

# Maths fundamentals for AI and deep learning

# Linear Algebra

Scope: Focused on the subset of linear algebra most relevant to deep learning.

# Scalars

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

$a, n, x$

# Vectors

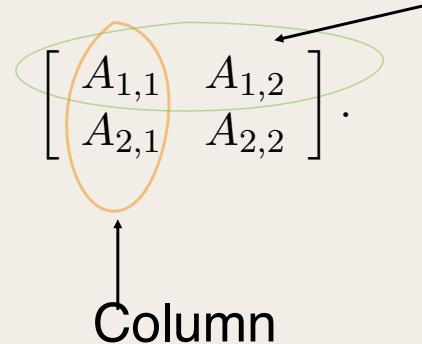
- A vector is a 1-D array of numbers:

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

- 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:

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


- Example notation for type and shape:  $A \in \mathbb{R}^{m \times n}$

# 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

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

The diagram shows a 3x2 matrix  $A$  with elements  $A_{1,1}, A_{1,2}, A_{2,1}, A_{2,2}, A_{3,1}, A_{3,2}$ . A curved arrow points from the element  $A_{1,1}$  to its transpose position  $A_{1,1}$  in the resulting 2x3 matrix  $\mathbf{A}^\top$ . The resulting matrix  $\mathbf{A}^\top$  is shown with elements  $A_{1,1}, A_{2,1}, A_{3,1}, A_{1,2}, A_{2,2}, A_{3,2}$ .

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

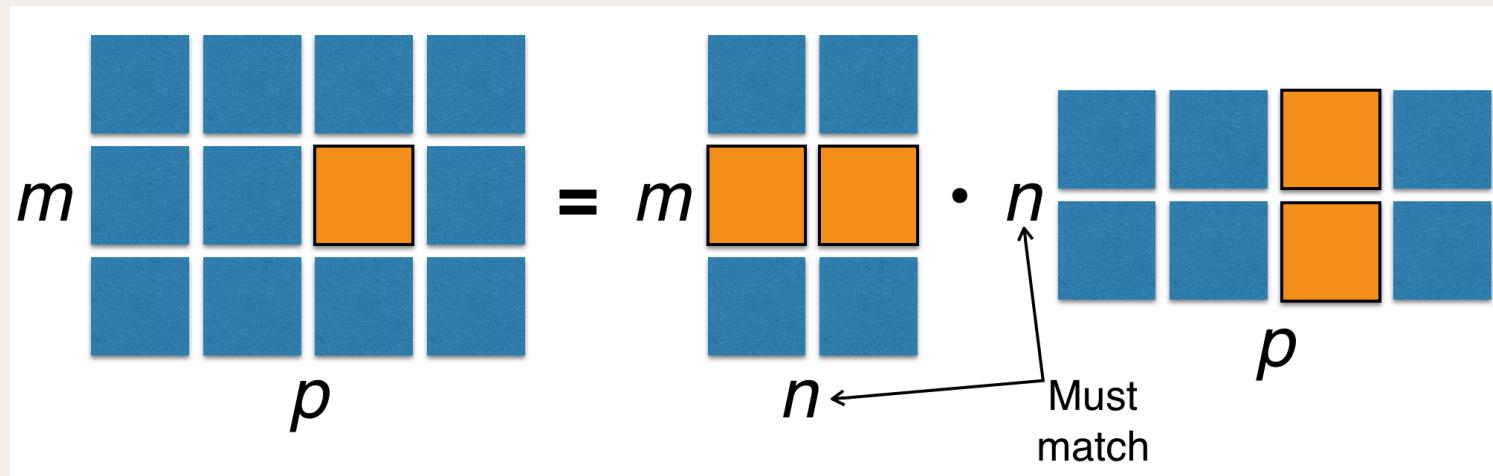
Figure 2.1: The transpose of the matrix can be thought of as a mirror image across the main diagonal.

