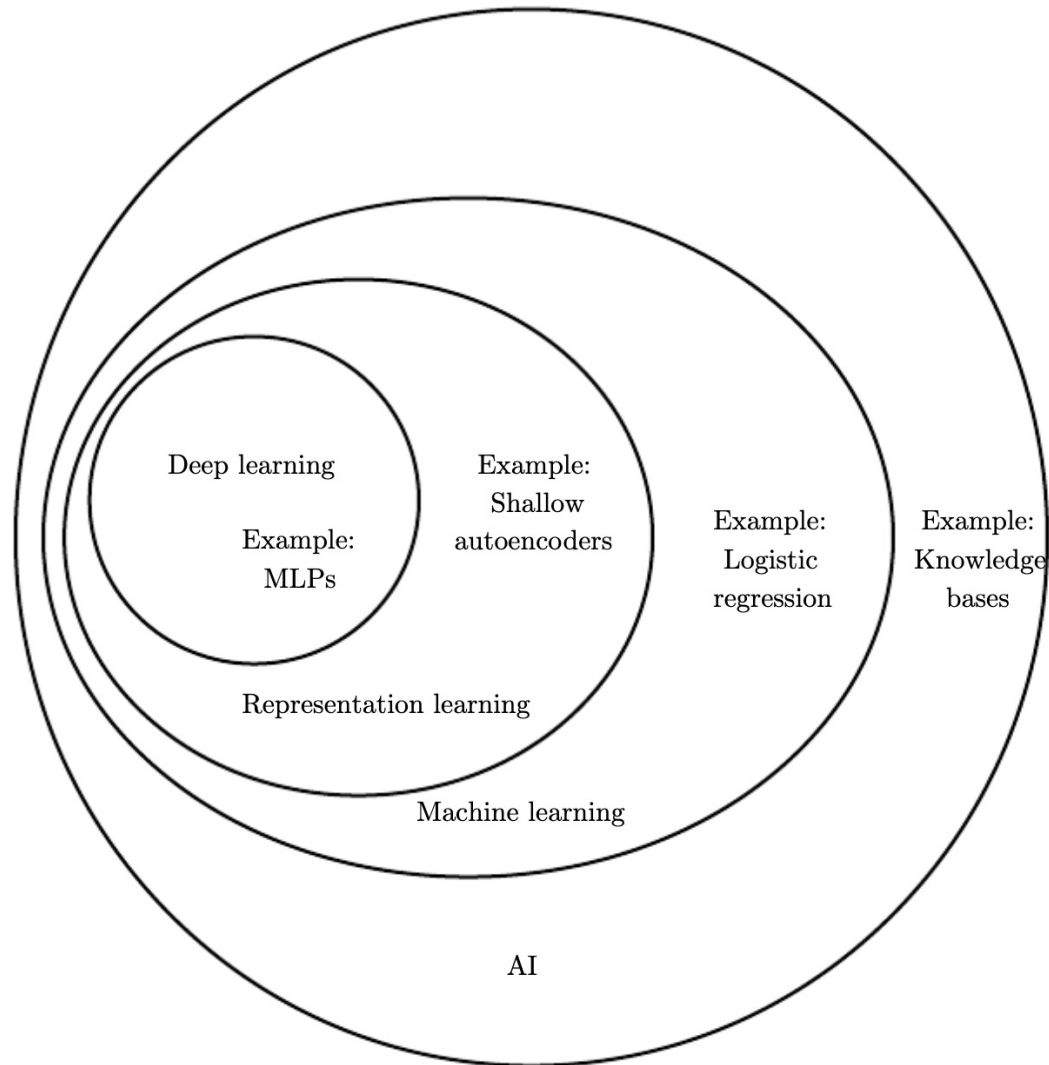




# CSE 176 Introduction to Machine Learning

## Lecture 2: Linear Algebra, Probability and Statistics

# Recap: Different AI systems



# Recap: Major Types of machine learning

- ❑ Supervised learning: Given pairs of input-output, learn to map the input to output
  - ❑ Image classification
  - ❑ Speech recognition
  - ❑ Regression (continuous output)
- ❑ Unsupervised learning: Given unlabeled data, uncover the underlying structure or distribution of the data
  - ❑ Clustering
  - ❑ Dimensionality reduction
- ❑ Reinforcement learning: training an agent to make decisions within an environment to maximize a cumulative reward
  - ❑ Game playing (e.g., AlphaGo)
  - ❑ Robot control



