

# Applied Math for Deep Learning

Prof. Kuan-Ting Lai

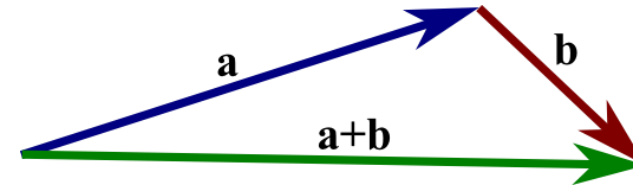
2021/3/8

# Applied Math for Deep Learning

- Linear Algebra
- Probability
- Calculus
- Optimization

# Linear Algebra

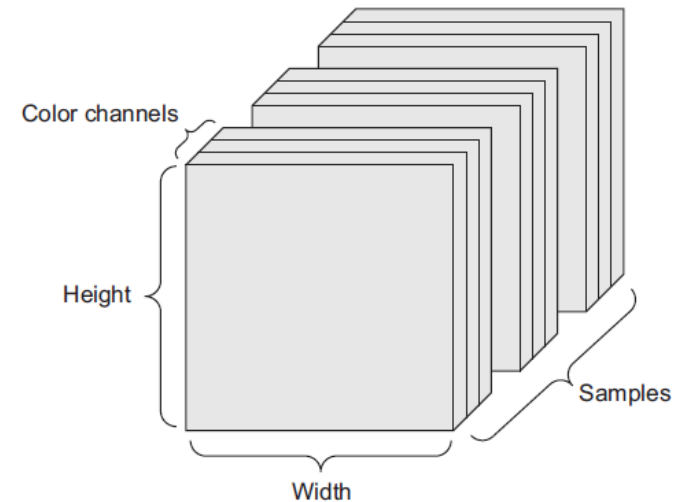
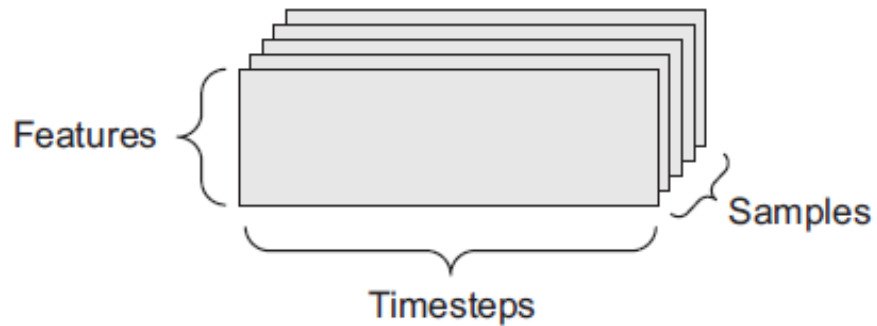
- Scalar
  - real numbers
- Vector (1D)
  - Has a magnitude & a direction
- Matrix (2D)
  - An array of numbers arranged in rows & columns
- Tensor ( $\geq 3D$ )
  - Multi-dimensional arrays of numbers



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

# 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)



# Vector Dimension vs. Tensor Dimension

- The number of data in a vector is also called “dimension”
- In deep learning , the dimension of Tensor is also called “rank”
- Matrix = 2d array = 2d tensor = rank 2 tensor
- Axis means the specific dimension of a Tensor



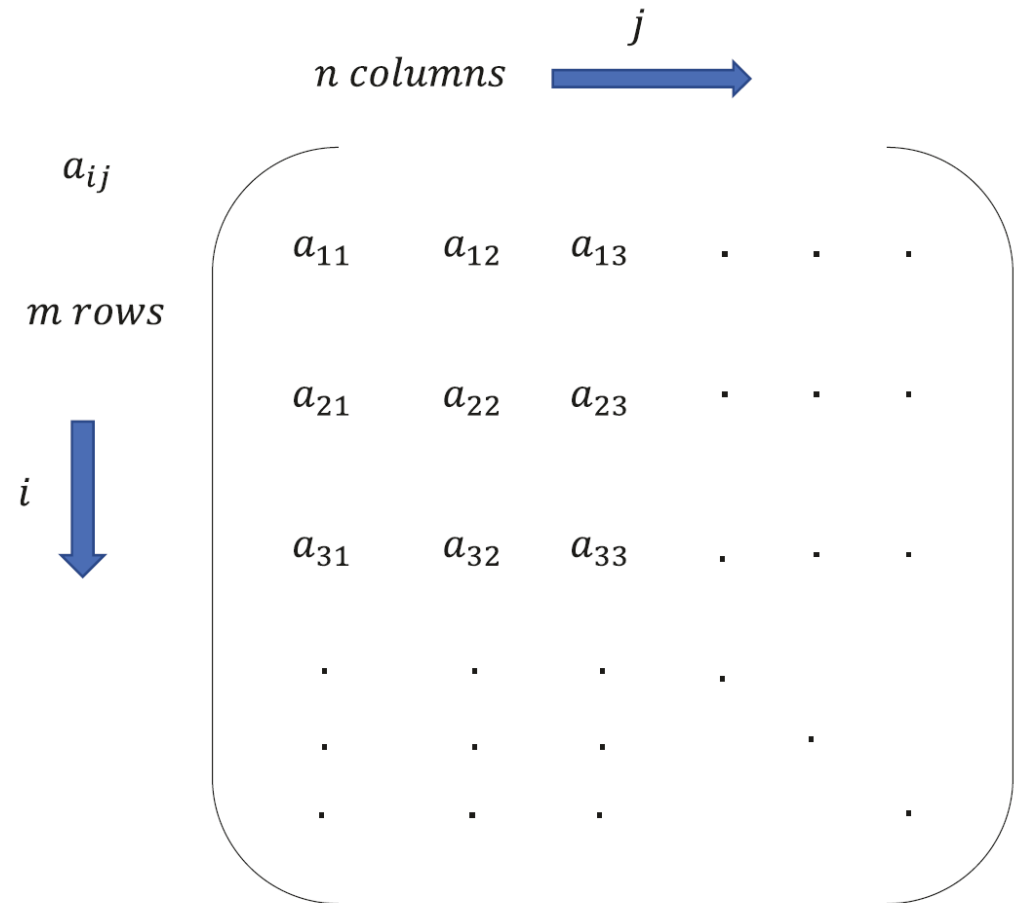
A full-page image of Keanu Reeves as Neo from the movie The Matrix. He is wearing a black turtleneck and dark sunglasses, looking directly at the camera. The background is a green digital rain effect, with white text resembling code or data floating around him. The title 'The Matrix' is centered in white text over his chest.

