# DDA4210/AIR6002 Advanced Machine Learning
## Lecture 01 Introduction and Review

Tongxin Li

School of Data Science, CUHK-Shenzhen

Spring, 2024

# Overview

1. About this course

2. Review for basic machine learning methods

# Logistics

- Instructor: Tongxin Li
  - Email: litongxin@cuhk.edu.cn
  - Personal website: https://tongxin.me/
  - Office: Daoyuan Building 323A
  - Office hours: Thu 9:50-10:50 am

## Logistics

- Instructor: Tongxin Li
  - Email: litongxin@cuhk.edu.cn
  - Personal website: https://tongxin.me/
  - Office: Daoyuan Building 323A
  - Office hours: Thu 9:50-10:50 am



SCAN ME

- Course website: https://tongxin.me/DDA4210/

# Logistics

Teaching assistants:

- Fangchen Yu, 220019040@link.cuhk.edu.cn
- Yifei Wu, 223040255@link.cuhk.edu.cn
- Fanzeng Xia, 223040232@link.cuhk.edu.cn

USTFs:

- Huihan Yang, 120090438@link.cuhk.edu.cn
- Yiming Xiong, 120090721@link.cuhk.edu.cn
- Rongxiao Qu, 120020144@link.cuhk.edu.cn
- Jiayi Yao, 120040070@link.cuhk.edu.cn

# Assessment (DDA4210)

- Homework (30%)
  - Three assignments (tri-weekly)
  - Involves theory, analysis, computation, and programming.

# Assessment (DDA4210)

- Homework (30%)
    - Three assignments (tri-weekly)
    - Involves theory, analysis, computation, and programming.
- Course project (35%)
    - Format: Python programming for advanced machine learning
    - Topic: determined by yourself (given a few examples or choices)
    - Teamwork: 1 to 4 members per team
    - Outcome evaluation: [report(25%)+presentation(75%)]$\times r_{ind}$
        - presentation: $75\% = 10\%$peer $+ 25\%$TA $+ 40\%$instructor
        - $r_{ind} \in \{0.5, 0.8, 1\}$: it is rated by your teammates on your contribution; 0.5 (or 0.8) means your contribution is less than 20% (or 50%) of the expected workload.
    - Evaluation criteria: significance, novelty, technical soundness, completeness

# Assessment (DDA4210)

- Homework (30%)
  - Three assignments (tri-weekly)
  - Involves theory, analysis, computation, and programming.
- Course project (35%)
  - Format: Python programming for advanced machine learning
  - Topic: determined by yourself (given a few examples or choices)
  - Teamwork: 1 to 4 members per team
  - Outcome evaluation: [report(25%)+presentation(75%)]$\times r_{ind}$
    - presentation: $75\% = 10\%$peer $+ 25\%$TA $+ 40\%$instructor
    - $r_{ind} \in \{0.5, 0.8, 1\}$: it is rated by your teammates on your contribution; 0.5 (or 0.8) means your contribution is less than 20% (or 50%) of the expected workload.
  - Evaluation criteria: significance, novelty, technical soundness, completeness
- Final exam (35%)
  - Single-choice questions, calculation, math derivation/proofs, etc.

# Assessment (AIR6002)

- Homework (40%)
    - Three assignments (tri-weekly)
    - Involves theory, analysis, computation, and **more** programming.

# Assessment (AIR6002)

- Homework (40%)
    - Three assignments (tri-weekly)
    - Involves theory, analysis, computation, and **more** programming.
- Course project (60%)
    - Format: Cutting-edge topics in advanced machine learning
    - Topic: determined by yourself (given a few examples or choices)
    - Teamwork: 1 to 3 members per team
    - Outcome evaluation: [mid-term proposal(10%)+ report(25%)+presentation(65%)]$\times r_{ind}$
        - presentation: $75\% = 10\%$peer $+ 25\%$TA $+ 40\%$instructor
        - $r_{ind} \in \{0.5, 0.8, 1\}$: it is rated by your teammates on your contribution; 0.5 (or 0.8) means your contribution is less than 20% (or 50%) of the expected workload.
    - Evaluation criteria: significance, novelty, technical soundness, completeness

# Some remarks

- Plagiarism violates the university policy of "**Academic Integrity**"
  - Plagiarism in homework assignments, course projects, and final exam will be dealt with **severity**.
  - For example, assignments with plagiarism will be graded as zero.
  - Repeated plagiarism will lead to an "F" for the entire course.
- Attendance requirement
  - Attending lectures/tutorials onsite is highly encouraged.
  - Please answer or raise questions actively.
  - Let the instructor/TAs be able to recognize you as a student in the class.
- Participation in Course&Teaching Evaluation (CTE)
  - Your feedback (either positive or negative) helps improve the course and make it even better.