# Linear Algebra

# Linear Algebra Topics

- ❑ Scalars, Vectors, Matrices and Tensors
- ❑ Multiplying Matrices and Vectors
- ❑ Identity and Inverse Matrices
- ❑ Linear Dependence and Span
- ❑ Norms
- ❑ Special kinds of matrices and vectors
- ❑ Eigen decomposition
- ❑ Singular value decomposition

# Scalar, Vector, Matrix, Tensor

□ **Scalar**: A single number (real-valued or integer)

□ **Vector**: An array of numbers arranged in order

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

□ **Matrix**: 2D Array of numbers

$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$$

□ **Tensor**: Sometimes need an array with more than two axes

□ E.g., an RGB color image has three axes

# Types of matrices

$$\begin{bmatrix} 1 & 3 \\ -4 & 7 \end{bmatrix}$$

Square matrix  
 $2 \times 2$

$$\begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix}$$

Rectangular  
matrix  
 $3 \times 2$

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

Zero matrix  
 $3 \times 5$

$$[1 \quad 9 \quad -3 \quad 0]$$

Row matrix  
 $1 \times 4$

$$\begin{bmatrix} 1 \\ 2 \\ 6 \end{bmatrix}$$

Column matrix  
 $3 \times 1$

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

Identity matrix  
 $3 \times 3$

# Matrix times matrix

- If  $A$  is of shape  $m \times n$  and  $B$  is of shape  $n \times p$  then *matrix product*  $C$  is of shape  $m \times p$

$$C = AB \Rightarrow C_{i,j} = \sum_k A_{ik} B_{kj}$$

- Note that the standard product of two matrices is not just the product of two individual elements
  - Such a product does exist and is called the element-wise product or the Hadamard product  $A \odot B$



# Matrix times vector: Linear transformation

- $A\mathbf{x}=\mathbf{b}$

- where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{x}, \mathbf{b} \in \mathbb{R}^n$

- More explicitly

$$\begin{aligned} A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n &= b_1 \\ A_{21}x_1 + A_{22}x_2 + \dots + A_{2n}x_n &= b_2 \\ &\vdots \\ A_{n1}x_1 + A_{n2}x_2 + \dots + A_{nn}x_n &= b_n \end{aligned}$$

$n$  equations in  
 $n$  unknowns

$$\begin{array}{ccc} \mathbf{A} = \begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{n1} & \dots & A_{nn} \end{bmatrix} & \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} & \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \\ n \times n & n \times 1 & n \times 1 \end{array}$$

Can view  $A$  as a *linear transformation*  
of vector  $\mathbf{x}$  to vector  $\mathbf{b}$

- Sometimes we wish to solve for the unknowns  $\mathbf{x} = \{x_1, \dots, x_n\}$  when  $A$  and  $\mathbf{b}$  provide constraints

# Matrix inverse

- Inverse of square matrix  $A$  defined as  $A^{-1}A = I_n$
- We can now solve  $A\mathbf{x} = \mathbf{b}$  as follows:

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

$$A^{-1}A\mathbf{x} = A^{-1}\mathbf{b}$$

$$I_n\mathbf{x} = A^{-1}\mathbf{b}$$

$$\mathbf{x} = A^{-1}\mathbf{b}$$

- This depends on being able to find  $A^{-1}$

# Matrix inverse

□ For a 2x2 matrix:

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

□ The inverse is :

$$A^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

□ Quiz: find the inverse of  $AB^T$

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

□ Answer:

$$AB^T = \begin{bmatrix} 3 & -1 & 1 \\ 2 & 0 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 0 \\ 2 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 6 & 2 \end{bmatrix}$$

$$(AB^T)^{-1} = \frac{1}{5 \cdot 2 - 0 \cdot 6} \cdot \begin{bmatrix} 2 & 0 \\ -6 & 5 \end{bmatrix} = \frac{1}{10} \cdot \begin{bmatrix} 2 & 0 \\ -6 & 5 \end{bmatrix} = \begin{bmatrix} 0.2 & 0 \\ -0.6 & 0.5 \end{bmatrix}$$