# The Matrix

# Matrix

- Define a matrix with  $m$  rows and  $n$  columns:

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



# Matrix Operations

- Addition and Subtraction

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

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

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

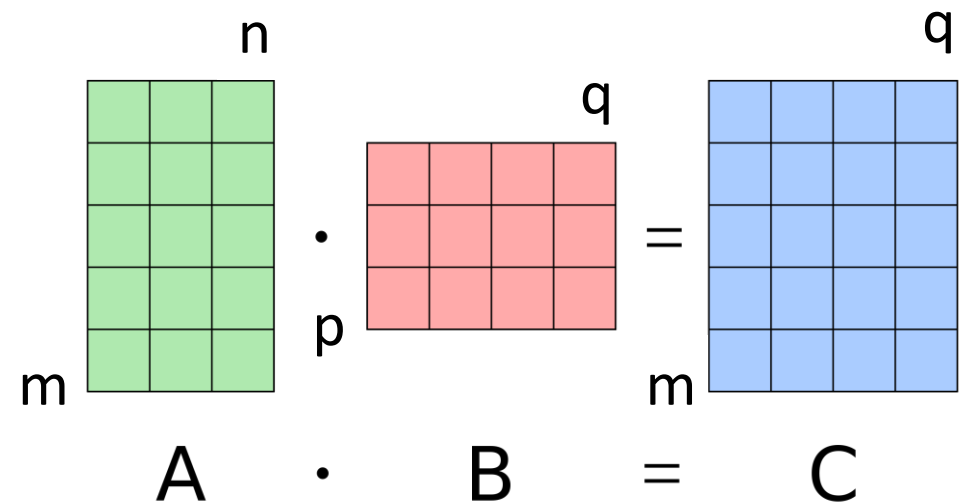


# Matrix Multiplication

- 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)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 & 64 \\ 82 & 91 \end{bmatrix}$$

## Example of Matrix Multiplication (3-3)

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



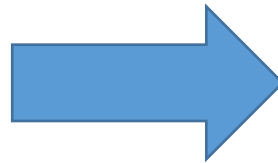
# Matrix Transpose

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

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

**A**

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



**A<sup>T</sup>**

1	3	5
2	4	6

# Dot Product

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

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

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

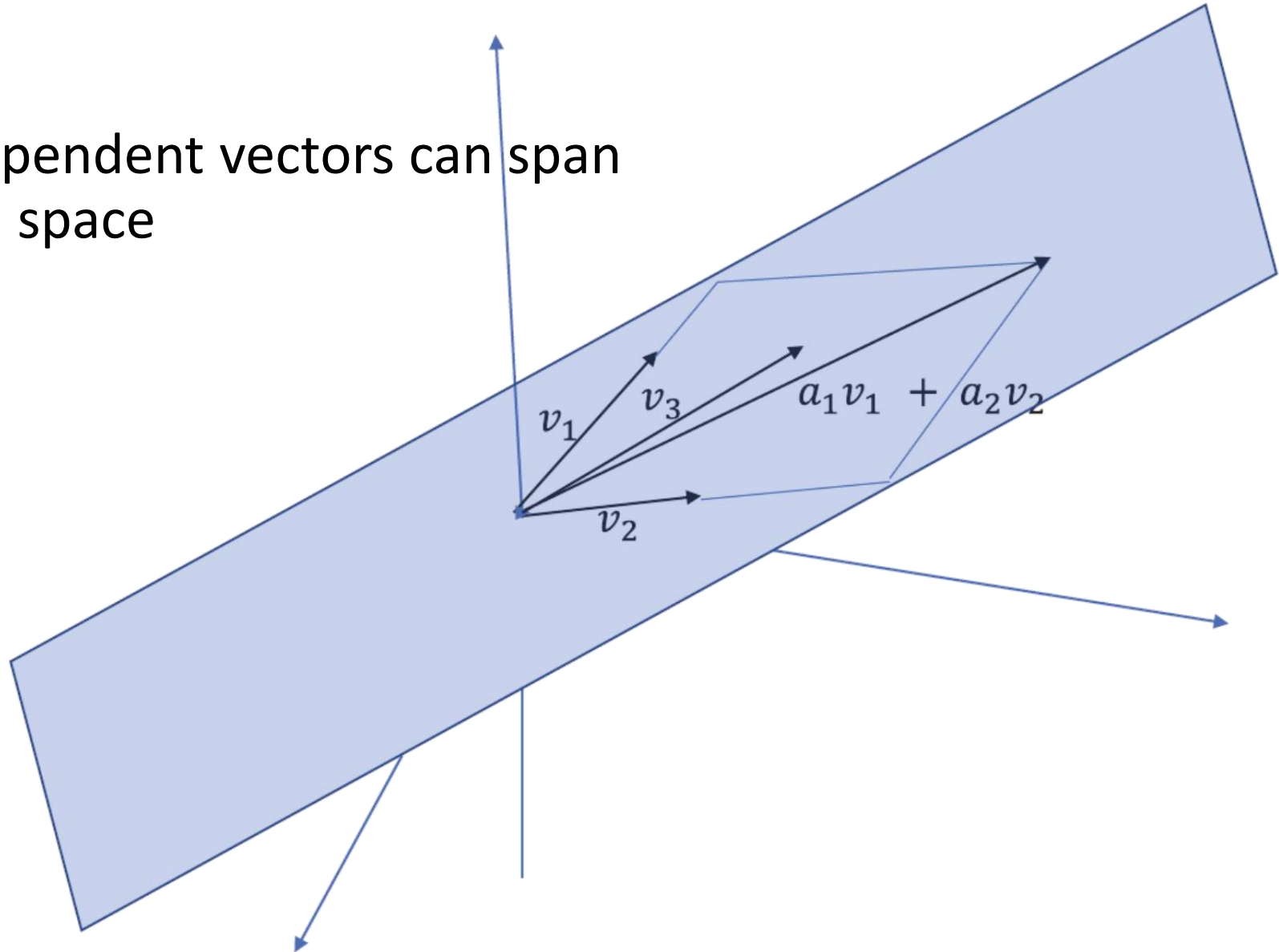
# Linear Independence

- A vector is **linearly dependent** on other vectors if it can be expressed as the linear combination of other vectors
- A set of vectors  $v_1, v_2, \dots, v_n$  is **linearly independent** if  $a_1 v_1 + a_2 v_2 + \dots + a_n v_n = 0$  implies all  $a_i = 0, \forall i \in \{1, 2, \dots, n\}$

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

# 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:

$$\mathbf{A} = \begin{bmatrix} \mathbf{1} & \mathbf{2} & \mathbf{1} \\ \mathbf{-2} & \mathbf{-3} & \mathbf{1} \\ \mathbf{3} & \mathbf{5} & \mathbf{0} \end{bmatrix} \xrightarrow{\text{Row-echelon form}} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{-5} \\ \mathbf{0} & \mathbf{1} & \mathbf{3} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}$$

# Identity Matrix I

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

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

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

# 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 
$$\|x\|_2 = \left( |x_1|^2 + |x_2|^2 + \dots + |x_n|^2 \right)^{1/2} = (x \cdot x)^{1/2} = (x^T x)^{1/2}$$