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

## Matrix (Dot) Product

$$C = AB. \quad (2.4)$$

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



(Goodfellow 2016)

## Identity Matrix

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

Figure 2.2: *Example identity matrix:* This is  $I_3$ .

$$\forall \mathbf{x} \in \mathbb{R}^n, I_n \mathbf{x} = \mathbf{x}. \quad (2.20)$$

## Systems of Equations

$$\mathbf{A}\mathbf{x} = \mathbf{b} \quad (2.11)$$

expands to

$$\mathbf{A}_{1,:}\mathbf{x} = b_1 \quad (2.12)$$

$$\mathbf{A}_{2,:}\mathbf{x} = b_2 \quad (2.13)$$

$$\dots \quad (2.14)$$

$$\mathbf{A}_{m,:}\mathbf{x} = b_m \quad (2.15)$$

# 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

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

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

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

- Max norm, infinite  $p$ :

$$\|x\|_\infty = \max_i |x_i|. \quad (2.32)$$

# Special Matrices and Vectors

- Unit vector:

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

- Symmetric Matrix:

$$A = A^\top. \quad (2.35)$$

- Orthogonal matrix:  $A^\top A = AA^\top = I.$  (2.37)

$$A^{-1} = A^\top$$

# 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

# Probability Mass Function

- The domain of  $P$  must be the set of all possible states of  $x$ .
- $\forall x \in \mathcal{X}, 0 \leq P(x) \leq 1$ . 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.
- $\sum_{x \in \mathcal{X}} P(x) = 1$ . 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:

$$P(\mathbf{x} = \mathbf{x}_i) = \frac{1}{k}$$

# Probability Density Function

- The domain of  $p$  must be the set of all possible states of  $x$ .
- $\forall x \in \mathcal{X}, p(x) \geq 0$ . 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

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

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

# Conditional Probability

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

# Chain Rule of Probability

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

# Independence

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

# Conditional Independence

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

# Expectation

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

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

linearity of expectations:

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

# 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:

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

## Bernoulli Distribution

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

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

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

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

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

# 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

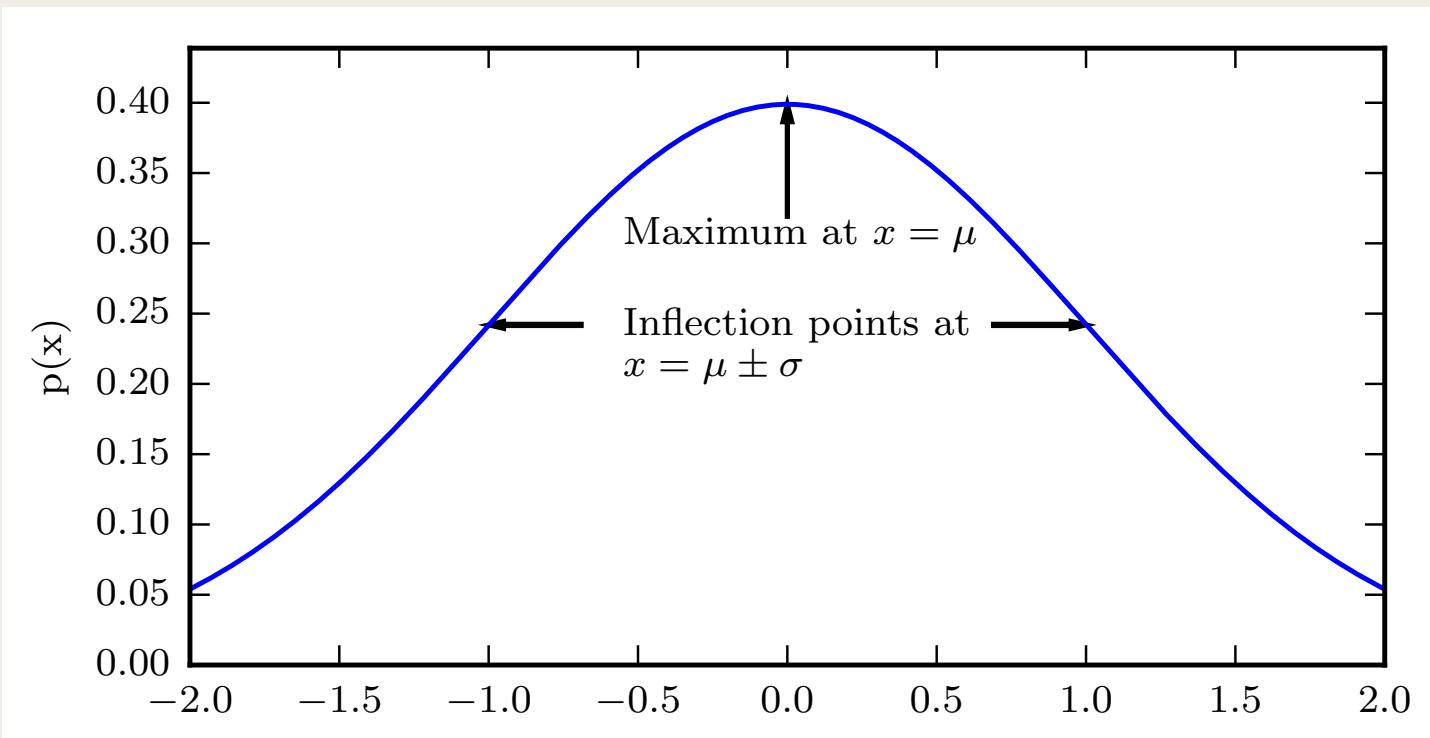


Figure 3.1

(Goodfellow 2016)

# 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:

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

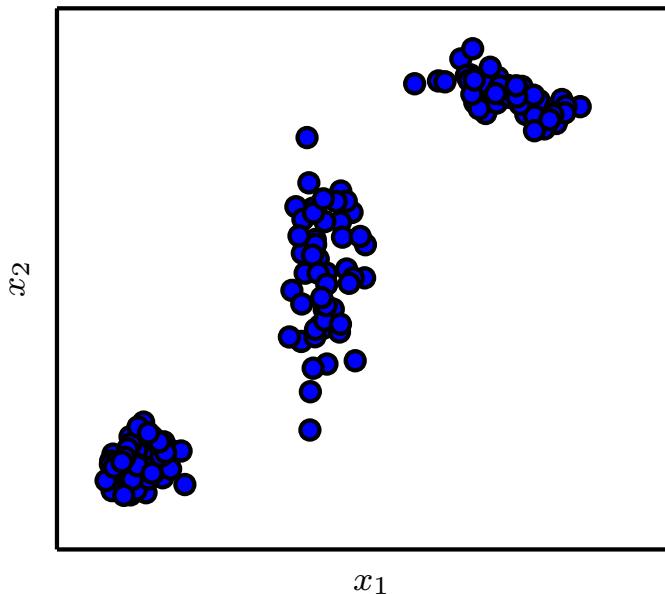
# Empirical Distribution

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

# Mixture Distributions

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

Gaussian  
mixture with  
three  
components



(Goodfellow 2016)