# Norms

- Used for measuring the size of a vector
- Norms map vectors to non-negative values
- Norm of vector  $\mathbf{x} = [x_1, \dots, x_n]^T$  is distance from origin to  $\mathbf{x}$ 
  - It is any function  $f$  that satisfies:

$$f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$$

$$f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y}) \quad \text{Triangle Inequality}$$

$$\forall \alpha \in \mathbb{R} \quad f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$$

# $L^p$ Norm

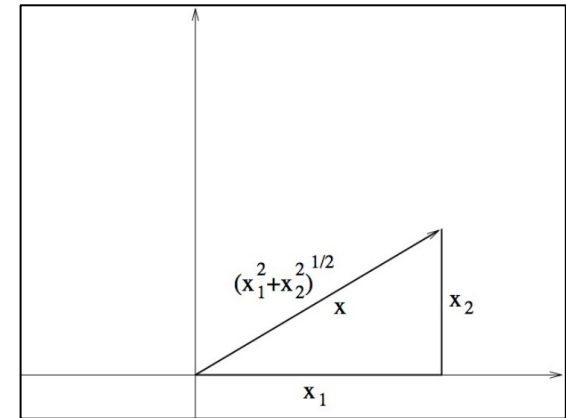
- Definition:

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

- $L^2$  Norm

- Called Euclidean norm

- Simply the Euclidean distance between the origin and the point  $\mathbf{x}$
    - written simply as  $\|\mathbf{x}\|$
    - Squared Euclidean norm is same as  $\mathbf{x}^T \mathbf{x}$



- $L^1$  Norm

- Sum of absolute value for each  $x_i$

- $L^\infty$  Norm

$$\|\mathbf{x}\|_\infty = \max_i |x_i|$$

- Called max norm

# Special kind of vectors

- Unit Vector

- A vector with unit norm

- Orthogonal Vectors

- A vector  $\mathbf{x}$  and a vector  $\mathbf{y}$  are orthogonal to each other if  $\mathbf{x}^T \mathbf{y} = 0$

- Orthonormal Vectors

- Vectors are orthogonal & have unit norm

- Orthogonal Matrix

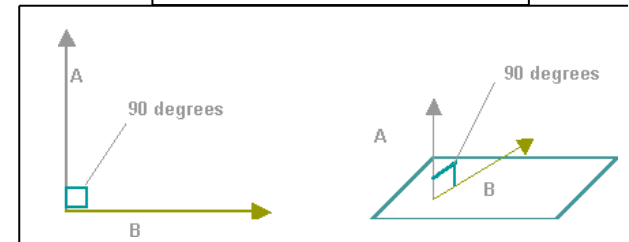
- A square matrix whose rows are mutually

- orthonormal:  $A^T A = A A^T = I$

- $A^{-1} = A^T$

$$\|\mathbf{x}\|_2 = 1$$

$$\begin{bmatrix} 2 \\ -3 \\ -2 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 4 \\ -2 \\ 7 \end{bmatrix}$$



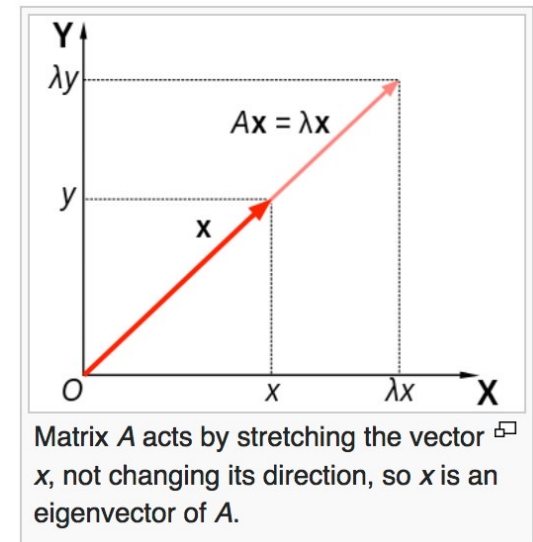
# Eigenvector

- An eigenvector of a square matrix **A** is a non-zero vector **v** such that multiplication by **A** only changes the scale of **v**

$$A\mathbf{v} = \lambda\mathbf{v}$$

– The scalar  $\lambda$  is known as eigenvalue

- If **v** is an eigenvector of **A**, so is any rescaled vector **sv**. Moreover **sv** still has the same eigen value. Thus look for a unit eigenvector