- $l_1$  norm 
$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$$

- $l_p$  norm 
$$\left( |x_1|^p + |x_2|^p + \dots + |x_n|^p \right)^{1/p}$$

- $l_\infty$  norm

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

# Eigen Vectors

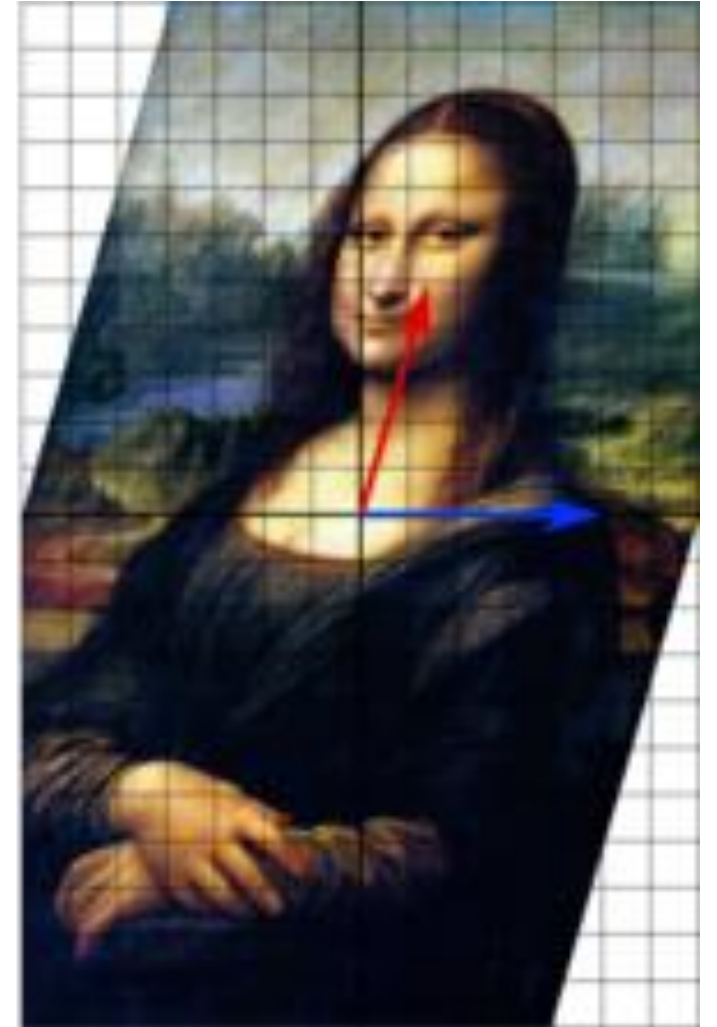
- Eigenvector is a non-zero vector that changed by only a scalar factor  $\lambda$  when linear transform  $A$  is applied to:

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

- $x$  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

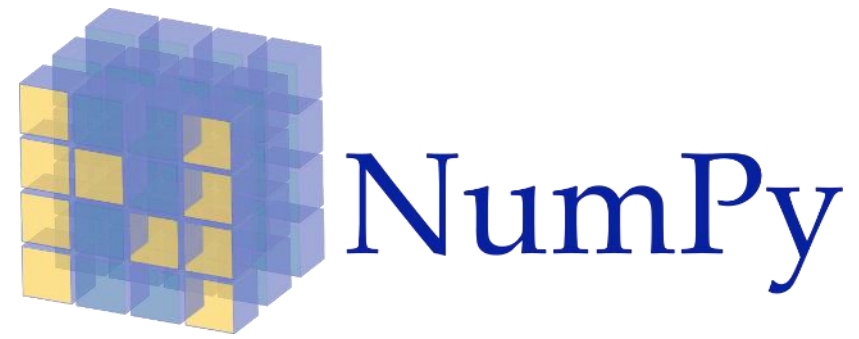
- Horizontal axis is the Eigenvector



# 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/>

## Table of contents:

- Python
  - Basic data types
  - Containers
    - Lists
    - Dictionaries
    - Sets
    - Tuples
  - Functions
  - Classes
- Numpy
  - Arrays
  - Array indexing
  - Datatypes
  - Array math
  - Broadcasting
- SciPy
  - Image operations
  - MATLAB files
  - Distance between points
- Matplotlib
  - Plotting
  - Subplots
  - Images



# 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

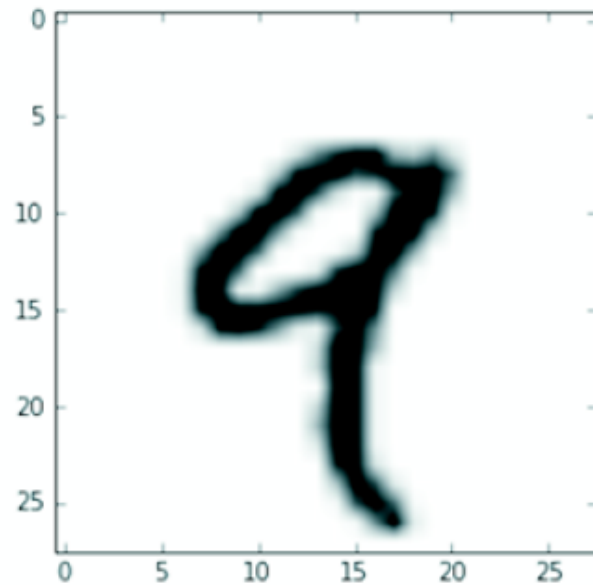
```
>>> my_slice = train_images[10:100]
>>> print(my_slice.shape)
(90, 28, 28)
```

```
>>> 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

# Numpy Multiplication

```
In [ ]:  ▶ import numpy as np
```

```
In [4]:  ▶ x = np.array([[1, 2, 3], [4, 5, 6]])  
x
```

```
Out[4]: array([[1, 2, 3],  
              [4, 5, 6]])
```

```
In [5]:  ▶ y = np.array([[7, 8], [9, 10], [11, 12]])  
y
```

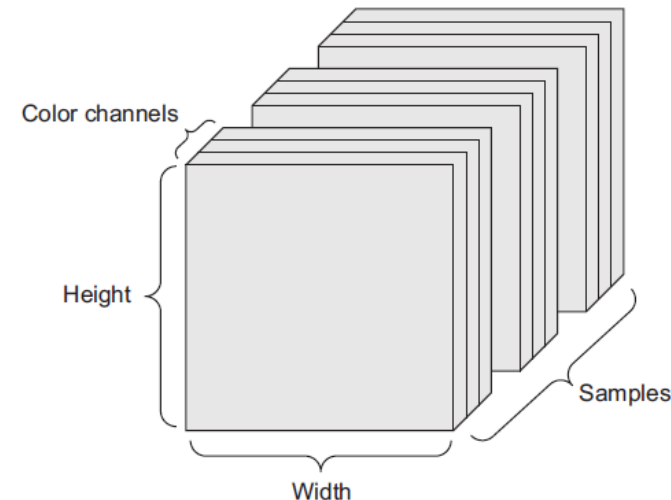
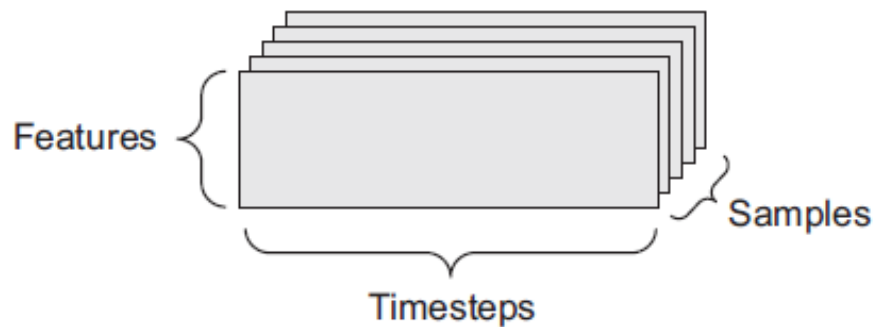
```
Out[5]: array([[ 7,  8],  
              [ 9, 10],  
              [11, 12]])
```

