# Applied Math for Deep Learning

Prof. Kuan-Ting Lai
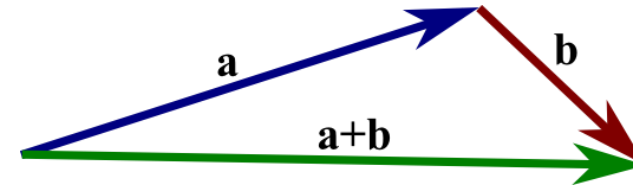
2020/3/10

# Applied Math for Deep Learning

- Linear Algebra
- Probability
- Calculus
- Optimization
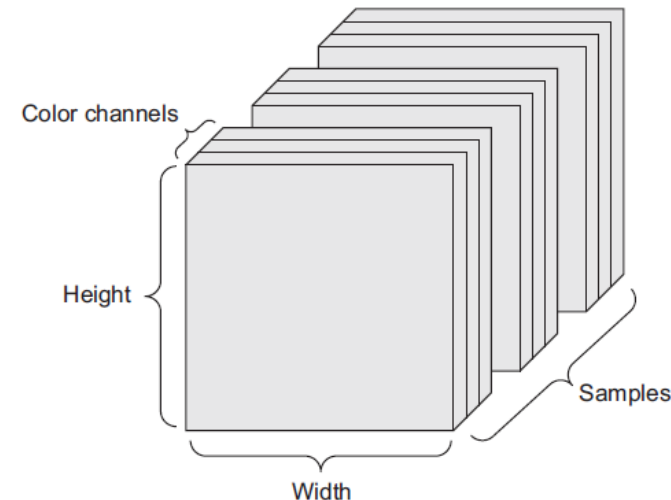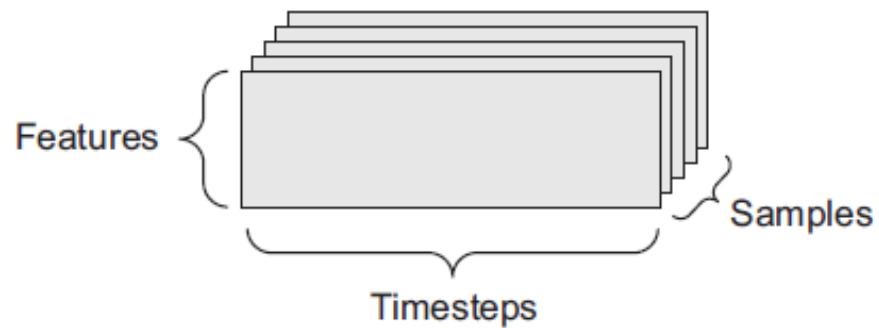
# Linear Algebra

- Scalar
  - real numbers

- Vector (1D)
  - Has a magnitude & a direction

- Matrix (2D)
  - An array of numbers arranges in rows & columns

- Tensor (>=3D)
  - Multi-dimensional arrays of numbers

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

# Real-world examples of Data Tensors

- Timeseries Data – 3D (samples, timesteps, features)
- Images – 4D (samples, height, width, channels)
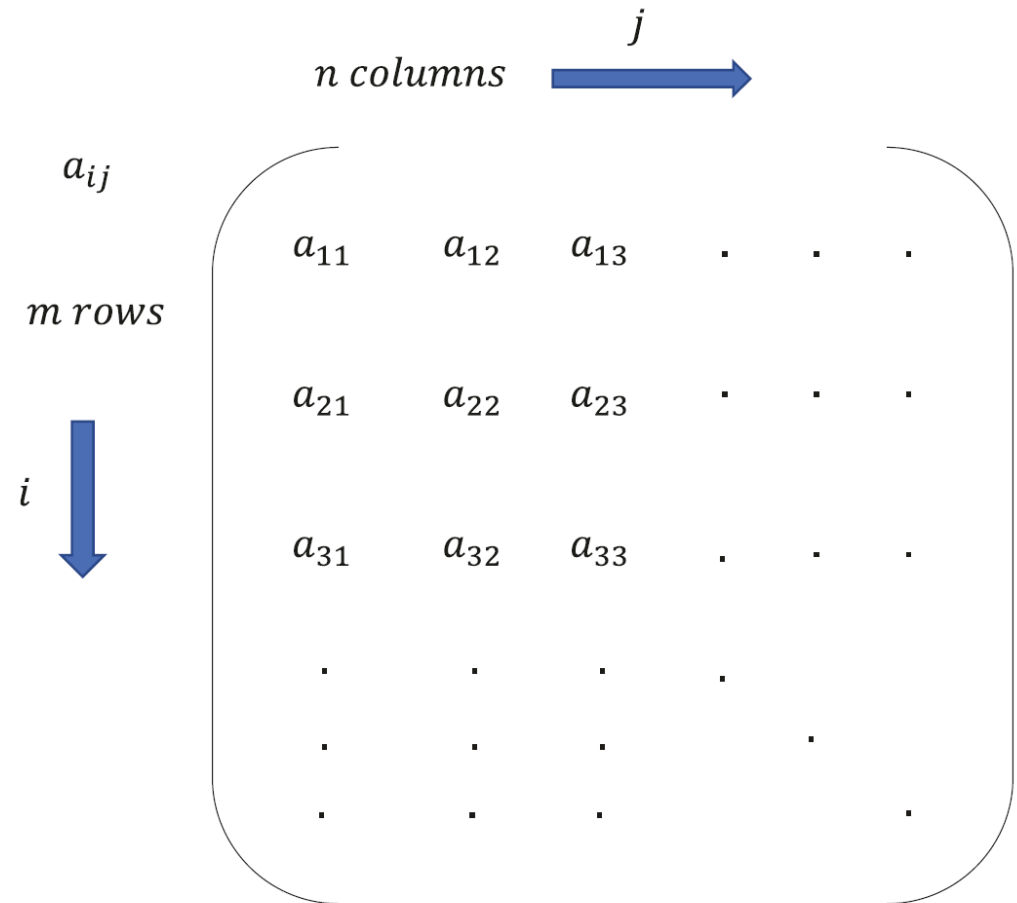- Video – 5D (samples, frames, height, width, channels)

# The Matrix

# Matrix

- Define a matrix with m rows and n columns:

$$A_{m \times n} \in \mathbb{R}^{m \times n}$$



$n\ columns \xrightarrow{j}$

$a_{ij}$

$m\ rows$

$i$

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & . & . & . \\ a_{21} & a_{22} & a_{23} & . & . & . \\ a_{31} & a_{32} & a_{33} & . & . & . \\ . & . & . & . & & \\ . & . & . & . & & . \\ . & . & . & & & . \end{pmatrix}$$

Santanu Pattanayak, "Pro Deep Learning with TensorFlow," Apress, 2017

# Matrix Operations

- Addition and Subtraction

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

$$A + B = \begin{bmatrix} 1+5 & 2+6 \\ 3+7 & 4+8 \end{bmatrix} = \begin{bmatrix} 6 & 8 \\ 10 & 12 \end{bmatrix}$$
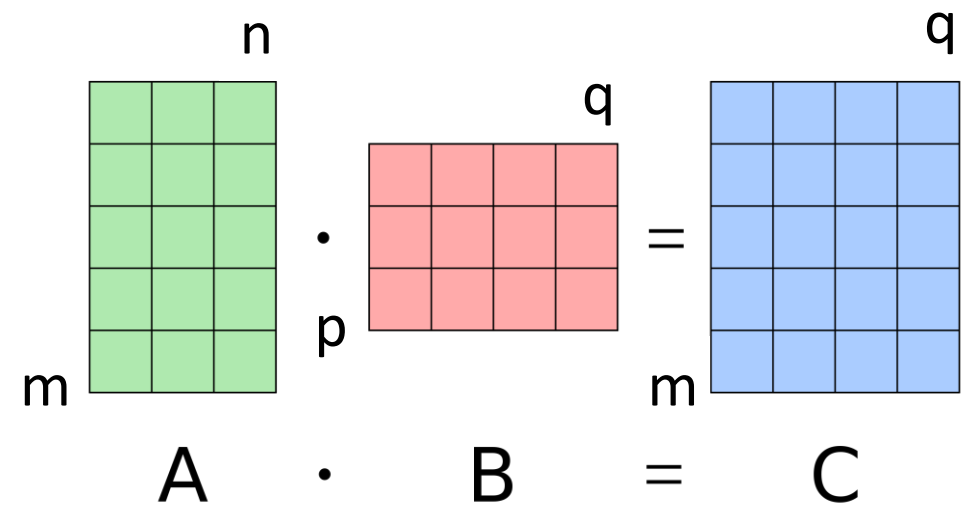
$$A - B = \begin{bmatrix} 1-5 & 2-6 \\ 3-7 & 4-8 \end{bmatrix} = \begin{bmatrix} -4 & -4 \\ -4 & -4 \end{bmatrix}$$

# Matrix Multiplication

- Two matrices A and B, where $A \in \mathbb{R}^{m \times n}$  $B \in \mathbb{R}^{p \times q}$

- The columns of A must be equal to the rows of B, i.e.  n == p

- A * B = C, where  $C \in \mathbb{R}^{m \times q}$

- $c_{ij} = \sum_{k=1}^{n} a_{ik} b_{kj}$

# Example of Matrix Multiplication (3-1)

"Dot Product"

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 & \end{bmatrix}$$

# Example of Matrix Multiplication (3-2)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 & 64 \\ & \end{bmatrix}$$

# Example of Matrix Multiplication (3-3)

$$
\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 & 64 \\ 139 & 154 \end{bmatrix} \checkmark
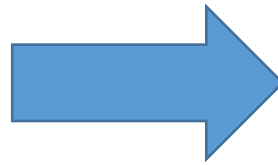$$