# Some remarks

*The building blocks of machine learning are **data**, **models**, and **algorithms**.*

*The building blocks of advanced machine learning (DDA4210/AIR6002) are more complicated data, more powerful models, and state-of-the-art algorithms in real-world applications.*

# Syllabus

1. Review of basic machine learning methods
2. Advanced ensemble learning
3. Learning theory
4. Advanced applications: recommendation and search
5. Spectral clustering and semi-supervised learning
6. Graph neural networks
7. Nonlinear dimensionality reduction and data visualization
8. Generative models (VAE, GAN, diffusion model)
9. Causal machine learning
10. Privacy in machine learning
11. Fairness in machine learning
12. Interpretability in machine learning
13. Course project presentation
14. Final exam

# Basic machine learning methods

- Linear regression and classification
- K-nearest neighbor method
- Decision tree, bagging, and random forest
- Support vector machine
- Neural networks (MLP, CNN, and RNN)
- K-means and Gaussian mixture models
- Principal component analysis

# Advanced Machine Learning: Boosting



*Boosting is an ensemble meta-algorithm for primarily reducing bias, and also variance in supervised learning, and a family of machine learning algorithms that convert weak learners to strong ones.*—-Wikipedia

# Advanced Machine Learning: Learning Theory

Machine Learning Theory[1]

- Also known as *Computational Learning Theory*
- Aims to understand the fundamental principles of learning as a computational process and combines tools from Computer Science and Statistics
  - Creating mathematical models that capture key aspects of machine learning, in which one can analyze the inherent ease or difficulty of different types of learning problems.
  - Proving guarantees for algorithms (under what conditions will they succeed, how much data and computation time is needed) and developing machine learning algorithms that provably meet desired criteria.
  - Mathematically analyzing general issues, such as: "When can one be confident about predictions made from limited data?", "What kinds of methods can learn even in the presence of large quantities of distracting information?"

---

[1] https://www.cs.cmu.edu/~avrim/Talks/mlt.pdf

- Collaborative filtering methods
- Content-based methods
- Hybrid methods

K-means

Spectral clustering

two circles, 2 clusters (K−means)

twocircles, 2 clusters

• what if you don't have linearly separated datapts ?

## Why Semi-Supervised Learning?

Classification on the two moons pattern [Zhou et al. 04]:

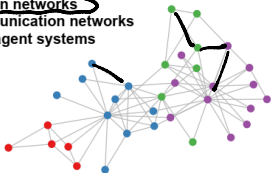(a) two labeled points; (b) SVM with an RBF kernel; (c) k-NN with k = 1.

Graph-Structured data cannot be well handled by conventional neural networks!

GNN

$$X = \{x_i\}_{i=1}^{n}$$

$$G = (V, E)$$

A lot of real-world data does not "live" on grids

$$\{x_1, \dots, x_n\}$$

**Social networks**
**Citation networks**
**Communication networks**
**Multi-agent systems**



Knowledge graphs

:country
U.S.A.

citizen_of

Mikhail Baryshnikov — educated_at → Vaganova Academy

:ballet_dancer
awarded

:university

:award
Jiicek prize

Molecules

Protein interaction networks

**Standard deep learning architectures like CNNs and RNNs don't work here!**

Road maps

{ • Classification
  • prediction .

# Advanced Machine Learning: Graph Neural Networks

Graph-Structured data cannot be well handled by conventional neural networks!