```
In [10]: ▶ np.matmul(x, y)
```

```
Out[10]: array([[ 58,  64],  
               [139, 154]])
```

# 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
  - $1 < \text{Batch size} < \text{all}$ : 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 and y are 2D  
Numpy tensors.**



```
    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
```

**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):  
    assert len(x.shape) == 2  
    assert len(y.shape) == 1  
    assert x.shape[1] == y.shape[0]  
  
    x = x.copy()  
    for i in range(x.shape[0]):  
        for j in range(x.shape[1]):  
            x[i, j] += y[j]  
    return x
```

**x is a 2D Numpy tensor.**

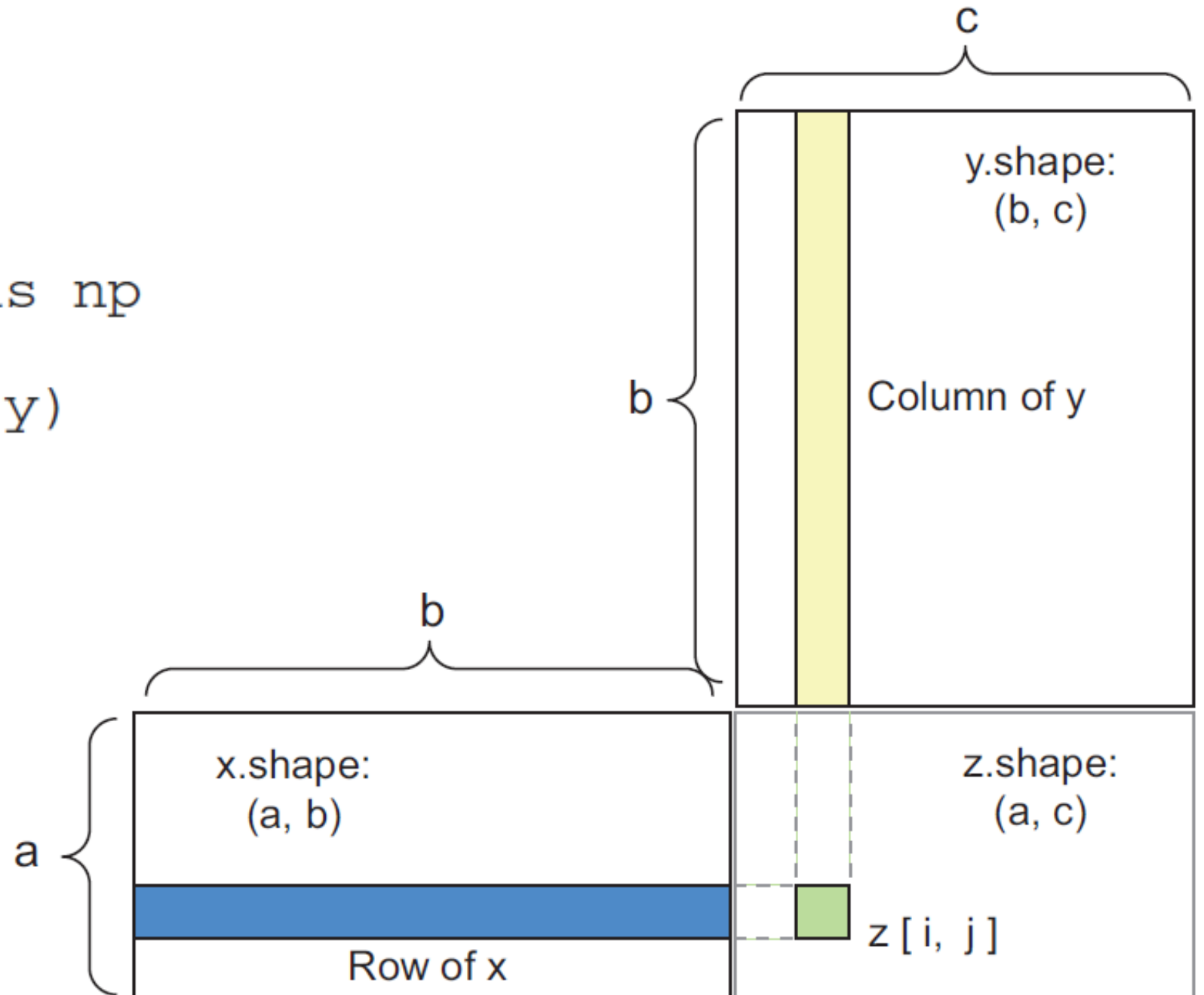
**y is a Numpy vector.**

**Avoid overwriting  
the input tensor.**

# Tensor Dot

```
import numpy as np
```

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



# Implementation of Dot Product

**x and y are Numpy matrices.**

```
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 Ghetana Agnesi

$$(\ln x)' = \frac{1}{x} \quad \int \frac{1}{x} dx = \ln|x| + C$$

$$f(x) = x^2$$
$$\int \sin x dx = -\cos x dx + C$$

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

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

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