# Matrix Transpose

$$A \in \mathbb{R}^{m \times n} \qquad A^{\mathrm{T}} \in \mathbb{R}^{n \times m}$$

$$a'_{ji} = a_{ij} \quad \forall\, i \in \{1, 2, ..m\}, \forall\, j \in \{1, 2, ..n\}$$

A

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

A

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

$A^{\mathrm{T}}$

$$\begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}$$

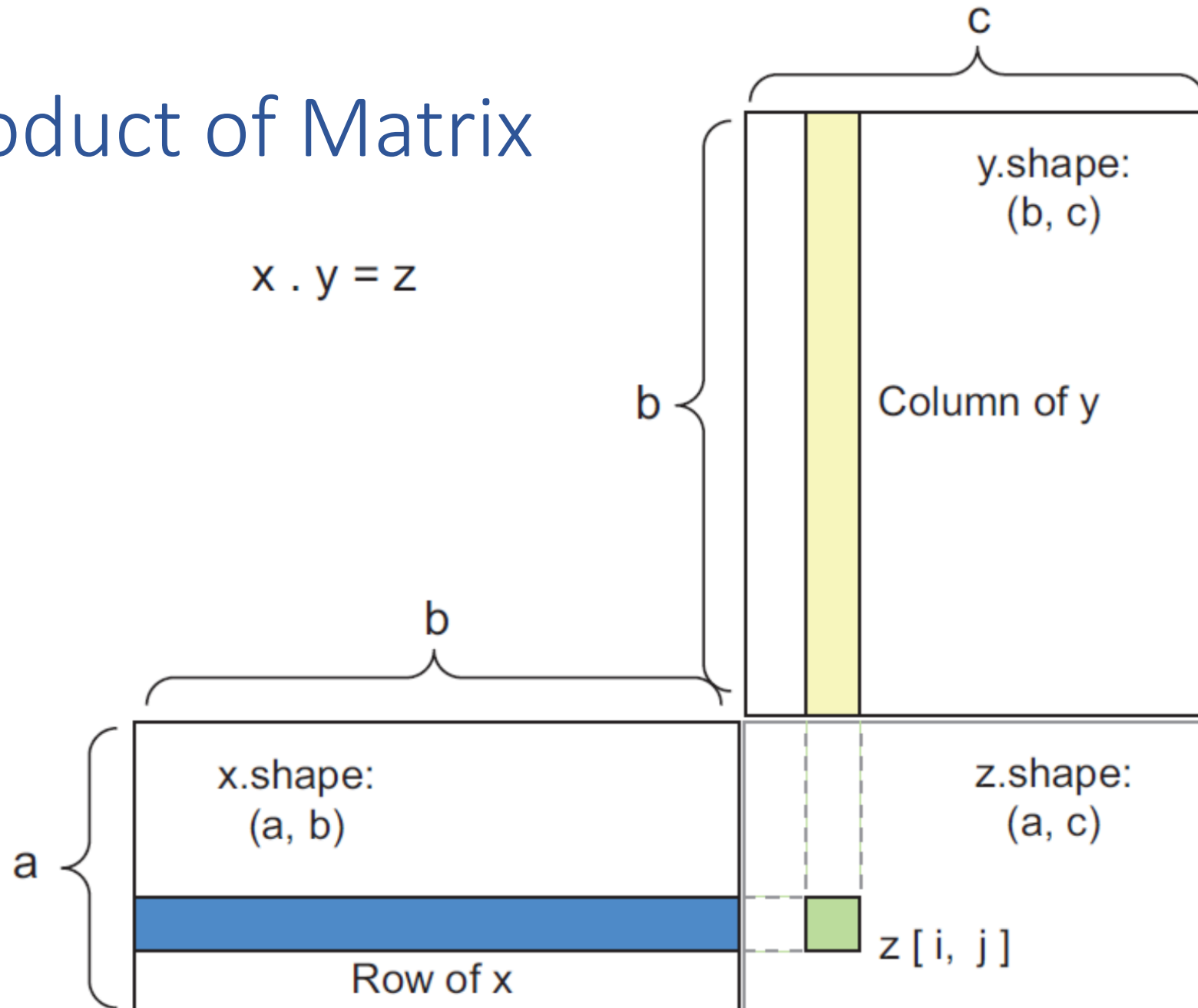https://en.wikipedia.org/wiki/Transpose

# Dot Product

- Dot product of two vectors become a **scalar**
- Notation: $v_1 \cdot v_2$ or $v_1{}^T v_2$

$$v_1 = \begin{bmatrix} v_{11} \\ v_{12} \\ \cdot \\ \cdot \\ \cdot \\ v_{1n} \end{bmatrix} \quad v_2 = \begin{bmatrix} v_{21} \\ v_{22} \\ \cdot \\ \cdot \\ \cdot \\ v_{2n} \end{bmatrix}$$

$$v_1 \cdot v_2 = v_1{}^T v_2 = v_2{}^T v_1 = v_{11}v_{21} + v_{12}v_{22} + .. + v_{1n}v_{2n} = \sum_{k=1}^{n} v_{1k}v_{2k}$$

# Dot Product of Matrix

$$x \cdot y = z$$



c

y.shape:
(b, c)

b

Column of y

b

a

x.shape:
(a, b)

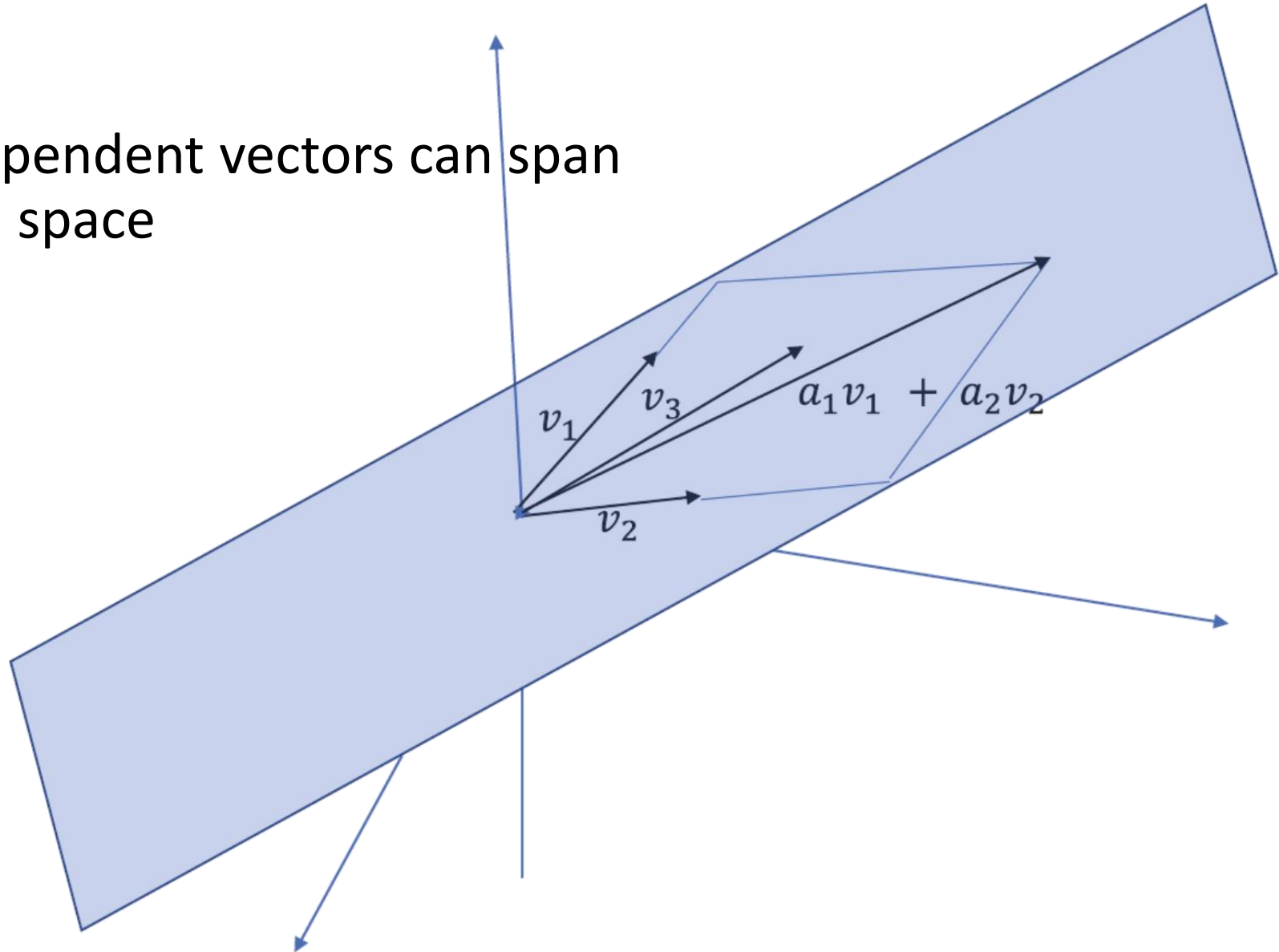z.shape:
(a, c)

z [ i,  j ]

Row of x

# Linear Independence

- A vector is **linearly dependent** on other vectors if it can be expressed as the linear combination of other vectors