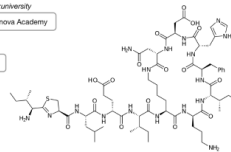A lot of real-world data does not "live" on grids

Social networks
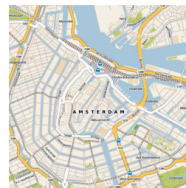Citation networks
Communication networks
Multi-agent systems

Knowledge graphs

Molecules

Protein interaction networks

**Standard deep learning architectures like CNNs and RNNs don't work here!**

Road maps

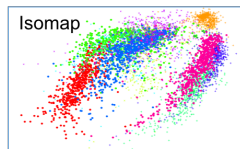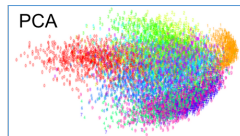Tasks: node classification, link prediction, graph classification
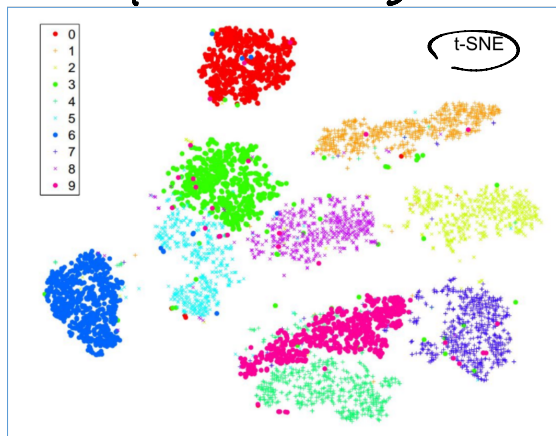
The image is from Thomas Kipf.

NLDR: nonlinear Dimensionality reduction
Example: visualizing MNIST handwritten digits (10 classes)

$$\{0, 1, 2, 3 \cdots 9\}$$

*VAE, GAN. Diffusion model.*

Use the model trained on training data to generate new data, such as images, text, audio, and videos.

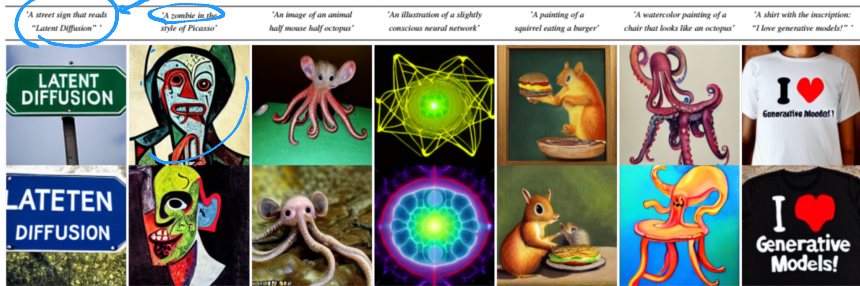Generating images [Rombach et al. 2022]

Use the model trained on training data to generate new data, such as images, text, audio, and videos.

Text to image [Rombach et al. 2022]



Text-to-Image Synthesis on LAION. 1.45B Model.

# First PART.

1. Ensemble learning.
2. Recommendation system.
3. Clustering.
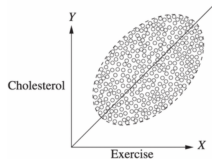4. Nonlinear DR.
5. GNN
6. Generative Models.
7. Learning theory

- Develop ML models. and algorithms. that are ~~accurate~~.
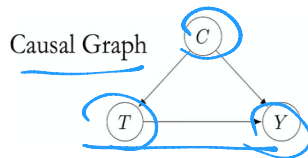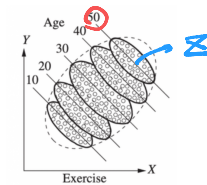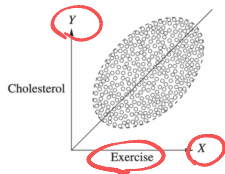
# Second PART.

Next, besides accuracy.

- explainable.
- fair
- privacy - preserving.
- safe / robust.

- Causal inference

- Causal inference

- Causal inference



Causal Graph

- Causal machine learning (for non-i.i.d problem)

$X = (x_1, \cdots x_n)$

$x_i \sim P$ iid.