Wikipedia

# Eigendecomposition

- Suppose that matrix  $A$  has  $n$  linearly independent eigenvectors  $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$  with eigenvalues  $\{\lambda_1, \dots, \lambda_n\}$
- Concatenate eigenvectors to form matrix  $V$
- Concatenate eigenvalues to form vector  $\lambda = [\lambda_1, \dots, \lambda_n]$
- Eigendecomposition of  $A$  is given by

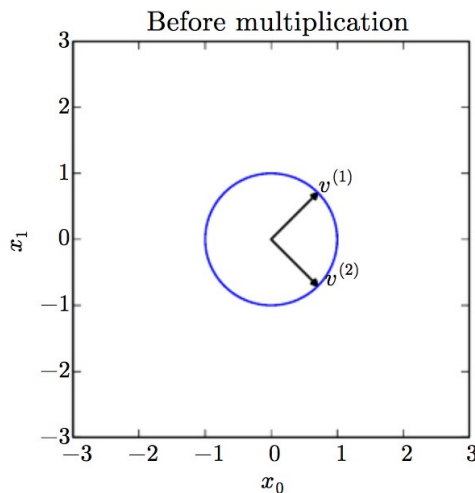
$$A = V \text{diag}(\lambda) V^{-1}$$



# Effect of eigenvalue and eigenvector

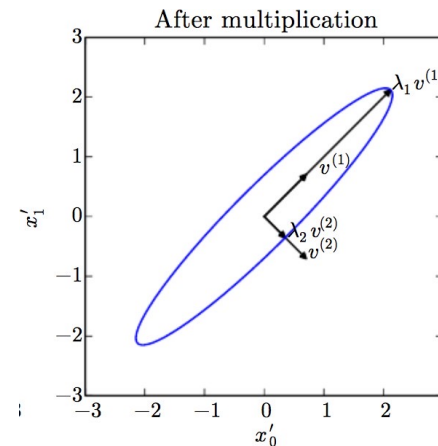
- Example of  $2 \times 2$  matrix
- Matrix  $A$  with two orthonormal eigenvectors
  - –  $v^{(1)}$  with eigenvalue  $\lambda_1$ ,  $v^{(2)}$  with eigenvalue  $\lambda_2$

Plot of unit vectors  $u \in U^2$   
(circle)



with two variables  $x_1$  and  $x_2$

Plot of vectors  $Au$   
(ellipse)



# Positive Semidefinite Matrix (PSD)

- A matrix whose eigenvalues are all positive is called *positive definite*
  - Positive or zero is called *positive semidefinite*
- If eigen values are all negative it is *negative definite*
  - Positive definite matrices guarantee that  $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$

# Singular Value Decomposition (SVD)

- Eigendecomposition has form:  $A = V \text{diag}(\lambda) V^{-1}$ 
  - If  $A$  is not square, eigendecomposition is undefined
- SVD is a decomposition of the form  $A = U D V^T$
- SVD is more general than eigendecomposition
  - Used with any matrix rather than symmetric ones
  - Every real matrix has a SVD
    - Same is not true of eigen decomposition



# Probability and Statistics

# Random Variable

- Variable that can take different values randomly
- Scalar random variable denoted  $x$
- Vector random variable is denoted in bold as  $\mathbf{x}$
- Values of r.v.s denoted in italics  $x$  or  $\mathbf{x}$ 
  - Values denoted as  $\text{Val}(\mathbf{x}) = \{x_1, x_2\}$
- Random variable must have a probability distribution to specify how likely the states are
- Random variables can be discrete or continuous
  - Discrete values need not be integers, can be named states
  - Continuous random variable is associated with a real value

# Probability Distribution

- ❑ A probability distribution is a description of how likely a random variable or a set of random variables is to take each of its possible states
- ❑ The way to describe the distribution depends on whether it is discrete or continuous

# Continuous Variables and PDFs

- When working with continuous variables, we describe probability distributions using probability density functions
- To be a pdf  $p$  must satisfy:
  - The domain of  $p$  must be the set of all possible states of  $x$ .
  - $\forall x \in \mathbf{x}, p(x) \geq 0$ . Note that we do not require  $p(x) \leq 1$ .
  - $\int p(x)dx = 1$ .