# Calculus

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

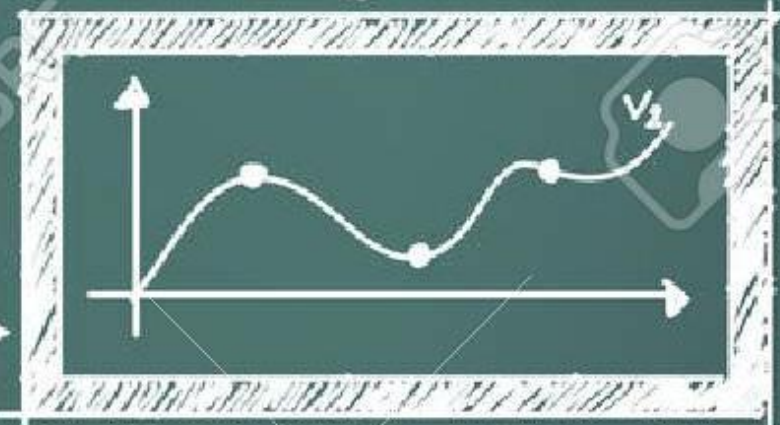
$$4x^2 - 3x - 1 = 0$$

$$\int f(x) dx$$



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

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

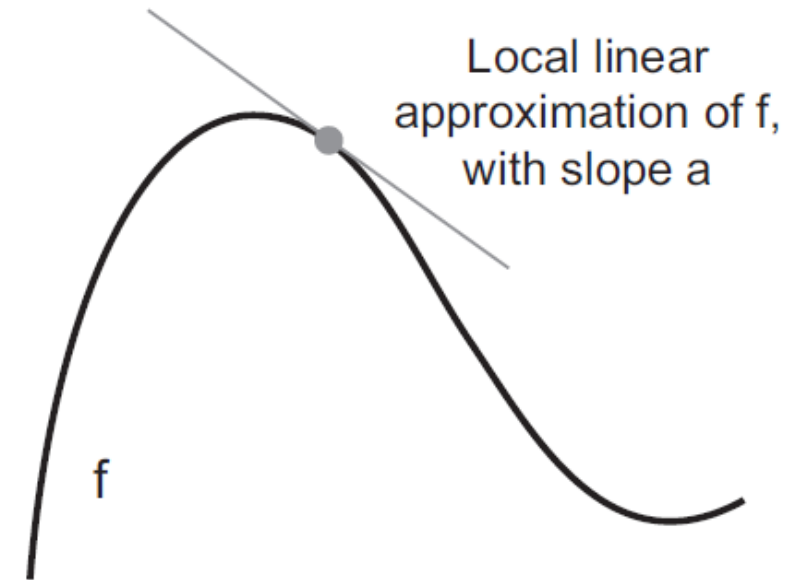


# Differentiation

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

OR

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



# Derivatives of Basic Function $\frac{dy}{dx}$

$$y = x^n \rightarrow \frac{dy}{dx} = nx^{n-1} \quad \frac{d}{dx}\left(\frac{1}{x}\right) \Rightarrow \frac{-1}{x^2}$$

$$y = e^x \rightarrow \frac{dy}{dx} = e^x$$

$$y = \ln x \rightarrow y' = \frac{1}{x}$$

# Gradient of a Function

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

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

- Example

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

$$\nabla f = 1 \times 2y \times 3z^2$$

# Hessian Matrix

- Second-order partial derivatives

$$Hf = \begin{bmatrix} \frac{\delta^2 f}{\delta x^2} & \frac{\delta^2 f}{\delta x \delta y} & \frac{\delta^2 f}{\delta x \delta z} \\ \frac{\delta^2 f}{\delta y \delta x} & \frac{\delta^2 f}{\delta y^2} & \frac{\delta^2 f}{\delta y \delta z} \\ \frac{\delta^2 f}{\delta z \delta x} & \frac{\delta^2 f}{\delta z \delta y} & \frac{\delta^2 f}{\delta z^2} \end{bmatrix}$$