Train:

"0"    "1"

Test:

"0" (misleading to "1")

- Causal inference



- Causal machine learning (for non-i.i.d problem)



- Causal discovery

Schematic overview of the relationships and interactions between data, algorithms, actors and techniques in the field of secure and private AI [Kaissis et al. 2020]:

# Advanced Machine Learning: Fairness

Where does the unfairness in machine learning algorithms come from? How can we address the unfairness?

Examples of how bias in machine learning can affect our daily lives [Grabski et al. 2020]:

Understanding the reasons behind decisions made by black-box machine learning models

Explaining individual outputs of a model that predicts that a patient has the flu [Ribeiro et al. 2016]:

Understanding the reasons behind decisions made by black-box machine learning models

Explaining individual outputs of a model that predicts that a patient has the flu [Ribeiro et al. 2016]:



Explaining an image classification prediction made by Google's Inception neural network [Ribeiro et al. 2016]:



(a) Original Image    (b) Explaining *Electric guitar*    (c) Explaining *Acoustic guitar*    (d) Explaining *Labrador*

# Review for basic machine learning methods

- Linear regression and classification
- K-nearest neighbor method
- Decision tree, bagging, and random forest
- Support vector machine
- Neural networks (MLP, CNN, and RNN)

Supervised Methods

- K-means and Gaussian mixture models
- Principal component analysis

unsupervised.

- Training data: $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \ldots (\mathbf{x}_N, \mathbf{y}_N)\}$
  - $\mathbf{x}_i \in \mathbb{R}^D$, $\mathbf{y}_i \in \mathbb{R}^K$, $i = 1, 2, \ldots, N$
  - with i.i.d assumption usually

$N \gg D$

$N \gg K$

- Learn a linear function $f_{\mathbf{W},\mathbf{b}}(\mathbf{x}) = \underbrace{\mathbf{W}\mathbf{x} + \mathbf{b}}_{\text{affine.}}$ from $\mathcal{D}$
  - $\mathbf{W} \in \mathbb{R}^{K \times D}$, $\mathbf{b} \in \mathbb{R}^K$

- Linear regression (least squares)

(i) $\sum_{i=1}^{N} \|y_i - \bar{W}\bar{x}_i\|^2 = \|Y - \bar{W}\bar{X}\|_F^2$

$\begin{cases} x_i \in \mathbb{R}^D \\ y_i \in \mathbb{R}^K \\ N := \# \text{ of samples.} \end{cases}$

putting $\{y_i\}, \{x_i\}$ here

$\bar{W} := [b \quad W]$

$\bar{x}_i := \begin{bmatrix} 1 \\ x_i \end{bmatrix}$

$$\min_{\mathbf{W}, \mathbf{b}} \sum_{i=1}^{N} \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2 \qquad (1)$$

2.

Q: • Space complexity?

$O(N(K+D))$ $(N \gg D)$.

- Ridge regression

$$\min_{\mathbf{W}, \mathbf{b}} \frac{1}{2} \sum_{i=1}^{N} \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2 + \frac{\lambda}{2} \|\mathbf{W}\|_F^2 \qquad (2)$$

Avoid overfitting

regularizer.

- LASSO

$\bar{W} = Y\bar{X}^T(\bar{X}\bar{X}^T)^{-1}$

"normal equation"

$$\min_{\mathbf{W}, \mathbf{b}} \frac{1}{2} \sum_{i=1}^{N} \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i - \mathbf{b}\|^2 + \lambda\|\mathbf{W}\|_1 \qquad (3)$$

enforces sparsity.

Q: $\ell_0$ norm?

\* $\|\mathbf{W}\|_F = \sqrt{\sum_i \sum_j w_{ij}^2}, \quad \|\mathbf{W}\|_1 = \sum_i \sum_j |w_{ij}|$

# Review: Linear Classification

- Training data: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots (\mathbf{x}_N, y_N)\}$
  - $\mathbf{x}_i \in \mathbb{R}^D$, $y_i \in \{+1, -1\}$, $i = 1, 2, \dots, N$
  - with i.i.d assumption usually

- Learn a linear classifier $f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ from $\mathcal{D}$

# Review: Linear Classification

*(handwritten: Q: Decision Boundary of LR : Linear, Non-linear.)*

- Logistic regression (binary classification, $y \in \{0, 1\}$)