# Logistic Sigmoid

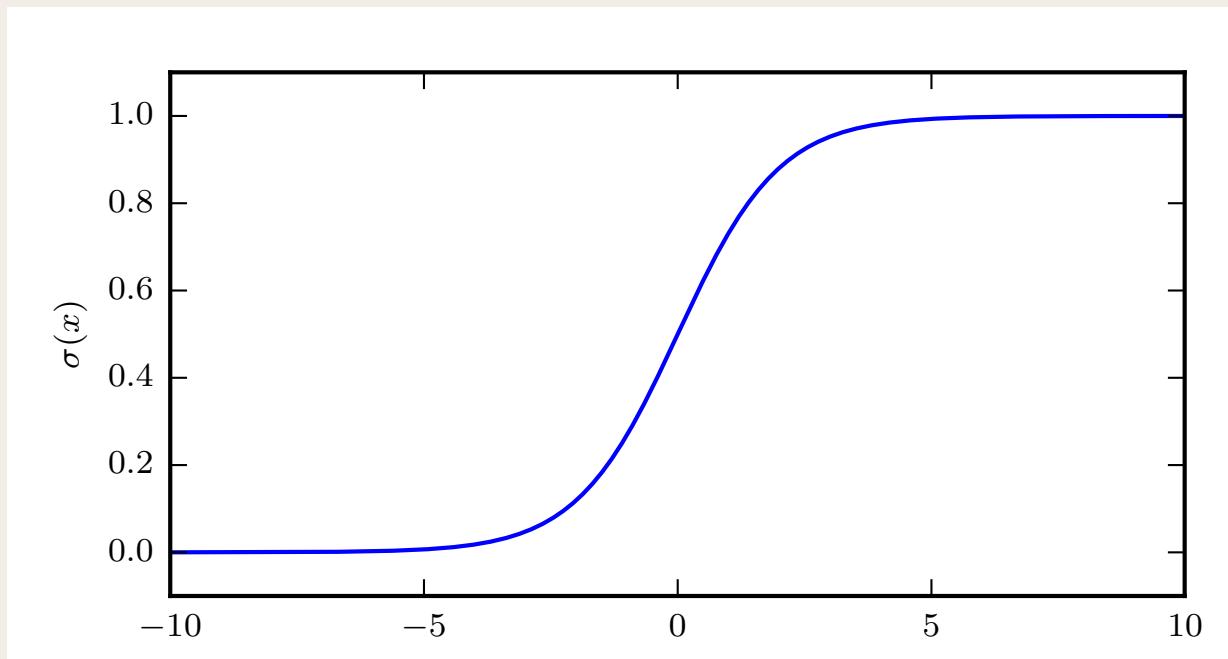


Figure 3.3: The logistic sigmoid function.

Commonly used to parametrize Bernoulli distributions

(Goodfellow 2016)

# Softplus Function

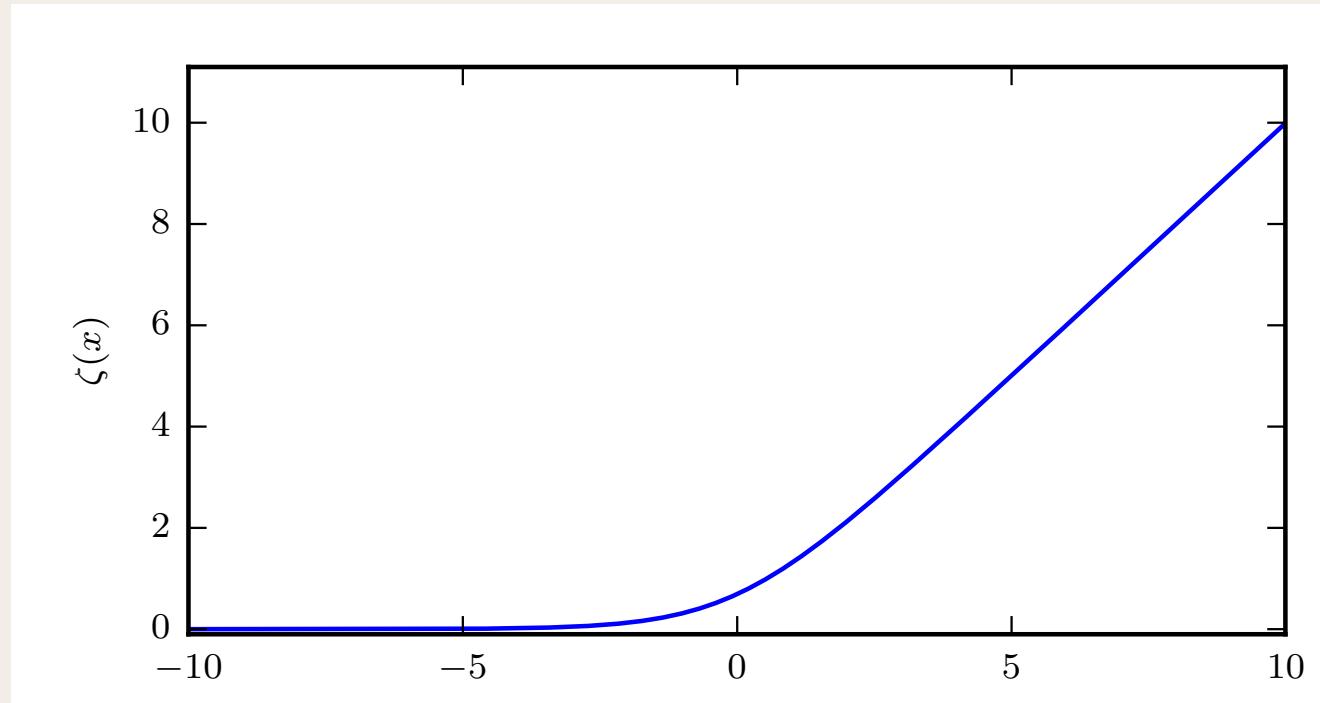


Figure 3.4: The softplus function.

