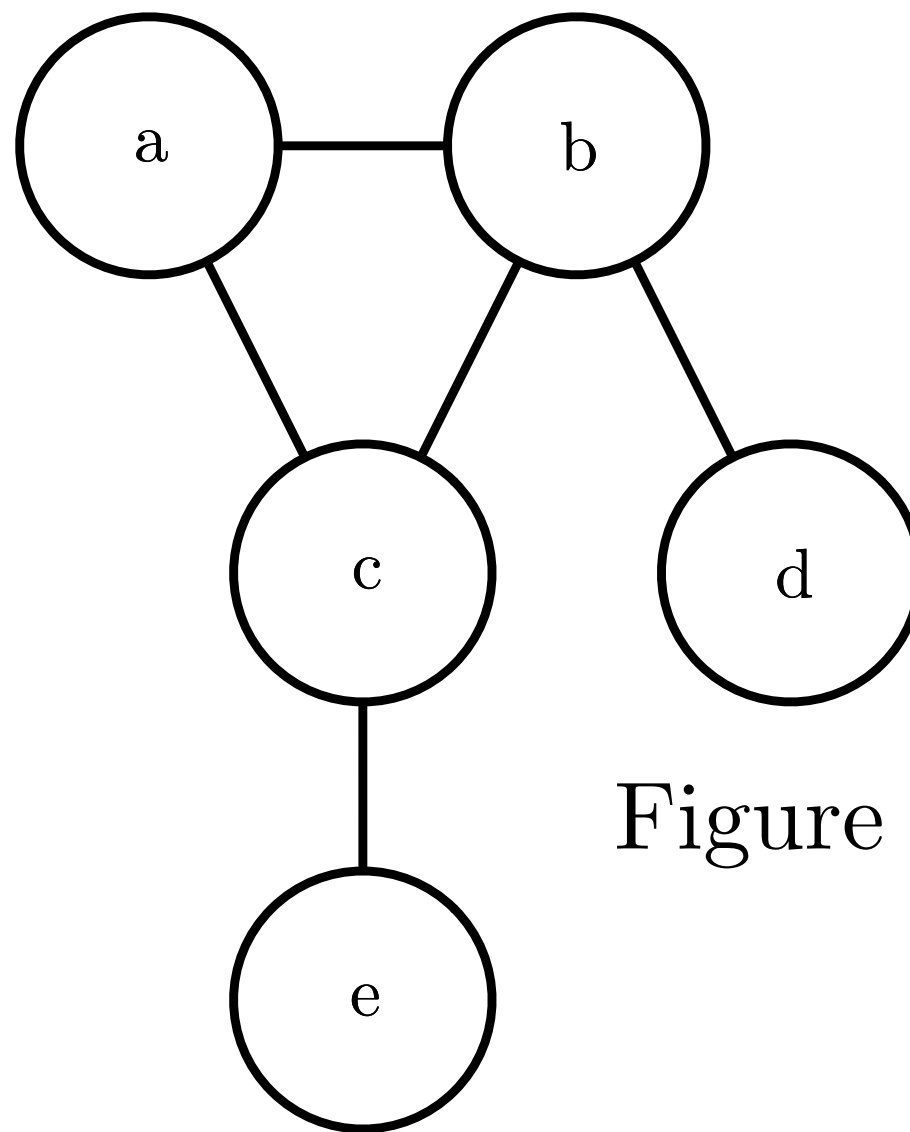


Figure 3.8

$$p(a, b, c, d, e) = \frac{1}{Z} \phi^{(1)}(a, b, c) \phi^{(2)}(b, d) \phi^{(3)}(c, e). \quad (3.56)$$