- A set of vectors $v_1, v_2, \cdots, v_n$ is **linearly independent** if $a_1 v_1 + a_2 v_2 + \cdots + a_n v_n = 0$ implies all $a_i = 0, \forall i \in \{1,2, \cdots n\}$

$$\begin{bmatrix} v_1 v_2 \dots v_n \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ . \\ . \\ a_n \end{bmatrix} = 0 \; where \; v_i \in \mathbb{R}^{m \times 1} \; \forall i \in \{1,2,\dots,n\}, \; \begin{bmatrix} a_1 \\ a_2 \\ . \\ . \\ a_n \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

# Span the Vector Space

- *n* linearly independent vectors can span *n*-dimensional space

# Rank of a Matrix

- Rank is:
  - The number of linearly independent row or column vectors
  - The dimension of the vector space generated by its columns

- Row rank = Column rank

- Example:

$$A = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$$

Row-echelon form $\longrightarrow$

$$\begin{bmatrix} 1 & 0 & -5 \\ 0 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix}$$

https://en.wikipedia.org/wiki/Rank_(linear_algebra)

# Identity Matrix I

- Any vector or matrix multiplied by I remains unchanged
- For a matrix $A_{m \times n}$, $AI_n = I_m A = A$

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 3}$$

$$Iv = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$$

# Inverse of a Matrix

- The product of a **square** matrix $A$ and its inverse matrix $A^{-1}$ produces the identity matrix $I$
- $AA^{-1} = A^{-1}A = I$
- Inverse matrix is square, but not all square matrices has inverses

# Pseudo Inverse

- Non-square matrix and have left-inverse or right-inverse matrix
- Example:

$$Ax = b, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^n$$

 – Create a square matrix $A^T A$

$$A^T A x = A^T b$$

 – Multiplied both sides by inverse matrix $(A^T A)^{-1}$

$$x = (A^T A)^{-1} A^T b$$

– $(A^T A)^{-1} A^T$ is the pseudo inverse function

# Norm

- Norm is a measure of a vector's magnitude
- $l_2$ norm

$$\|x\|_2 = \left( |x_1|^2 + |x_2|^2 + \ldots + |x_n|^2 \right)^{1/2} = (x.x)^{1/2} = \left(x^T x\right)^{1/2}$$
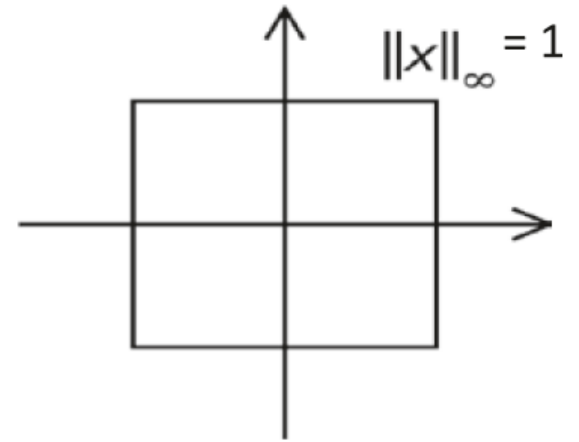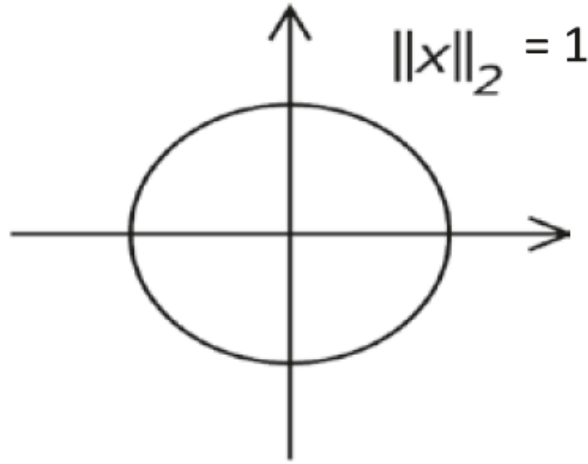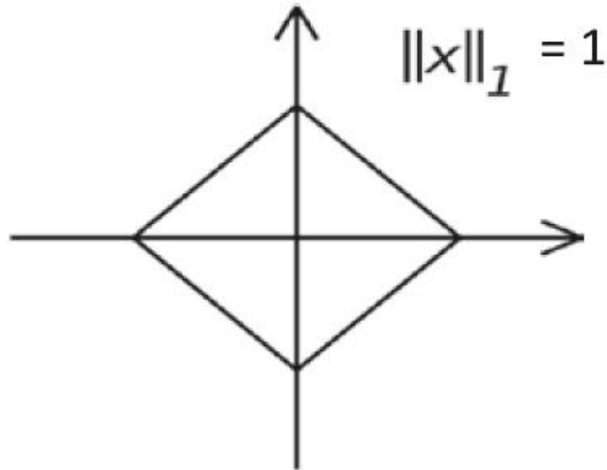
- $l_1$ norm

$$\|x\|_1 = |x_1| + |x_2| + \ldots + |x_n|$$

- $l_p$ norm

$$\left( |x_1|^p + |x_2|^p + \ldots + |x_n|^p \right)^{1/p}$$

- $l_\infty$ norm

$$\lim_{p \to \infty} \|x\|_p = \lim_{p \to \infty} \left( |x_1|^p + |x_2|^p + \ldots + |x_n|^p \right)^{1/p} = max(x_1, x_2, \ldots, x_n)$$

# Unit norms in 2D Vectors

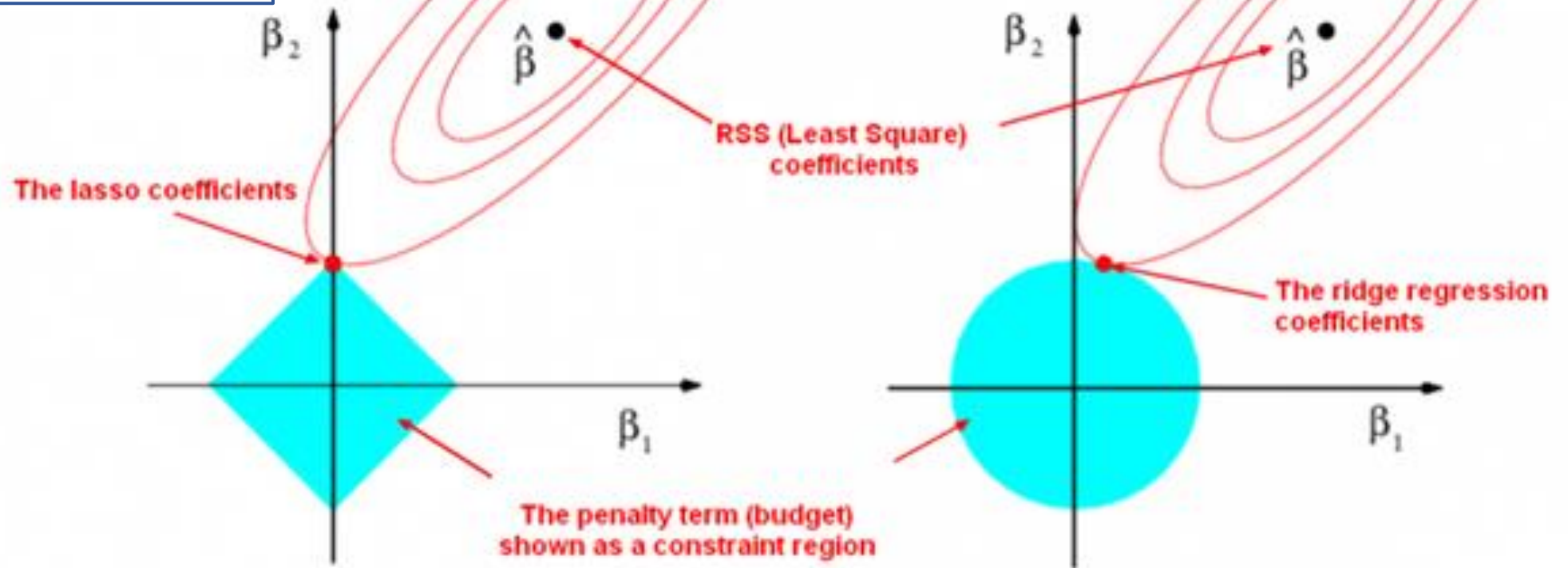- The set of all vectors of norm 1 in different 2D norms