*(handwritten annotations around the sigmoid plot: $1$, $\sigma(y)$, $0.5$, $0$, $y$)*

$$f_{\mathbf{w},b}(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + b) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x} - b)} \qquad (4)$$

*(handwritten: $\sigma(y)$; probability $\sigma(w^\top x + b) \to y \in \{1\}$; probability $1 - \sigma(w^\top x + b) \to y \in \{0\}$  affine; 2D example of DB)*

$$\min_{\mathbf{w},b} \ -\frac{1}{N} \sum_{i=1}^{N} \left( y_i \log f_{\mathbf{w},b}(\mathbf{x}_i) + (1 - y_i) \log(1 - f_{\mathbf{w},b}(\mathbf{x}_i)) \right) \qquad (5)$$

*(handwritten: cross-entropy)*

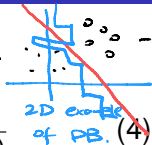- Softmax regression (multi-class classification, $\mathbf{y} \in \{0, 1\}^K$)

$$f_{\mathbf{W},\mathbf{b}}^{(j)}(\mathbf{x}) = \frac{\exp(\mathbf{w}_j^\top \mathbf{x} + b_j)}{\sum_{c=1}^{K} \exp(\mathbf{w}_c^\top \mathbf{x} + b_c)} \qquad (6)$$

$$\min_{\mathbf{w},\mathbf{b}} \ -\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{K} y_{ij} \log f_{\mathbf{w},b}^{(j)}(\mathbf{x}_i) \qquad (7)$$

- k-NN: a nonlinear regression or classification model
- Determine the following beforehand
  - distance metric ($\ell_2$ or $\ell_1$ norms, etc)
  - number ($k$) of nearest neighbors
- k-NN is a non-parametric model

Advantage: Non-parametric.

Disadvantage: In the (vanilla version) of k-NN. need to compare Xnew w/ all. existing data pts



Figure: A toy example (k=3)

Idea: mimicking the process in human-based decision-making.

Advantage: Natural deal-with general types of data



**Is Person Fit or Unfit?**

Age<30

Yes → Eat pizza?

No → Exercise

Eat pizza? → Yes → Unfit, No → Fit

Exercise → Yes → Fit, No → Unfit

Disadvantage. (1) overfitting if depth is high

(2). Variance can be high !

Person 1: age=20, eat pizza, exercise

Person 2: age=40, no eat pizza, exercise

reduce it !

Q: Decison boundary    Linear ?

Non linear ? ✓

Goal: Reduce variance of DT models (1995)
↳ ensemble learning

Dataset

$\mathbb{E}[X_i]$  $Var[X_i]$

$\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}X_i\right]$  $Var\left[\frac{1}{n}\sum_{i=1}^{n}X_i\right]$  $n \to \infty$

**Decision Tree-1**  **Decision Tree-2**  **Decision Tree-N**

**Result-1**  **Result-2**  **Result-N**

Majority Voting / Averaging

**Final Result**

(1964.)



- Margin width:

$$M = (\mathbf{x}^+ - \mathbf{x}^-) \cdot \mathbf{n}$$

$$= (\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$

- Maximum margin classifier

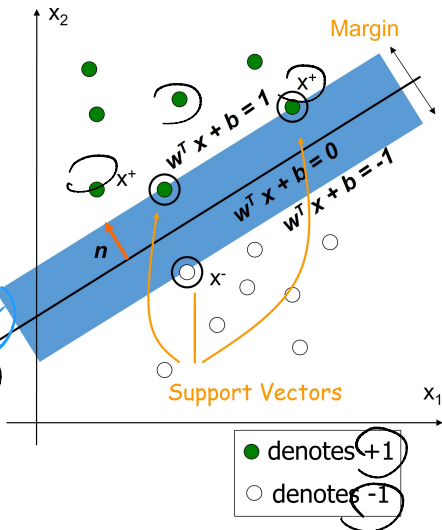(only works well when data is linearly separable)

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \qquad (8)$$

$$\text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \forall i$$

Q: DB        Linear    ✓

          Nonlinear  ✗

In figure: $x_2$, Margin, $x^+$, $w^T x + b = 1$, $x^+$, $w^T x + b = 0$, $w^T x + b = -1$, $n$, $x^-$, Support Vectors, $x_1$

- ● denotes +1
- ○ denotes -1

If data is not linearly separable.