(Goodfellow 2016)

## Bayes' Rule

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

# Information Theory

Information:

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

Entropy:

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

KL divergence:

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

# The KL Divergence is Asymmetric

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

Probability Density

$p(x)$   
 $q^*(x)$

$x$

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

Probability Density

$p(x)$   
 $q^*(x)$

$x$

Figure 3.6

(Goodfellow 2016)

# Numerical Computation

# Overflow and Underflow

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}. \quad (4.1)$$

The exponentiation can underflow when the argument is very negative, or overflow when it is very positive.

# Gradient Descent

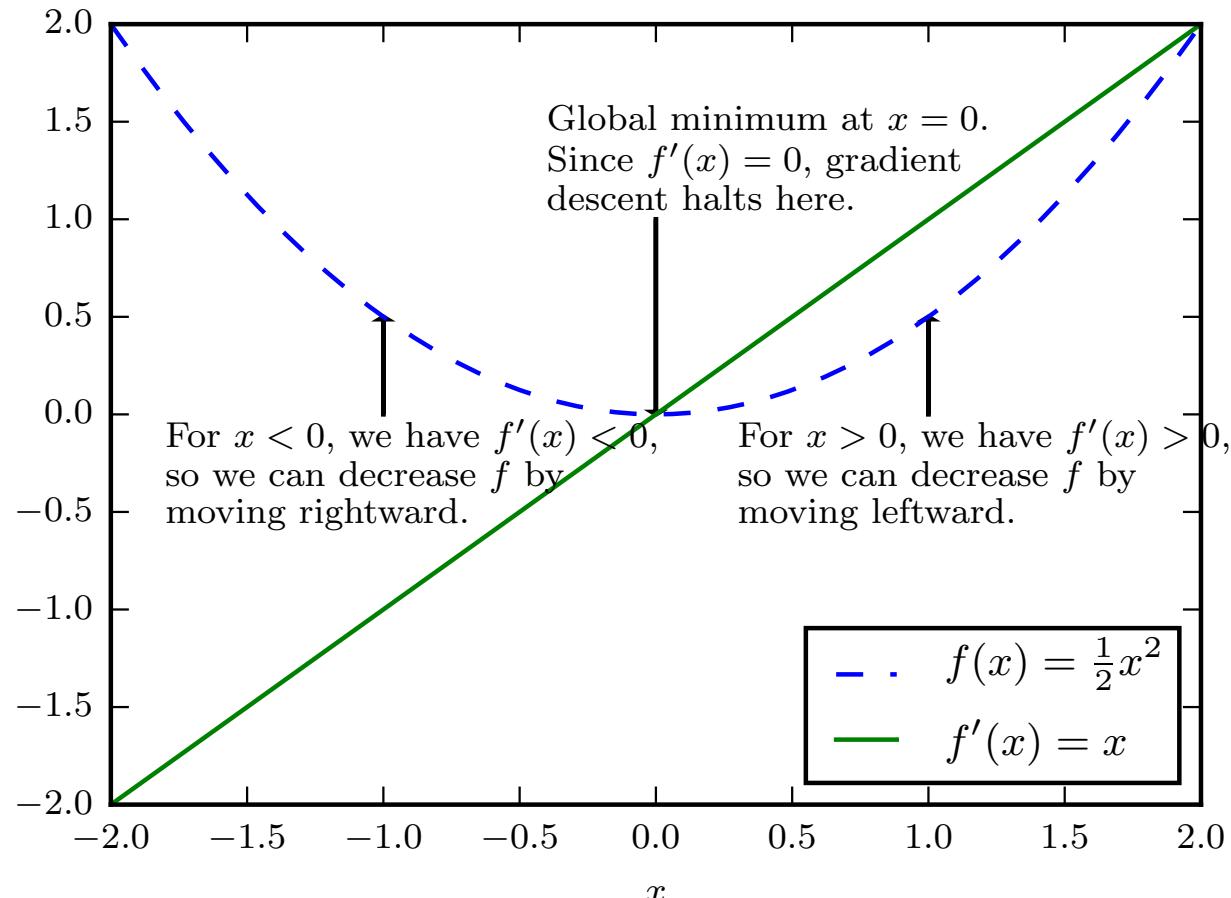


Figure 4.1

(Goodfellow 2016)