# L1 and L2 Regularization

$Min \sum_{i=1}^{N}(y_i - \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i})^2$

subject to

$\sum_{j=1}^{2} \beta_j^2 <= s$

contours of RSS as
it move away from
the minimum

$\beta_2$

$\hat{\beta}$ ●

RSS (Least Square)
coefficients

The lasso coefficients

$\beta_2$

$\hat{\beta}$ ●

The ridge regression
coefficients

$\beta_1$

The penalty term (budget)
shown as a constraint region

$\beta_1$
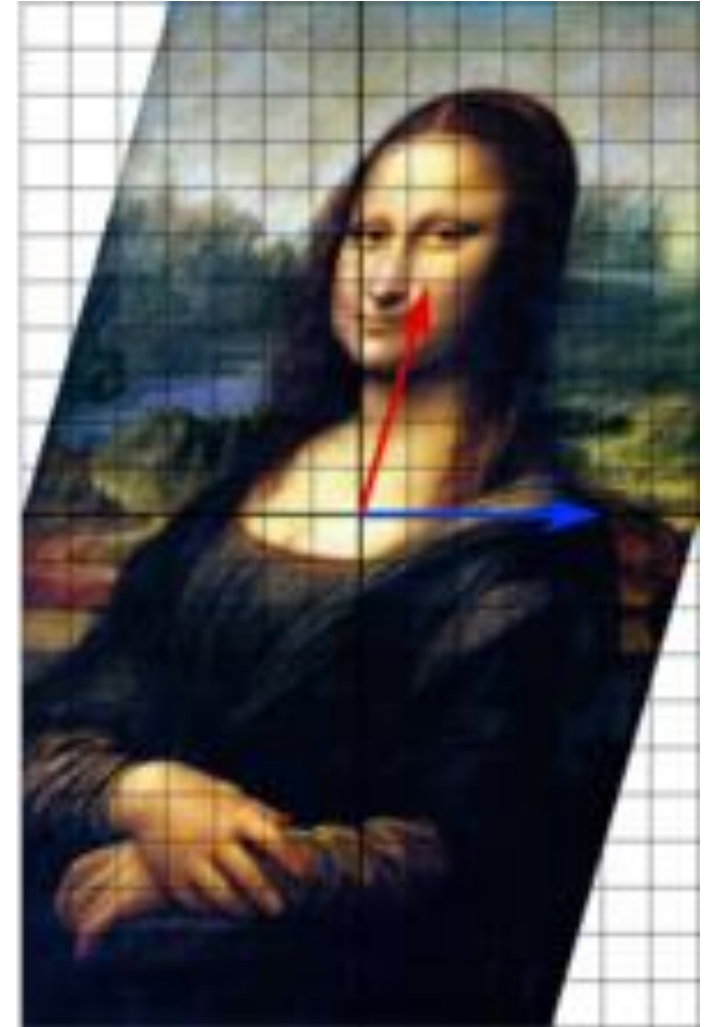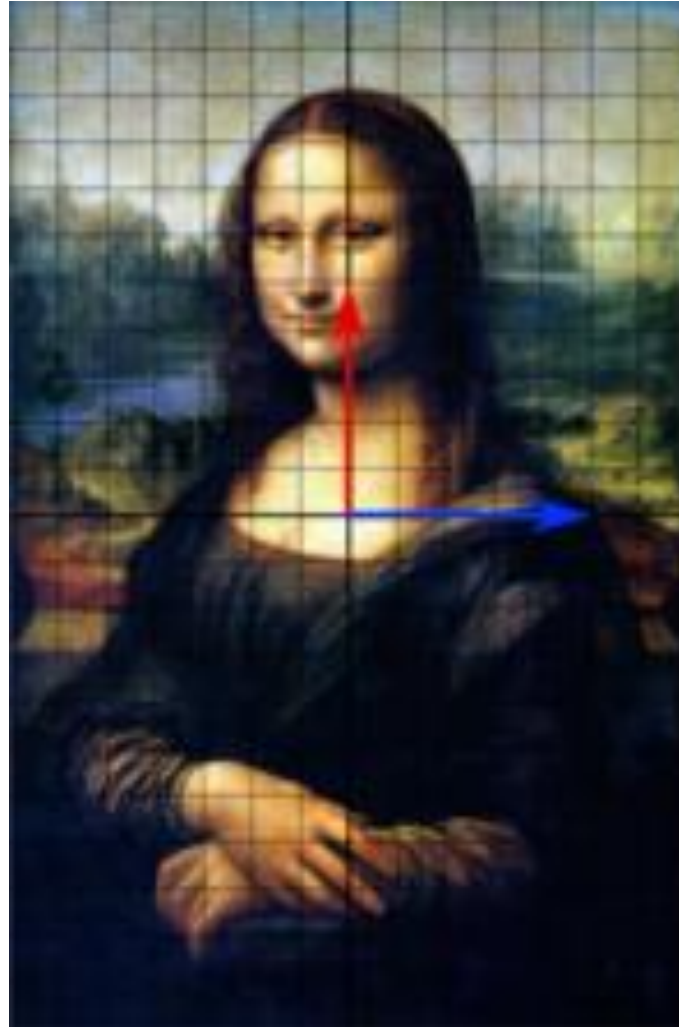
**LASSO**

**RIDGE REGRESSION**

# Eigen Vectors

- Eigenvector is a non-zero vector that changed by only a scalar factor λ when linear transform $A$ is applied to:

$$Ax = \lambda x, A \in \mathbb{R}^{n \times n}, x \in \mathbb{R}^n$$

- $x$ are Eigenvectors and $\lambda$ are Eigenvalues

- One of the most important concepts for machine learning, ex:
  - Principle Component Analysis (PCA)
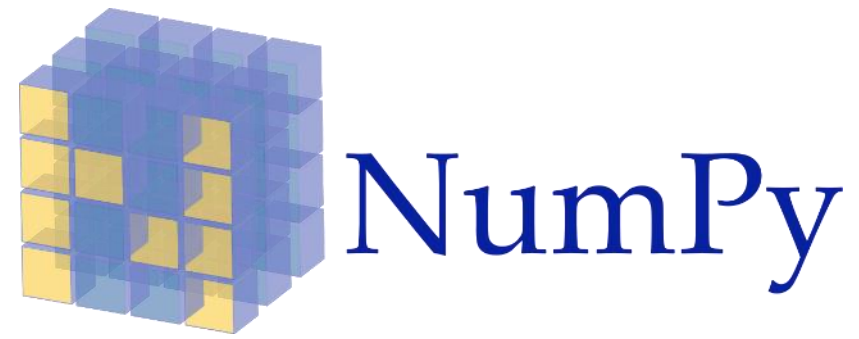  - Eigenvector centrality
  - PageRank
  - …

# Example: Shear Mapping

- Horizontal axis is the Eigenvector

# Power Iteration Method for Computing Eigenvector

1. Start with random vector $v$

2. Calculate iteratively: $v^{(k+1)} = A^k v$

3. After $v^k$ converges, $v^{(k+1)} \cong v^k$

4. $v^k$ will be the Eigenvector with largest Eigenvalue

# NumPy for Linear Algebra

- NumPy is the fundamental package for scientific computing with Python. It contains among other things:
  - a powerful N-dimensional array object
  - sophisticated (broadcasting) functions
  - tools for integrating C/C++ and Fortran code
  - useful linear algebra, Fourier transform, and random number capabilities

# Python & NumPy tutorial

- http://cs231n.github.io/python-numpy-tutorial/
- Stanford CS231n: Convolutional Neural Networks for Visual Recognition
  - http://cs231n.stanford.edu/

# Create Tensors

## Scalars (0D tensors)

```
>>> import numpy as np
>>> x = np.array(12)
>>> x
array(12)
>>> x.ndim
0
```

## Vectors (1D tensors)

```
>>> x = np.array([12, 3, 6, 14])
>>> x
array([12, 3, 6, 14])
>>> x.ndim
1
```

## Matrices (2D tensors)

```
>>> x = np.array([[5, 78, 2, 34, 0],
                  [6, 79, 3, 35, 1],
                  [7, 80, 4, 36, 2]])
>>> x.ndim
2
```

# Create 3D Tensor

```
>>> x = np.array([[[5, 78, 2, 34, 0],
                   [6, 79, 3, 35, 1],
                   [7, 80, 4, 36, 2]],
                  [[5, 78, 2, 34, 0],
                   [6, 79, 3, 35, 1],
                   [7, 80, 4, 36, 2]],
                  [[5, 78, 2, 34, 0],
                   [6, 79, 3, 35, 1],
                   [7, 80, 4, 36, 2]]])
>>> x.ndim
3
```

# Attributes of a Tensor

- **Number of axes (dimensions)**
  - x.ndim

- **Shape**
  - This is a tuple of integers showing how many data the tensor has along each axis

- **Data type**
  - uint8, float32 or float64

# Manipulating Tensors in Numpy

```
>>> my_slice = train_images[10:100]
>>> print(my_slice.shape)
(90, 28, 28)
```

```
>>> my_slice = train_images[10:100, :, :]        Equivalent to the
>>> my_slice.shape                               previous example
(90, 28, 28)
>>> my_slice = train_images[10:100, 0:28, 0:28]  Also equivalent to the
>>> my_slice.shape                               previous example
(90, 28, 28)
```

```
my_slice = train_images[:, 7:-7, 7:-7]
```