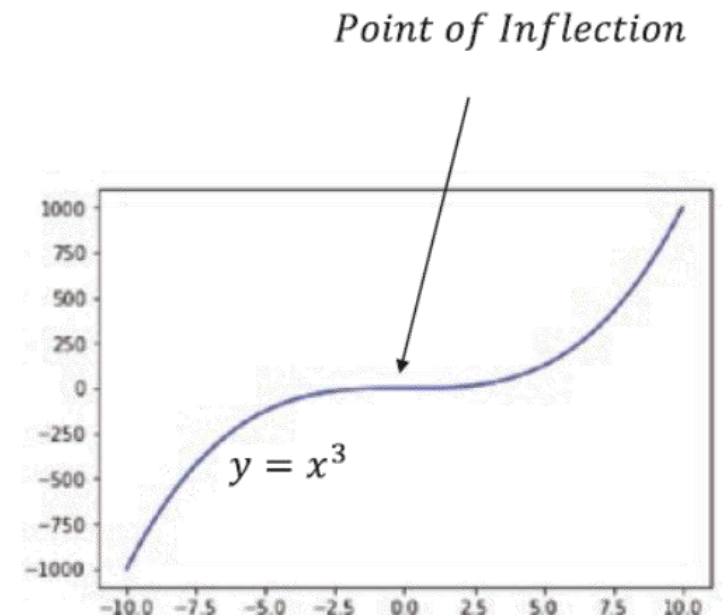
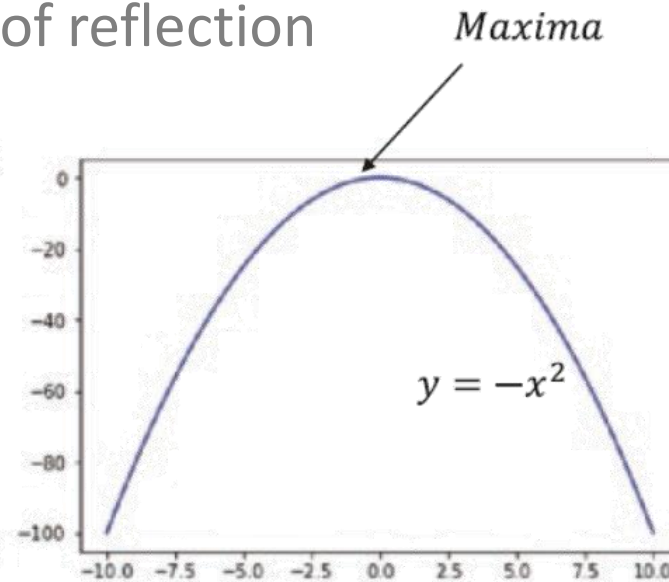
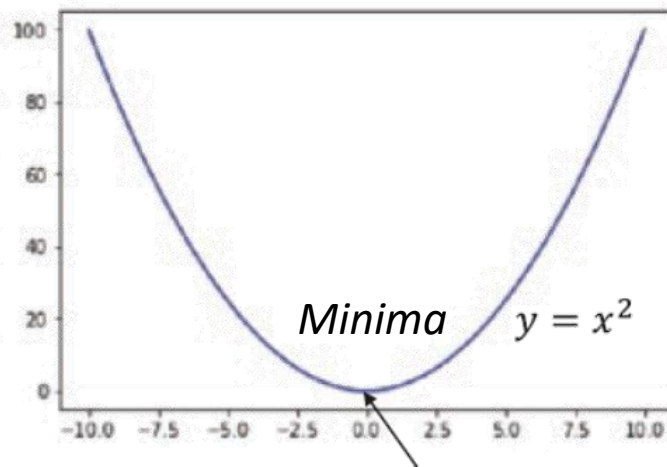
# Maxima and Minima for Univariate Function

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

– If  $\frac{d^2f(x)}{dx^2} < 0 \Rightarrow$  Maxima

– If  $\frac{d^2f(x)}{dx^2} > 0 \Rightarrow$  Minima

– If  $\frac{d^2f(x)}{dx^2} = 0 \Rightarrow$  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

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

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

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

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



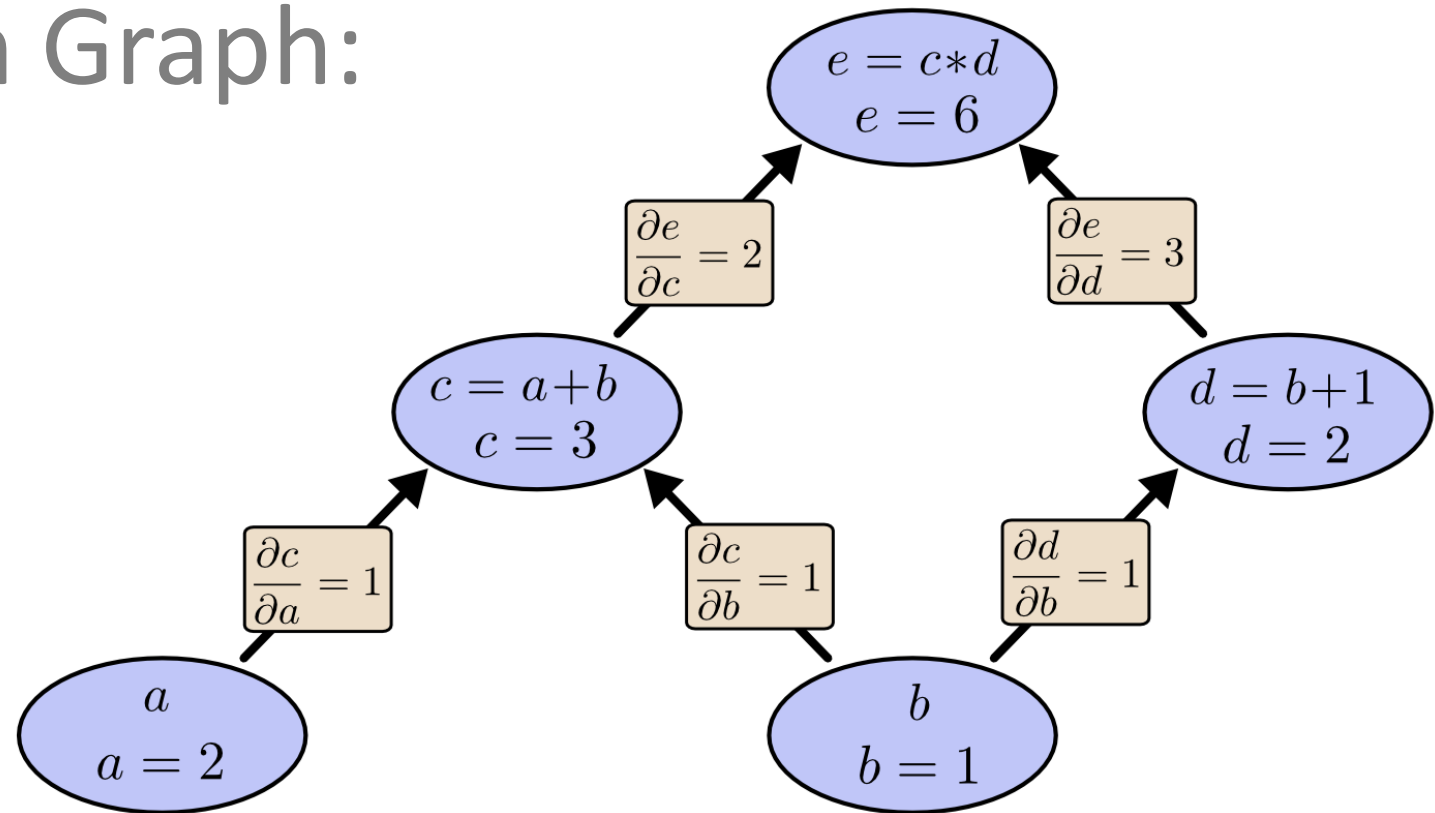
# Symbolic Differentiation

Computation Graph:

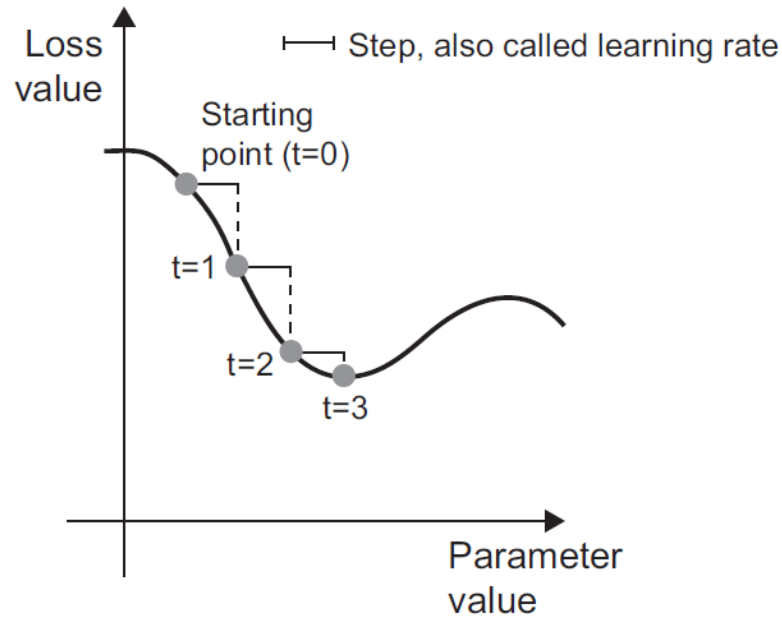
$$c = a + b$$

$$d = b + 1$$

$$e = c * d$$

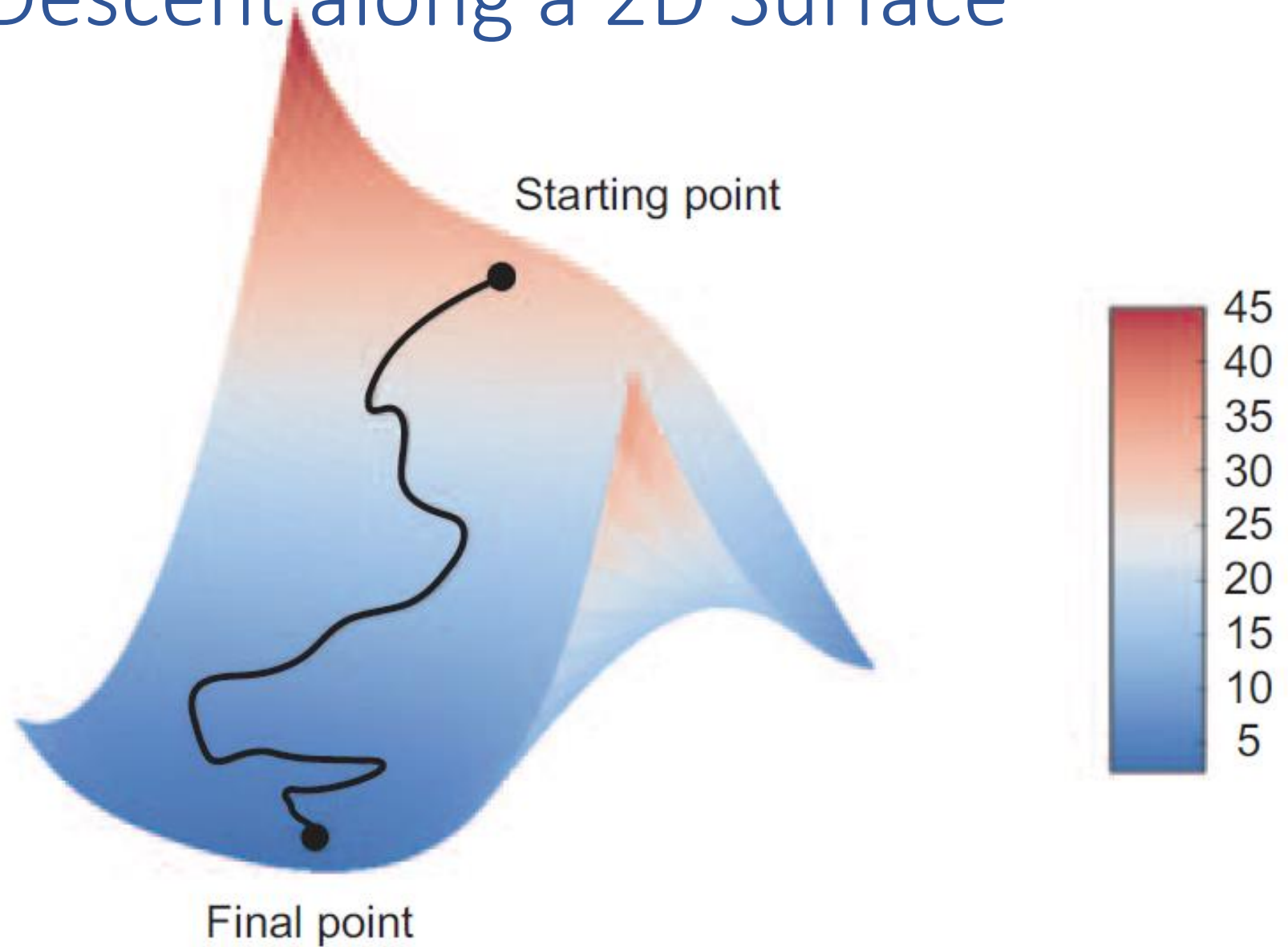


# Stochastic Gradient Descent



1. Draw a batch of training samples  $x$  and corresponding targets  $y$
2. Run the network on  $x$  to obtain predictions  $y_{\text{pred}}$
3. Compute the loss of the network on the batch, a measure of the mismatch between  $y_{\text{pred}}$  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:  $W -= \text{step} * \text{gradient}$

# Gradient Descent along a 2D Surface

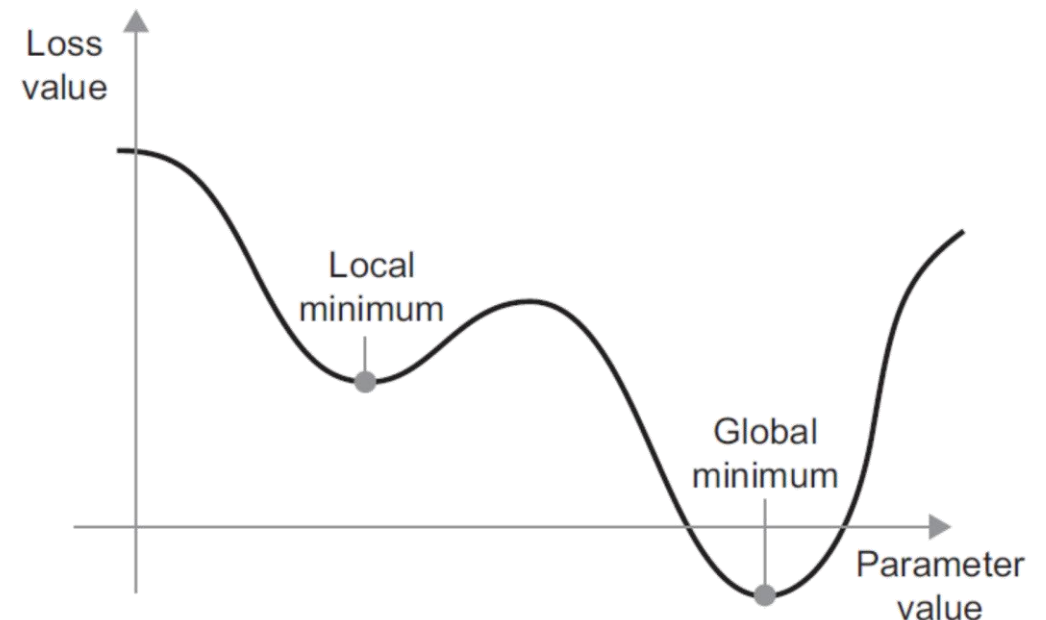


# Avoid Local Minimum using Momentum