- Margin width:

$$M = (\mathbf{x}^+ - \mathbf{x}^-) \cdot \mathbf{n}$$

$$= (\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$
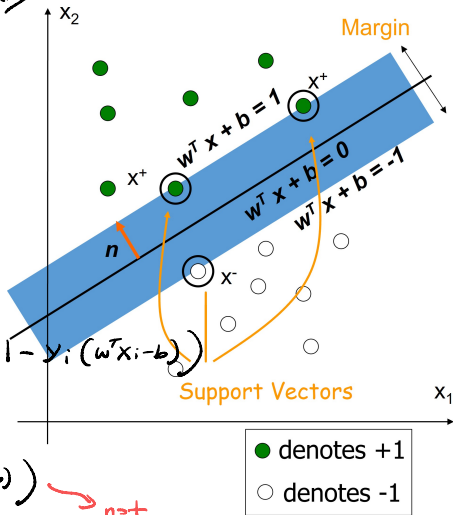
- Maximum margin classifier

Soft

$$\min_{\mathbf{w},\, b} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^{n} \max\left(0,\ 1 - y_i\left(w^T x_i - b\right)\right) \quad (8)$$

$$\text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \forall i$$

$$\geq 1 - \max\left(0,\ 1 - y_i\left(w^T x_i - b\right)\right)$$

hinge loss     not differentiable



Margin

$w^T x + b = 1$

$w^T x + b = 0$

$w^T x + b = -1$

$x^+$

$x^+$

$x^-$

$n$

Support Vectors

$x_2$

$x_1$

- ● denotes +1
- ○ denotes -1

# Review: Support Vector Machine

- Dual problem

  $$\max_{\alpha} \; \mathcal{L}_{\mathcal{D}}(\alpha) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^{\top} \mathbf{x}_j \quad \text{(9)}$$

  *quadratic programming*

  $$\text{s.t.} \; \sum_{i=1}^{N} \alpha_i y_i = 0, \; \alpha_i \geq 0, i = 1, \ldots, N$$

  $\alpha_i \leq \frac{1}{2\eta\lambda}$

  *generalize.*

  *other methods* { • sub-gradient • coordinate descent.

- Kernel SVM — ( Linear ⟶ nonlinear )
  - replace $\mathbf{x}$ with $\phi(\mathbf{x})$
  - $\phi(\mathbf{x}_i)^{\top} \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$
  - $k(\mathbf{x}_i, \mathbf{x}_j)$ is a kernel function, e.g, $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$
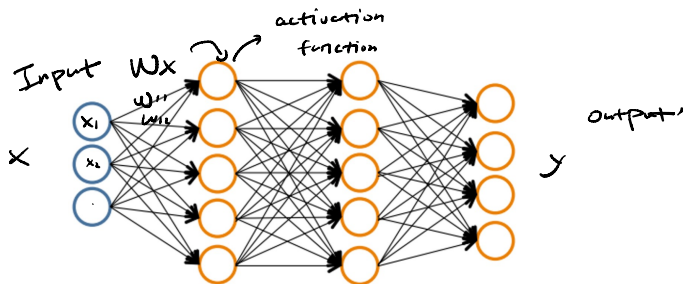
- Slacked SVM

  + + + / – –

  *outlier*

  $$\min_{\mathbf{w}, b, \xi} \; \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^{N} \xi_i \quad \text{(10)}$$

  $$\text{s.t.} \; y_i(\mathbf{w}^{\top} \mathbf{x}_i + b) \geq 1 - \xi_i, \; \xi_i \geq 0, \; i = 1, \ldots, N$$

- Fully connected feedforward network (multi-layer perceptron, MLP)