# Displaying the Fourth Digit

```
digit = train_images[4]

import matplotlib.pyplot as plt
plt.imshow(digit, cmap=plt.cm.binary)
plt.show()
```
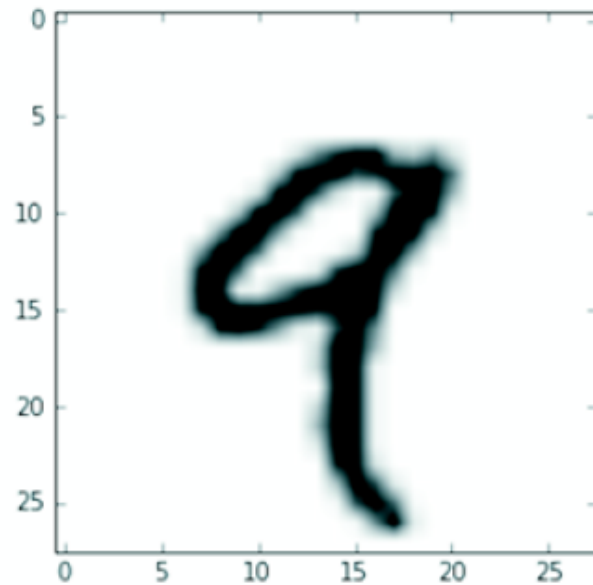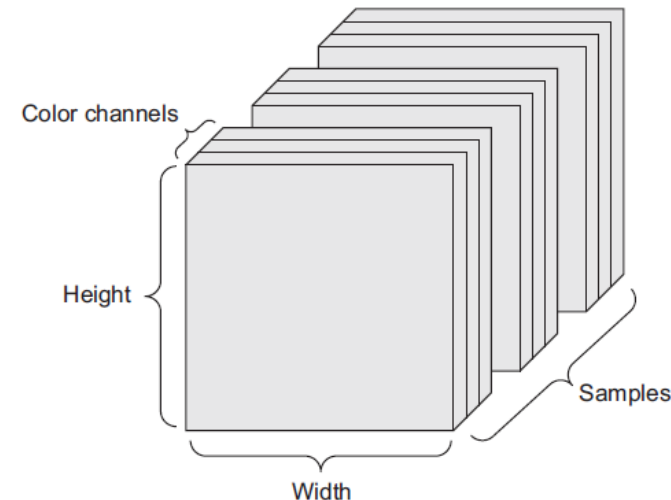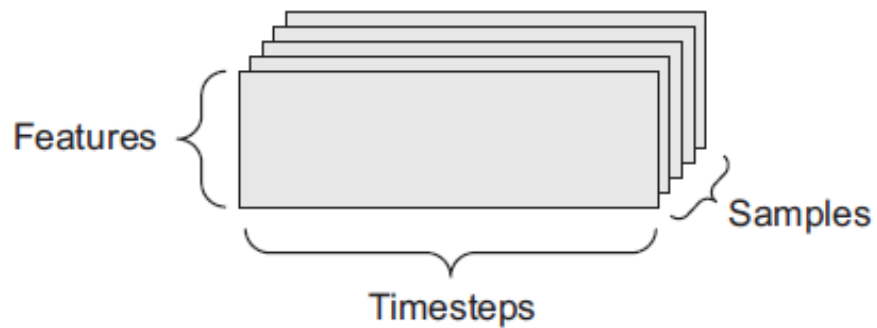


**Figure 2.2   The fourth sample in our dataset**

# Real-world examples of Data Tensors

- Vector data – 2D (samples, features)
- Timeseries Data – 3D (samples, timesteps, features)
- Images – 4D (samples, height, width, channels)
- Video – 5D (samples, frames, height, width, channels)

# Batch size & Epochs

- ## A sample
  - A sample is a single row of data

- ## Batch size
  - Number of samples used for one iteration of gradient descent
  - Batch size = 1: stochastic gradient descent
  - 1 < Batch size < all: mini-batch gradient descent
  - Batch size = all: batch gradient descent

- ## Epoch
  - Number of times that the learning algorithm work through all training samples

# Element-wise Operations for Matrix

- Operate on each element

```
def naive_add(x, y):
    assert len(x.shape) == 2
    assert x.shape == y.shape

    x = x.copy()
    for i in range(x.shape[0]):
        for j in range(x.shape[1]):
            x[i, j] += y[i, j]
    return x
```

**x and y are 2D Numpy tensors.**

**Avoid overwriting the input tensor.**

# NumPy Operation for Matrix

• Leverage the Basic Linear Algebra subprograms (BLAS)

• BLAS is optimized using C or Fortran

```
import numpy as np

z = x + y                    ⟵———— Element-wise addition

z = np.maximum(z, 0.)               ⟵——— Element-wise relu
```

# Broadcasting

- Apply smaller tensor repeated to the extra axes of the larger tensor

```
def naive_add_matrix_and_vector(x, y):        x is a 2D Numpy tensor.
    assert len(x.shape) == 2
    assert len(y.shape) == 1                  y is a Numpy vector.
    assert x.shape[1] == y.shape[0]

    x = x.copy()                              Avoid overwriting
    for i in range(x.shape[0]):               the input tensor.
        for j in range(x.shape[1]):
            x[i, j] += y[j]
    return x
```
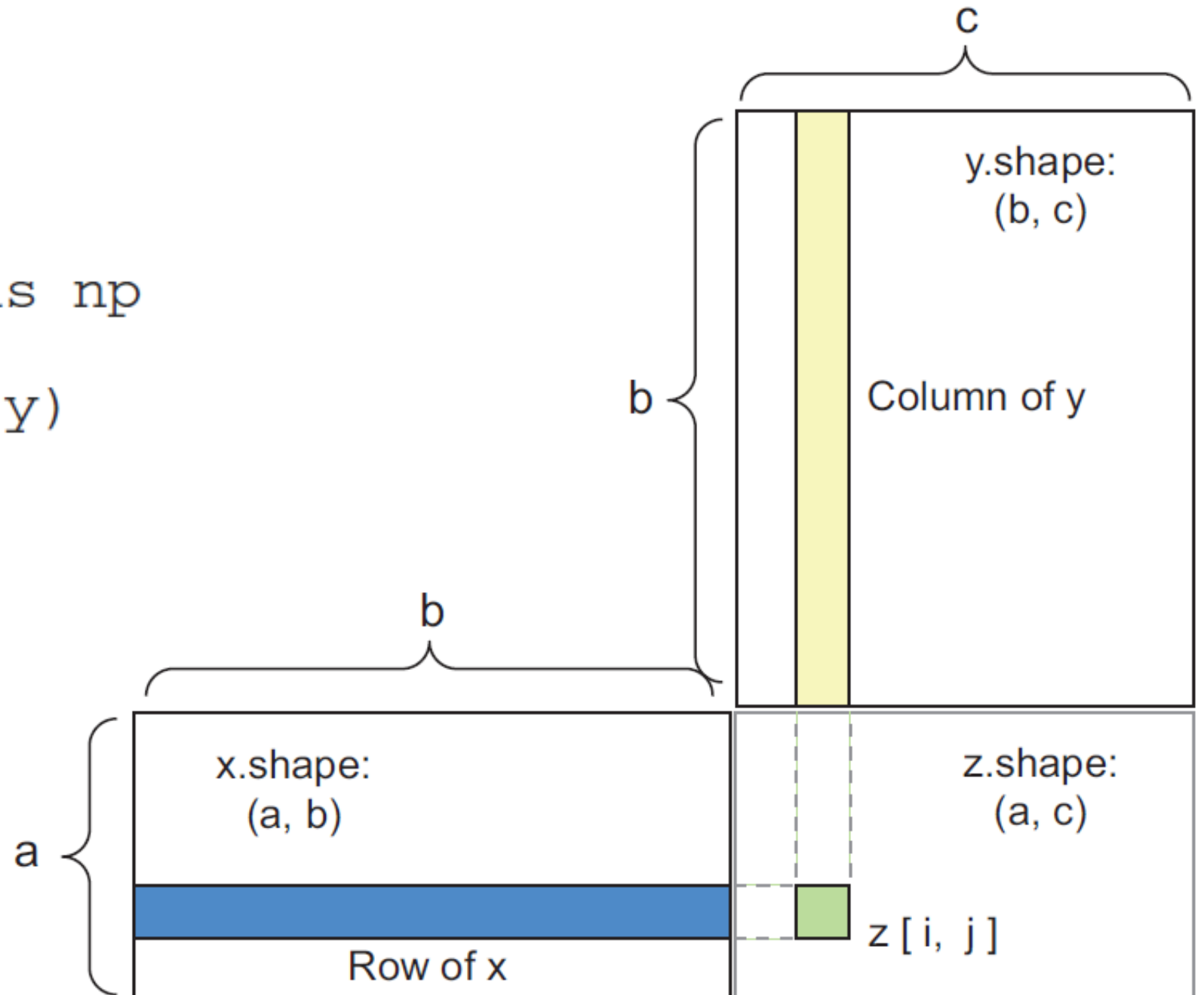
# Tensor Dot

```
import numpy as np

z = np.dot(x, y)
```

y.shape:
(b, c)

Column of y

x.shape:
(a, b)

z.shape:
(a, c)

z[i, j]

Row of x

# Implementation of Dot Product

**x and y are Numpy matrices.**

```python
def naive_matrix_dot(x, y):
    assert len(x.shape) == 2
    assert len(y.shape) == 2
    assert x.shape[1] == y.shape[0]

    z = np.zeros((x.shape[0], y.shape[1]))
    for i in range(x.shape[0]):
        for j in range(y.shape[1]):
            row_x = x[i, :]
            column_y = y[:, j]
            z[i, j] = naive_vector_dot(row_x, column_y)
    return z
```

**The first dimension of x must be the same as the 0th dimension of y!**

**This operation returns a matrix of 0s with a specific shape.**

**Iterates over the rows of x ...**

**... and over the columns of y.**

# Tensor Reshaping

- Rearrange a tensor's rows and columns to match a target shape

```
>>> x = np.array([[0., 1.],
                  [2., 3.],
                  [4., 5.]])
>>> print(x.shape)
(3, 2)

>>> x = x.reshape((6, 1))
>>> x
array([[ 0.],
       [ 1.],
       [ 2.],
       [ 3.],
       [ 4.],
       [ 5.]])

>>> x = x.reshape((2, 3))
>>> x
array([[ 0.,  1.,  2.],
       [ 3.,  4.,  5.]])
```

# Matrix Transposition

- Transposing a matrix means exchanging its rows and its columns

```
>>> x = np.zeros((300, 20))
>>> x = np.transpose(x)
>>> print(x.shape)
(20, 300)
```

Creates an all-zeros matrix of shape (300, 20)

# Unfolding the Manifold

- Tensor operations are complex geometric transformation in high-dimensional space
  - Dimension reduction

Maria Gaetana Agnesi

$(\ln x)' = \frac{1}{x}$  $\int \frac{1}{x} dx = |\ln|x|| + c$  $\int \sin x \, dx = -\cos x \, dx + c$