```
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

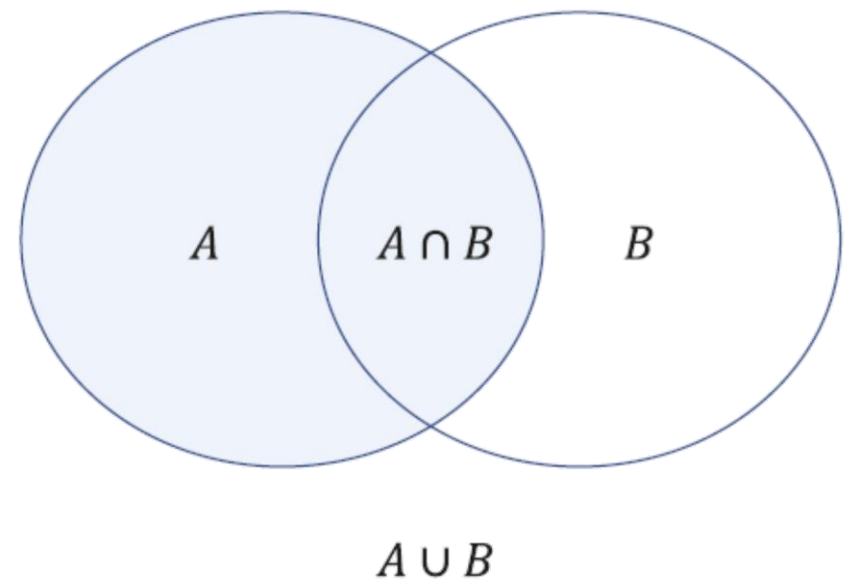
- Given an Event  $E$  in a sample space  $S$ ,  $S = \bigcup_{i=1}^N E_i$
- First axiom
  - $P(E) \in \mathbb{R}, 0 \leq P(E) \leq 1$
- Second axiom
  - $P(S) = 1$
- Third axiom
  - Additivity, any countable sequence of mutually exclusive events  $E_i$
  - $P(\bigcup_{i=1}^n E_i) = P(E_1) + P(E_2) + \cdots + P(E_n) = \sum_{i=1}^n P(E_i)$

# 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

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

$$- P(AB) = P(A|B)P(B) = P(B|A)P(A)$$



# Chain Rule of Probability

- The joint probability can be expressed as chain rule

$$\begin{aligned} P(A_1 A_2 A_3 \dots A_n) &= P(A_1) P(A_2 / A_1) P(A_3 / A_1 A_2) \dots P(A_n / A_1 A_2 \dots A_{(n-1)}) \\ &= P(A_1) \prod_{i=2}^n P(A_i / A_1 A_2 A_3 \dots A_{(n-1)}) \end{aligned}$$



# 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

(Training Data)

feature class (label)

likelihood

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

class

features (Image)

Evidence

prior (class)

Proof:

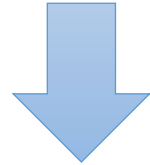
Remember  $P(A|B) = P\left(\frac{AB}{B}\right)$

So  $P(AB) = P(A|B)P(B) = P(B|A)P(A)$

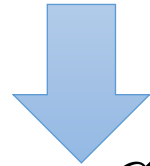
Then Bayes  $P(A|B) = P(B|A)P(A)/P(B)$

# Naïve Bayes Classifier

$$p(C_k \mid \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} \mid C_k)}{p(\mathbf{x})}$$



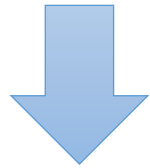
$$p(C_k \mid x_1, \dots, x_n)$$



$$\begin{aligned} p(C_k, x_1, \dots, x_n) &= p(x_1, \dots, x_n, C_k) \\ &= p(x_1 \mid x_2, \dots, x_n, C_k) p(x_2, \dots, x_n, C_k) \\ &= p(x_1 \mid x_2, \dots, x_n, C_k) p(x_2 \mid x_3, \dots, x_n, C_k) p(x_3, \dots, x_n, C_k) \\ &= \dots \\ &= p(x_1 \mid x_2, \dots, x_n, C_k) p(x_2 \mid x_3, \dots, x_n, C_k) \cdots p(x_{n-1} \mid x_n, C_k) p(x_n \mid C_k) p(C_k) \end{aligned}$$