# Approximate Optimization

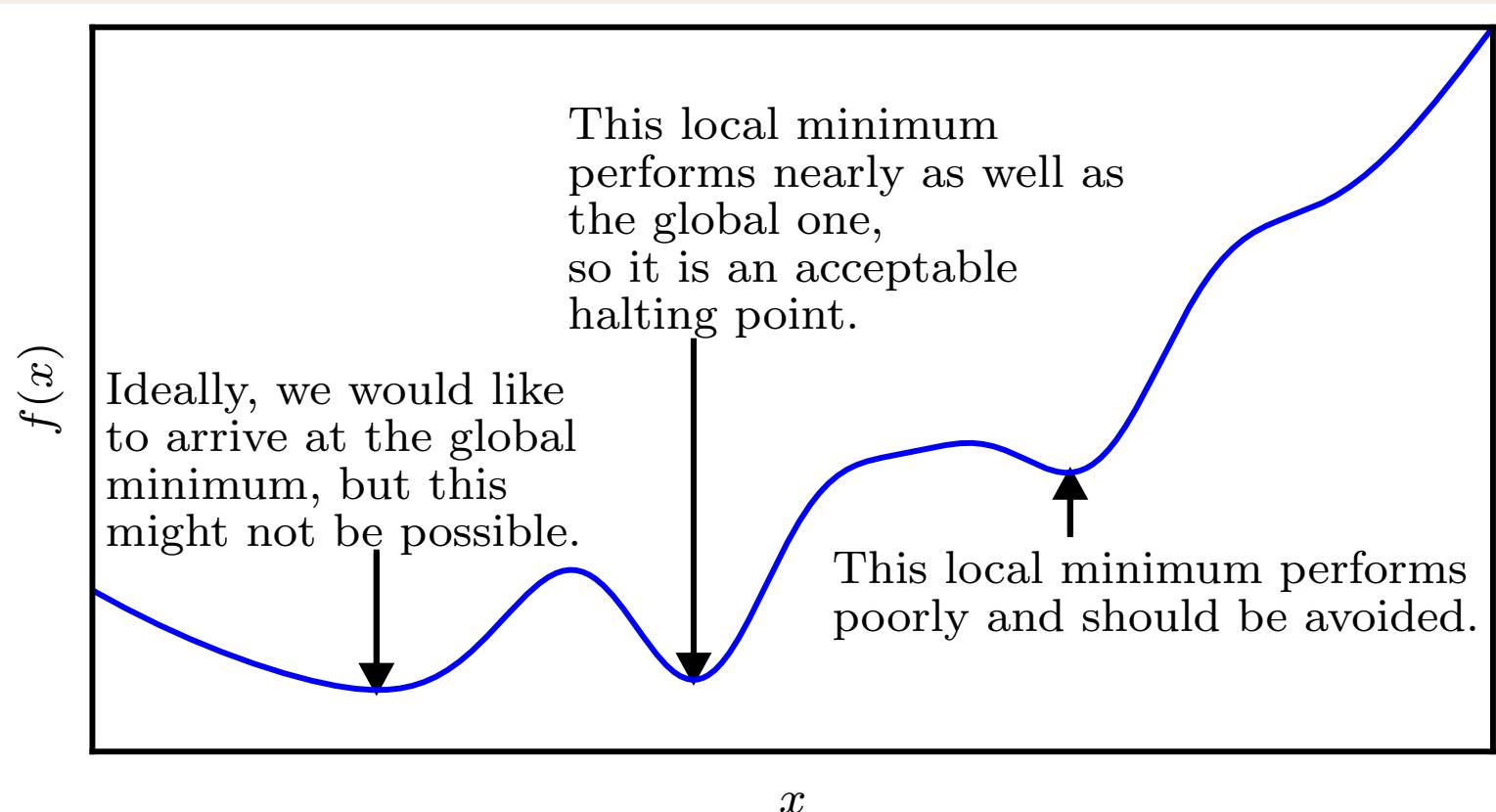


Figure 4.3

(Goodfellow 2016)