$f(x) = x^2$

$\int_a^b f'(x) dx = f(b) - f(a)$

$x^2 - 3x - 4 = 0$

$4x^2 - 3x - 1 = 0$

$m \frac{d^2 x}{dt^2} = -kx$

$\int f(x) dx$

# Calculus

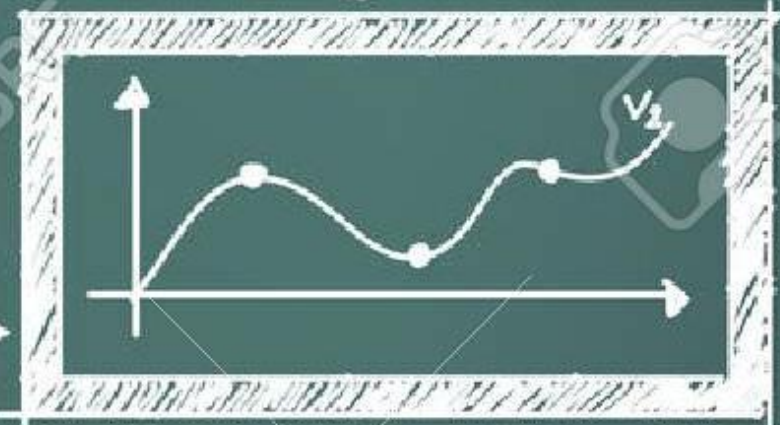$\frac{df(x)}{dz}$

$\frac{dA}{dt} = \frac{dB}{dt} = -\frac{dC}{dt} = -\frac{dD}{dt} = (d_1)T^{\frac{1}{2}}AB - (d_2)T^{\frac{1}{2}}CD$

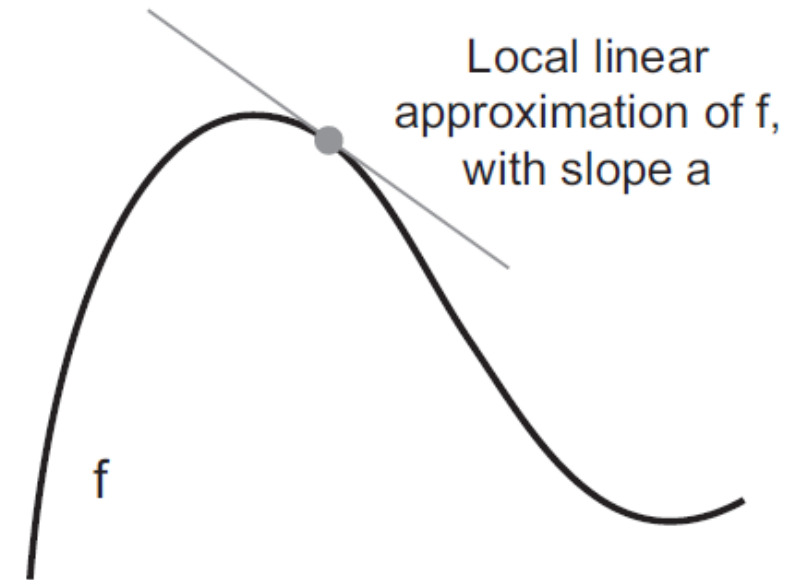$x^2 = A$  $\frac{dT}{dt} = (c_3)\frac{dA}{dt} - (c_4)(T_0 - T)$

# Differentiation

$$\frac{df}{dt} = \lim_{h \to 0} \frac{f(t+h) - f(t)}{h}$$

OR

$$\frac{df}{dt} = \lim_{h \to 0} \frac{f(t+h) - f(t-h)}{2h}$$

Local linear approximation of f, with slope a

f

# Gradient of a Function

- Gradient is a multi-variable generalization of the derivative
- Apply partial derivatives

$$\nabla f = \left[ \frac{\partial f}{\partial x_1} \frac{\partial f}{\partial x_2} \cdots \cdots \frac{\partial f}{\partial x_n} \right]^T$$

- Example

$$f(x, y, z) = x + y^2 + z^3$$

$$\nabla f = \left[ 1 \ 2y \ 3z^2 \right]^T$$

# Hessian Matrix

- Second-order partial derivatives

$$Hf = \begin{bmatrix} \dfrac{\delta^2 f}{\delta x^2} & \dfrac{\delta^2 f}{\delta x \delta y} & \dfrac{\delta^2 f}{\delta x \delta z} \\[2em] \dfrac{\delta^2 f}{\delta y \delta x} & \dfrac{\delta^2 f}{\delta y^2} & \dfrac{\delta^2 f}{\delta y \delta z} \\[2em] \dfrac{\delta^2 f}{\delta z \delta x} & \dfrac{\delta^2 f}{\delta z \delta y} & \dfrac{\delta^2 f}{\delta z^2} \end{bmatrix}$$
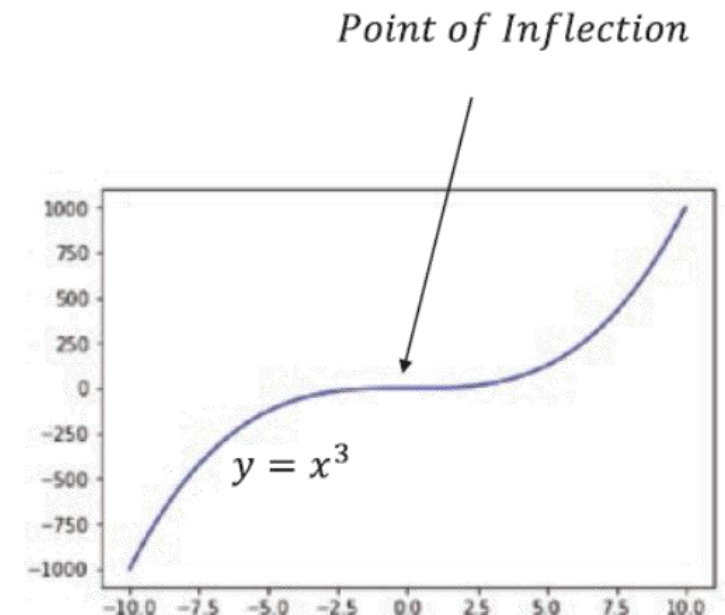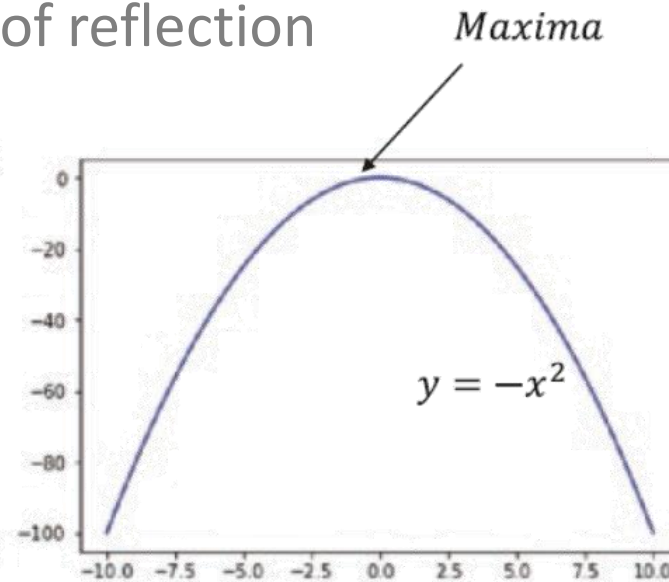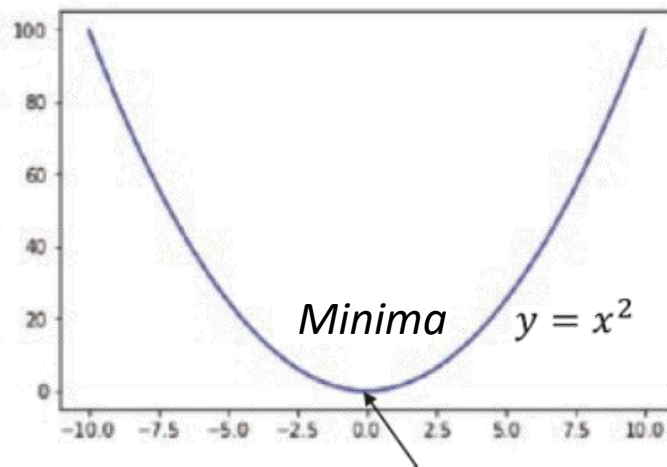
# Maxima and Minima for Univariate Function

- If $\dfrac{df(x)}{dx} = 0$, it's a minima or a maxima point, then we study the second derivative:

  - If $\dfrac{d^2f(x)}{dx^2} < 0$ => Maxima
  - If $\dfrac{d^2f(x)}{dx^2} > 0$ => Minima
  - If $\dfrac{d^2f(x)}{dx^2} = 0$ => Point of reflection

# Maxima and Minima for Multivariate Function

- Computing the gradient and setting it to zero vector would give us the list of stationary points.