# Naïve = Assume All Features Independent

$$p(x_i \mid x_{i+1}, \dots, x_n, C_k) = p(x_i \mid C_k)$$



$$\begin{aligned} p(C_k \mid x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &= p(C_k) p(x_1 \mid C_k) p(x_2 \mid C_k) p(x_3 \mid C_k) \cdots \\ &= p(C_k) \prod_{i=1}^n p(x_i \mid C_k), \end{aligned}$$

# 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

$$E[X] = x_1p_1 + x_2p_2 + \cdots + x_np_n = \sum_{i=1}^n x_ip_i$$

- Expectation of a continuous random variable

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

# 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

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

where  $\mu_x = E[X]$ ,  $\mu_y = E[Y]$

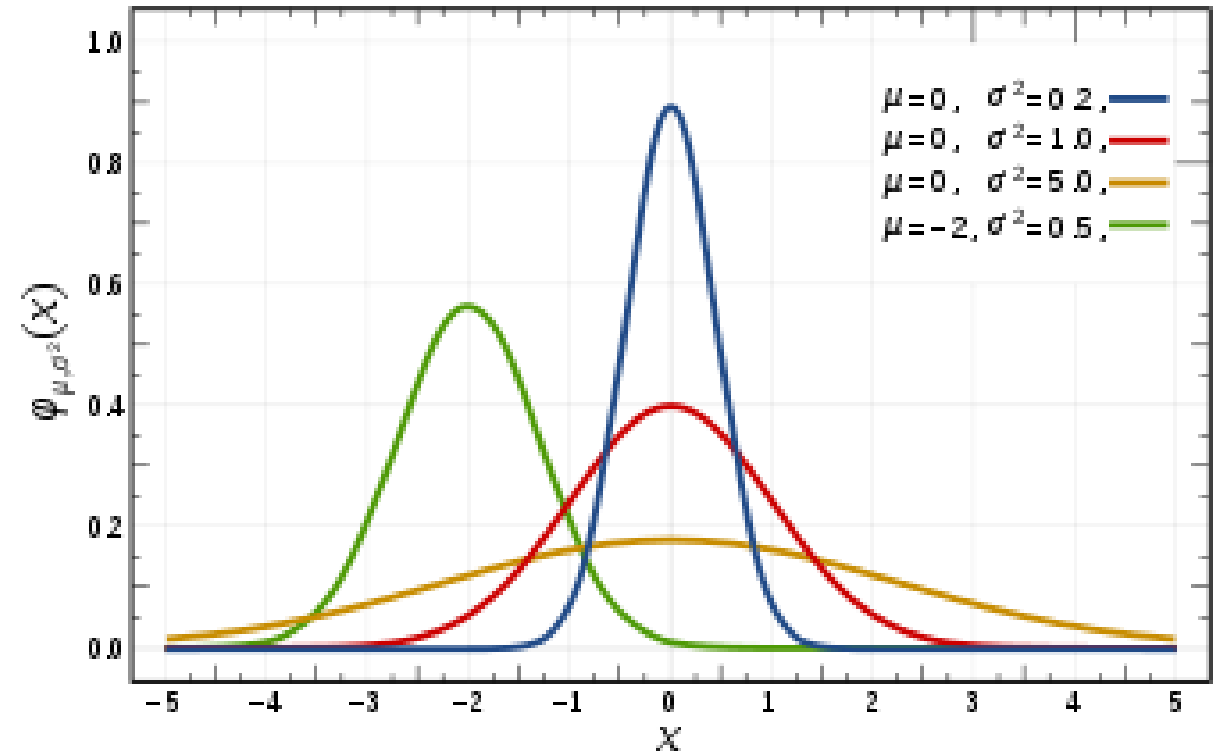
- Correlation coefficient

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

# 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<sup>[1]</sup>

$$\begin{array}{ll}\underset{x}{\text{minimize}} & f(x) \\ \text{subject to} & g_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_j(x) = 0, \quad j = 1, \dots, p\end{array}$$

where

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

# Formulate Your Problem

- Linear model:  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$
- Least-squared Error:  $(f(\mathbf{x}) - \mathbf{y})^2$
- Regularization:  $\|\mathbf{w}\|$
- Objective function:

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

# Principle Component Analysis (PCA)

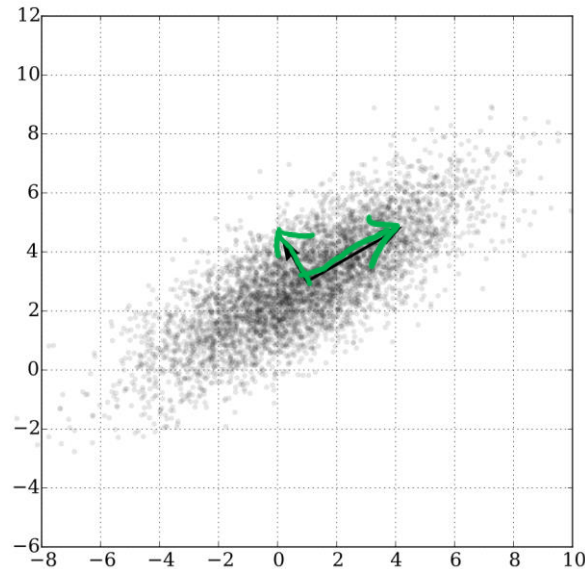
- Assumptions

- *Linearity*
- *Mean and Variance are sufficient statistics*
- *The principal components are orthogonal*

$$y = W^T X$$

$$\max. \text{Cov}(y, y)$$

$$\text{s.t. } W^T W = I$$



# Principle Component Analysis (PCA)

$$\begin{array}{ll} \max_{\mathbf{W}} & \text{cov}(\mathbf{Y}, \mathbf{Y}) \\ \text{s.t.} & \mathbf{W}^T \mathbf{W} = \mathbf{I} \end{array}$$

$$\mathbf{W}^T \mathbf{x} \quad \parallel \quad f = \text{cov}(\mathbf{Y}, \mathbf{Y}) + \lambda(\mathbf{W}^T \mathbf{W} - \mathbf{I})$$

$$\text{Cov}(\mathbf{Y}, \mathbf{Y}) = \frac{1}{N-1} (\mathbf{Y} - \mu_{\mathbf{Y}})^T (\mathbf{Y} - \mu_{\mathbf{Y}}) = \frac{1}{N-1} (\mathbf{W}^T \mathbf{x} - \mathbf{W}^T \mu_{\mathbf{x}})^T (\mathbf{W}^T \mathbf{x} - \mathbf{W}^T \mu_{\mathbf{x}}) = \mathbf{W} \overset{\text{Cov}(\mathbf{x}, \mathbf{x})}{\parallel} \Sigma_{\mathbf{x}} \mathbf{W}^T$$

$$\frac{df}{d\mathbf{W}} = 0 \Rightarrow \frac{d}{d\mathbf{W}} (\mathbf{W} \Sigma_{\mathbf{x}} \mathbf{W}^T + \lambda(\mathbf{W}^T \mathbf{W} - \mathbf{I}))$$

$$\Rightarrow 2\Sigma_{\mathbf{x}} \mathbf{W} + 2\lambda \mathbf{W} = 0 \Rightarrow \Sigma_{\mathbf{x}} \mathbf{W} = \lambda \mathbf{W}$$

# 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