# Saddle Points

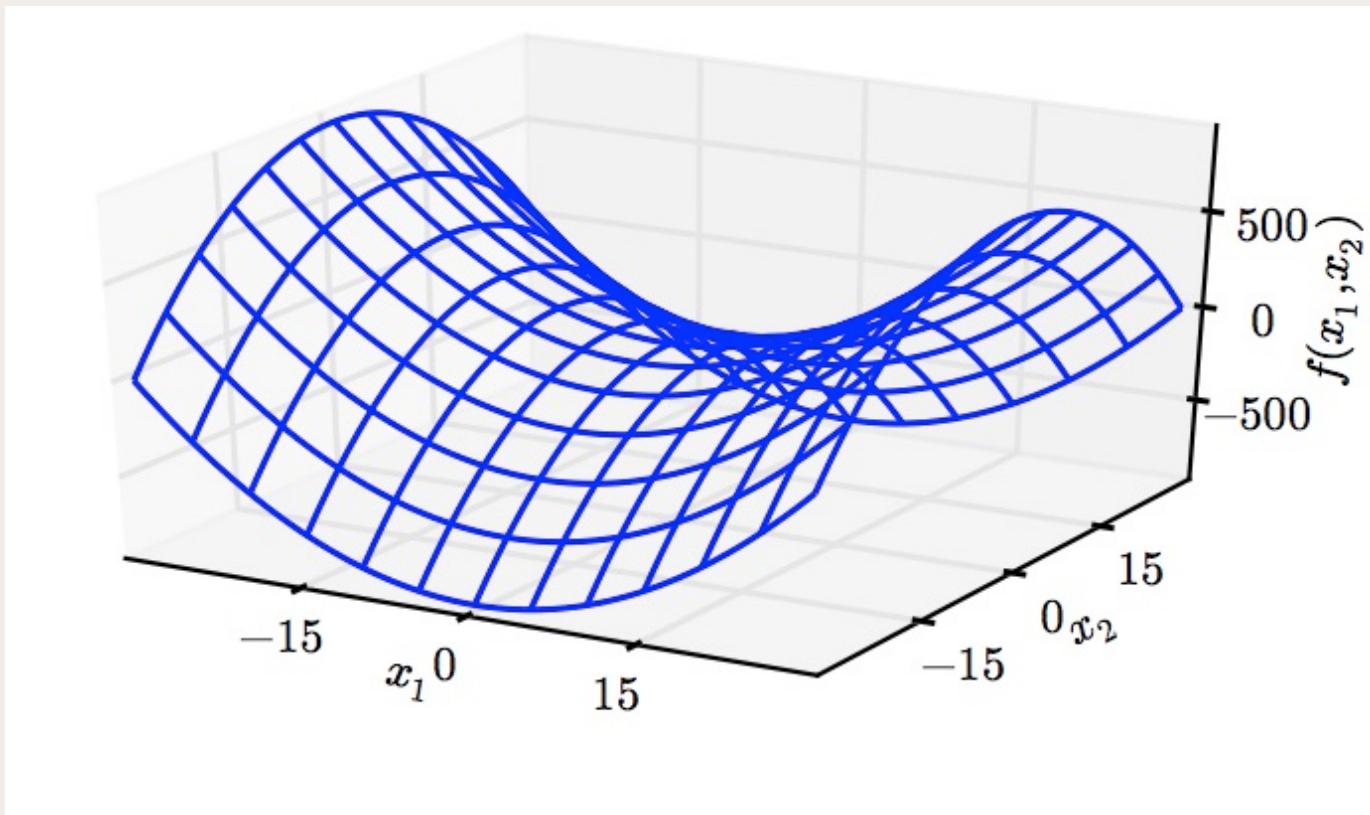


Figure 4.5 (Goodfellow 2016)