- For a stationary point $x_0 \in \mathbb{R}^n$
  - If the Hessian matrix of the function at $x_0$ has both positive and negative eigen values, then $x_0$ is a saddle point
  - If the eigen values of the Hessian matrix are all positive then the stationary point is a local minima
  - If the eigen values are all negative then the stationary point is a local maxima

# Chain Rule

$$\frac{dy}{dx} = \frac{dy}{du}\frac{du}{dx}$$

$$\frac{d^2y}{dx^2} = \frac{d^2y}{du^2}\left(\frac{du}{dx}\right)^2 + \frac{dy}{du}\frac{d^2u}{dx^2}$$

$$\frac{d^3y}{dx^3} = \frac{d^3y}{du^3}\left(\frac{du}{dx}\right)^3 + 3\frac{d^2y}{du^2}\frac{du}{dx}\frac{d^2u}{dx^2} + \frac{dy}{du}\frac{d^3u}{dx^3}$$

$$\frac{d^4y}{dx^4} = \frac{d^4y}{du^4}\left(\frac{du}{dx}\right)^4 + 6\frac{d^3y}{du^3}\left(\frac{du}{dx}\right)^2\frac{d^2u}{dx^2} + \frac{d^2y}{du^2}\left(4\frac{du}{dx}\frac{d^3u}{dx^3} + 3\left(\frac{d^2u}{dx^2}\right)^2\right) + \frac{dy}{du}\frac{d^4u}{dx^4}.$$

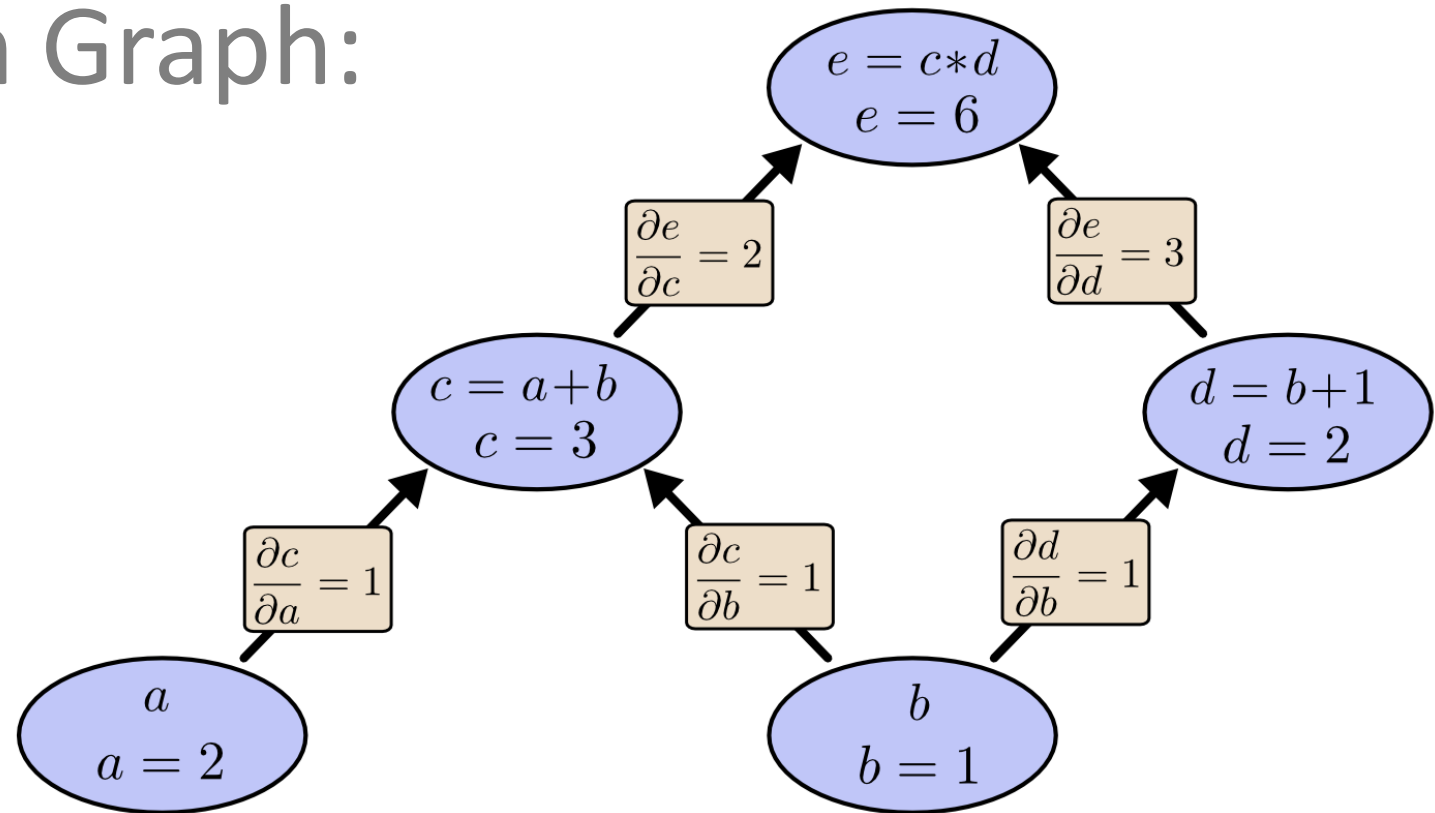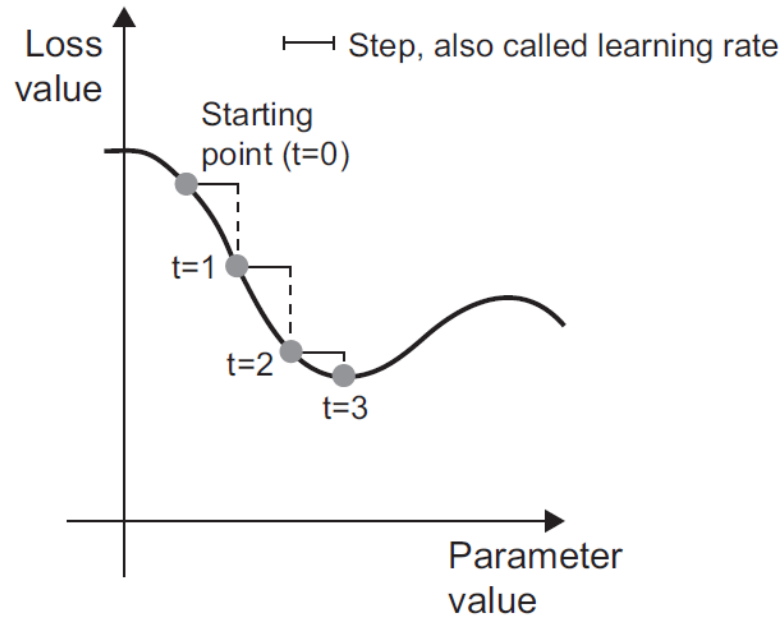# Symbolic Differentiation

Computation Graph:

c = a + b

d = b + 1

e = c*d

# Stochastic Gradient Descent



1. Draw a batch of training samples x and corresponding targets y

2. Run the network on x to obtain predictions y_pred

3. Compute the loss of the network on the batch, a measure of the mismatch between y_pred and y

4. Compute the gradient of the loss with regard to the network's parameters (a backward pass).

5. Move the parameters a little in the opposite direction from the gradient: W -= step * gradient

# Gradient Descent along a 2D Surface

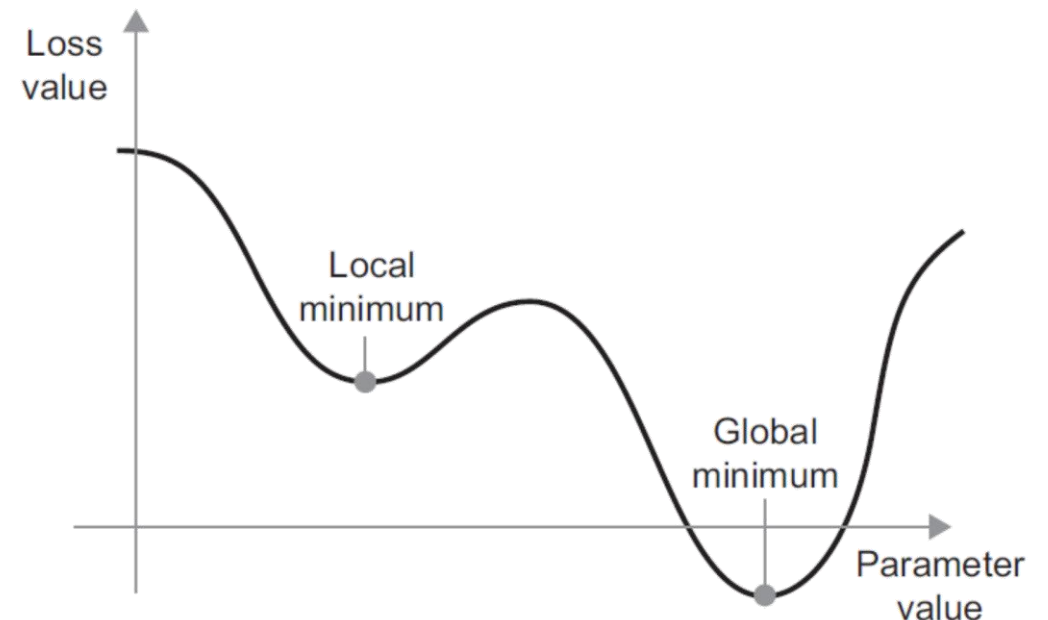# Avoid Local Minimum using Momentum

```
past_velocity = 0.
momentum = 0.1
while loss > 0.01:
    w, loss, gradient = get_current_parameters()
    velocity = past_velocity * momentum + learning_rate * gradient
    w = w + momentum * velocity - learning_rate * gradient
    past_velocity = velocity
    update_parameter(w)
```