# Marginal distribution

- Sometimes we know the joint distribution of several variables
- And we want to know the distribution over some of them
- It can be computed using

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

$$p(x) = \int p(x, y) dy$$



# Conditional probability

- We are often interested in the probability of an event given that some other event has happened
- This is called conditional probability
- It can be computed using

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

# Chain rule of conditional probability

- Any probability distribution over many variables can be decomposed into conditional distributions over only one variable

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

- An example with three variables

$$P(a, b, c) = P(a \mid b, c)P(b, c)$$

$$P(b, c) = P(b \mid c)P(c)$$

$$P(a, b, c) = P(a \mid b, c)P(b \mid c)P(c)$$

# Independence and conditional independence

- Independence:  $x \perp y$

- Two variables  $x$  and  $y$  are independent if their probability distribution can be expressed as a product of two factors, one involving only  $x$  and the other involving only  $y$

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

- Conditional Independence:  $x \perp y \mid z$

- Two variables  $x$  and  $y$  are independent given variable  $z$ , if the conditional probability distribution over  $x$  and  $y$  factorizes in this way for every  $z$

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

# Common probability distribution

- Several simple probability distributions are useful in many contexts in machine learning
  - Bernoulli over a single binary random variable
  - Multinoulli distribution over a variable with  $k$  states
  - Gaussian distribution
  - Mixture distribution

# Mixture of Distribution

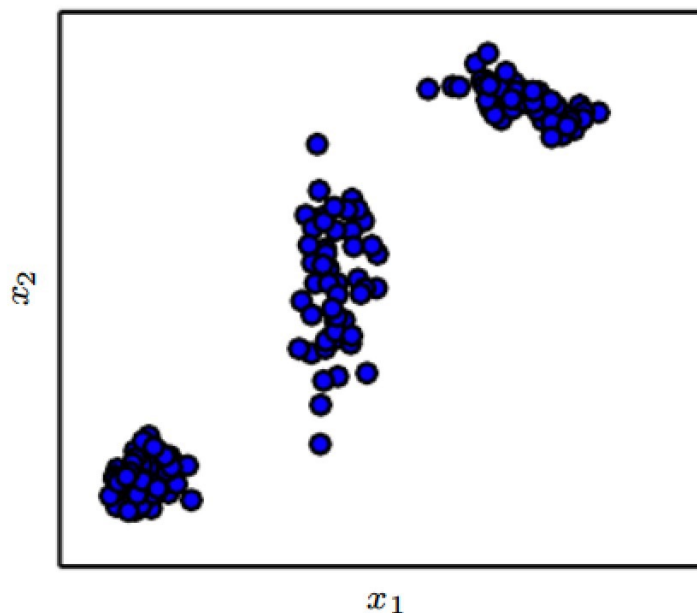
- A mixture distribution is made up of several component distributions
- On each trial, the choice of which component distribution generates the sample is determined by sampling a component identity from a multinoulli distribution:

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

– where  $P(c)$  is a multinoulli distribution

# Gaussian mixture model

- Components  $p(\mathbf{x}|\mathbf{c}=i)$  are Gaussian
- Each component has a separately parameterized mean  $\mu^{(i)}$  and covariance  $\Sigma^{(i)}$
- Any smooth density can be approximated with enough components
- Samples from a GMM:
  - 3 components



# Quiz

□ A random variable,  $X$ , has the probability distribution table as shown.

$x$	-2	-1	0	1	2
$P(X = x)$			0.4	0.1	0.1

Assume that  $P(X = -2) = P(X = -1)$ . Compute the expectation and variance of  $X$ .

# Bayes's rule

- **Bayes' theorem** (alternatively **Bayes' law** or **Bayes' rule**), named after [Thomas Bayes](#), describes the [probability](#) of an [event](#), based on prior knowledge of conditions that might be related to the event.
- For example, if the risk of health problems is known to increase with age, Bayes' theorem allows the risk to an individual of a known age to be assessed more accurately by conditioning it relative to their age.

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



# Quiz

Suppose that  $P(A \cap B) = 0.4$  and  $P(B) = 0.9$ . Find  $P(A|B)$ .

*Solution:*  $\frac{4}{9}$

From the definition of conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.4}{0.9} = \frac{4}{9} = 0.\bar{4}$$

# Quiz

- A motor insurance company insures drivers in age group A, B and C. 40% of the customers are in group A, 25% are in B, and 35% are in group C. The company's record shows that each year, 2% of customers in age group A, 1% in group B and 1.5% in group C made a claim. Given that a driver made a claim, what is the probability that the driver is from age group C?