## Gradient Descent and Poor Conditioning

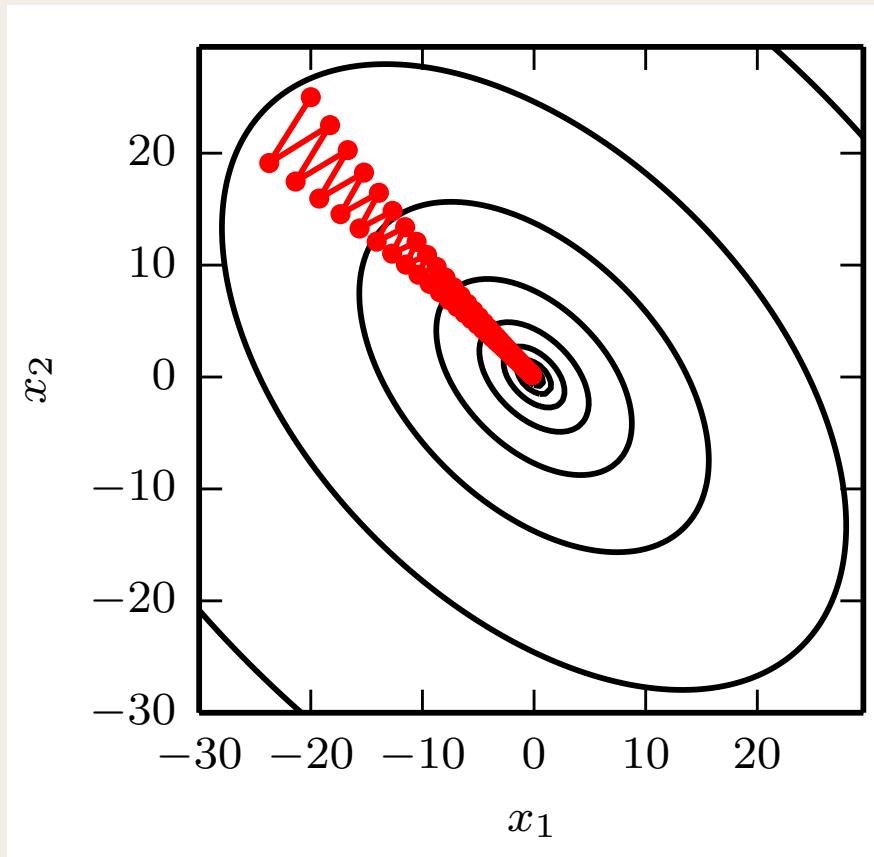


Figure 4.6  
(Goodfellow 2016)

## Backpropagation: a simple example

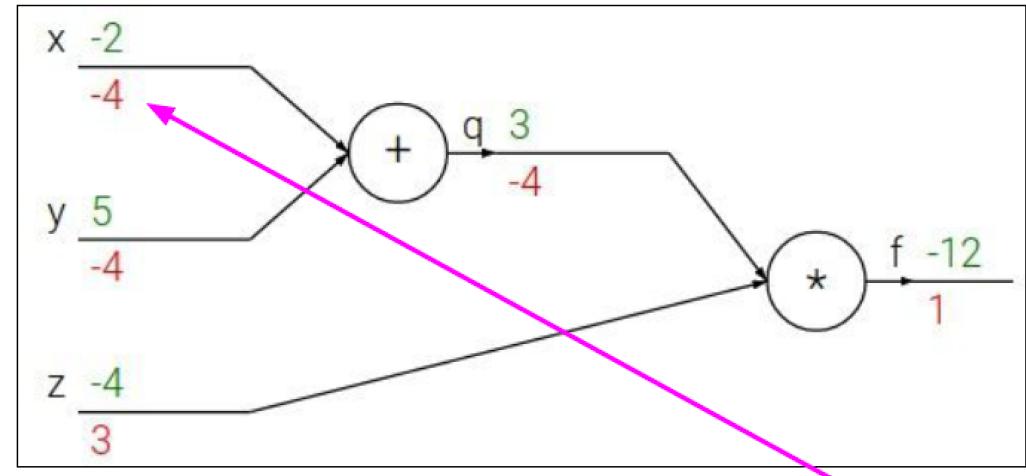
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream  
gradient

Local  
gradient

# References

- *Deep Learning*, [www.deeplearningbook.org](http://www.deeplearningbook.org), Ian Goodfellow et. al.
- *Lecture slides for Chapter 2, 3 , 4 of Deep Learning*, Ian Goodfellow, 2016-06-24
- *Linear Algebra*, Georgi Shilov, Publisher : Dover Publications Inc.; New edition (1 Mar. 1978), ISBN-10 : 048663518X, ISBN-13 : 978-0486635187
- *Stanford CS class* [CS231n: Convolutional Neural Networks for Visual Recognition](https://cs231n.github.io/)