**Constant momentum factor**

**Optimization loop**

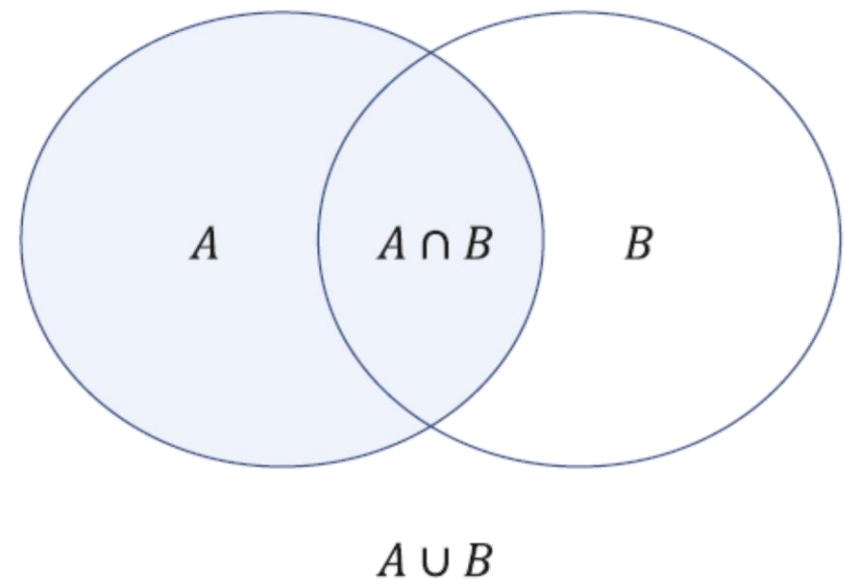Basics of Probability

# Three Axioms of Probability

- Given an Event $E$ in a sample space $S$, $S = \bigcup_{i=1}^{N} E_i$

- First axiom
  - $P(E) \in \mathbb{R}, 0 \leq P(E) \leq 1$

- Second axiom
  - $P(S) = 1$

- Third axiom
  - Additivity, any countable sequence of mutually exclusive events $E_i$
  - $P(\bigcup_{i=1}^{n} E_i) = P(E_1) + P(E_2) + \cdots + P(E_n) = \sum_{i=1}^{n} P(E_i)$

# Union, Intersection, and Conditional Probability

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(A \cap B)$ is simplified as $P(AB)$
- Conditional Probability $P(A|B)$, the probability of event A given B has occurred

  - $P(A|B) = P\left(\dfrac{AB}{B}\right)$
  - $P(AB) = P(A|B)P(B) = P(B|A)P(A)$



$A$ $\quad$ $A \cap B$ $\quad$ $B$

$A \cup B$

# Chain Rule of Probability

- The joint probability can be expressed as chain rule

$$P\left(A_1 A_2 A_3 \ldots A_n\right) = P\left(A_1\right) P\left(A_2 / A_1\right) P\left(A_3 / A_1 A_2\right) \ldots P\left(A_n / A_1 A_2 .. A_{(n-1)}\right)$$

$$= P\left(A_1\right) \prod_{i=2}^{n} P\left(A_i / A_1 A_2 A_3 \ldots A_{(n-1)}\right)$$

# Mutually Exclusive

- $P(AB) = 0$
- $P(A \cup B) = P(A) + P(B)$

# Independence of Events

- Two events A and B are said to be independent if the probability of their intersection is equal to the product of their individual probabilities
  - $P(AB) = P(A)P(B)$
  - $P(A|B) = P(A)$

# Bayes Rule

- $P(A|B) = \dfrac{P(B|A)P(A)}{P(B)}$

- Proof:
  - Remember $P(A|B) = P\left(\dfrac{AB}{B}\right)$
  - So $P(AB) = P(A|B)P(B) = P(B|A)P(A)$
  - Then Bayes $P(A|B) = P(B|A)P(A)/P(B)$

# Probability Mass Function and Dense Function

- ## Probability mass function (PMF)
  - Function that gives the probability that a discrete random variable is exactly equal to some value

$$P(X = i) = \frac{1}{6}, i \in \{1,2,3,4,5,6\}$$

- ## Probability dense function (PDF)
  - Specify the probability of the random variable falling within a particular range of values

$$\int_D P(x)dx = 1$$

# Expectation of a Random Variable

- Expectation of a discrete random variable

$$E[X] = x_1 p_1 + x_2 p_2 + \cdots + x_n p_n = \sum_{i=1}^{n} x_i p_i$$

- Expectation of a continuous random variable

$$E[X] = \int_D x P(x) dx$$

# Variance of a Random Variable

- Expectation of a discrete random variable

$$Var[X] = E[(X - \mu)^2], \text{where } \mu = E[X]$$

- Expectation of a continuous random variable

$$Var[X] = \int_D (x - \mu)^2 P(x) dx$$

- Standard deviation $\sigma$ is the square root of variance

# Covariance and Correlation Coefficient

- Expectation of a discrete random variable

$$Cov(X, Y) = E\big[(X - \mu_x)(Y - \mu_y)\big],$$
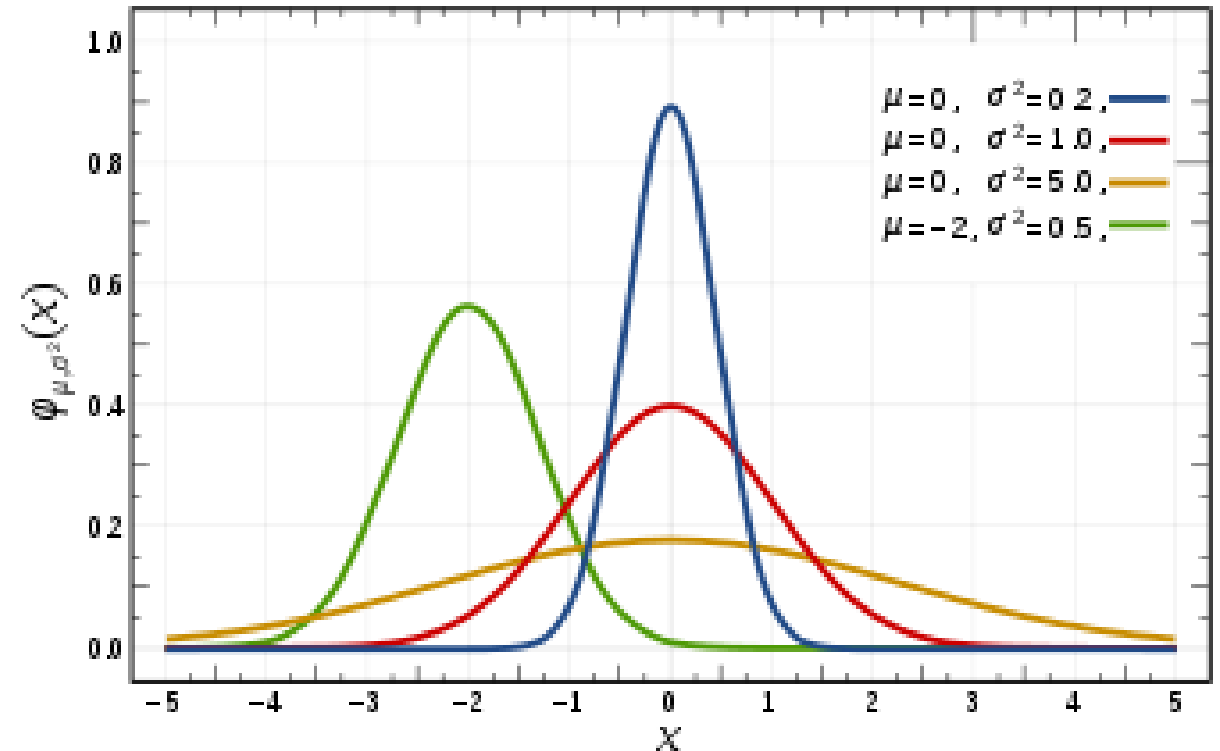$$\text{where } \mu_x = E[X], \mu_y = E[Y]$$

- Correlation coefficient

$$\rho = \frac{cov(X, Y)}{\sigma_x \sigma_y}$$

# Normal (Gaussian) Distribution

- One of the most important distributions
- Central limit theorem
  - Averages of samples of observations of random variables independently drawn from independent distributions converge to the normal distribution

$$f(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

# Optimization

The *standard form* of a continuous optimization problem is[1]

$$\underset{x}{\text{minimize}} \quad f(x)$$

$$\text{subject to} \quad g_i(x) \le 0, \quad i = 1, \ldots, m$$

$$h_j(x) = 0, \quad j = 1, \ldots, p$$

where

- $f : \mathbb{R}^n \to \mathbb{R}$ is the **objective function** to be minimized over the *n*-variable vector $x$,
- $g_i(x) \le 0$ are called **inequality constraints**
- $h_j(x) = 0$ are called **equality constraints**, and
- $m \ge 0 \; and \; p \ge 0$.

https://en.wikipedia.org/wiki/Optimization_problem

# Formulate Your Problem

- Linear model: $f(x) = w^T x + b$

- Least-squared Error: $(f(x) - y)^2$

- Regularization: $\|w\|$

- Objective function:

$$\min_{w} \left(w^T x - y\right)^2 + \lambda \|w\|$$

# References

- Francois Chollet, "Deep Learning with Python," Chapter 2 "Mathematical Building Blocks of Neural Networks"

- Santanu Pattanayak, "Pro Deep Learning with TensorFlow," Apress, 2017

- [Machine Learning Cheat Sheet](#)

- https://machinelearningmastery.com/difference-between-a-batch-and-an-epoch/

- https://www.quora.com/What-is-the-difference-between-L1-and-L2-regularization-How-does-it-solve-the-problem-of-overfitting-Which-regularizer-to-use-and-when

- Wikipedia