$\quad$ $\mathbf{h}^{(1)} = f^{(1)}(\mathbf{x})$ $\quad$ $\mathbf{h}^{(2)} = f^{(2)}(\mathbf{h}^{(1)})$ $\quad$ $\ldots$ $\quad$ $\mathbf{y} = f^{(L)}(\mathbf{h}^{(L-1)})$

$$\text{Or} \quad \mathbf{y} = f^{(L)} \circ \cdots \circ f^{(1)}(\mathbf{x})$$
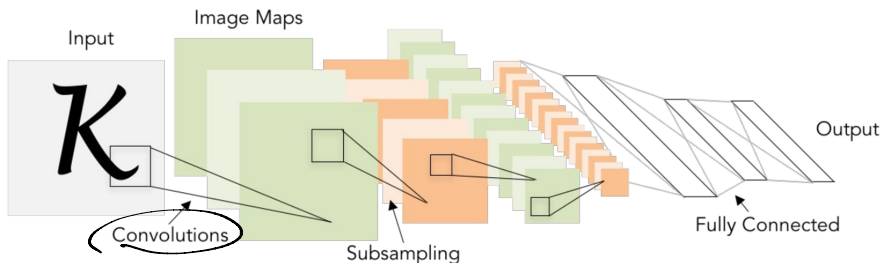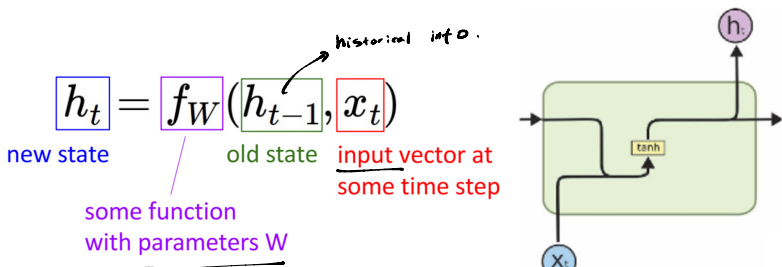
- Convolutional neural network (CNN)



Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

LeNet-5

$$(f * g)[n]$$
$$= \sum_{m=-\infty}^{\infty} f(m)\, g(n-m)$$

- Recurrent neural network (RNN) (create a control system)

$$h_t = f_W(h_{t-1}, x_t)$$

historical info.

new state

old state

input vector at some time step

some function with parameters W



- Works well for sequential data.

NLP
text
human voice

$x = (x_1, x_2 \cdots, x_N)$
↓
letter

- Recurrent neural network (RNN)

$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state

some function
with parameters W

old state   input vector at
some time step



- Other models for sequential data
  - LSTM, GRU, etc
  - Transformer
    - Widely used in LLMs such as GPTs
    - Not covered by DDA3020 and DDA4210

Classification on MNIST handwritten digits dataset

http://yann.lecun.com/exdb/mnist/



Figure: Samples of MNIST
($28 \times 28$ gray-scale images, 60k
for training, 10 k for testing)

| classifier | test error rate (%) |
|---|---|
| linear classifier (least squares) | 12.0 |
| k-nearest-neighbors | 5.0 |
| generalized linear classifier (Gaussian basis 1000) | 3.6 |
| neural network (MLP) 500-300 HU, softmax | 1.53 |
| CNN LeNet-5 | 0.95 |
| SVM (Gaussian kernel) | 1.4 |

## Classification on Fashion-MNIST dataset

`https://cloudxlab.com/blog/fashion-mnist-using-machine-learning/`



| Label | Description | Examples |
|-------|-------------|----------|
| 0 | T-Shirt/Top | |
| 1 | Trouser | |
| 2 | Pullover | |
| 3 | Dress | |
| 4 | Coat | |
| 5 | Sandals | |
| 6 | Shirt | |
| 7 | Sneaker | |
| 8 | Bag | |
| 9 | Ankle boots | |

Figure: Samples of Fashion-MNIST (28 × 28 gray-scale images, 60k for training, 10 k for testing)

| classifier | test error rate (%) |
|-----------|---------------------|
| softmax | 15.3 |
| decision tree | 21.06 |
| random forest | 15.18 |
| neural network (MLP) (256-128-100 HU) | 12.6 |
| CNN | <8 |
| HOG+SVM | 7.4 |
| Google AutoML | 6.1 |

More results are at `https://github.com/zalandoresearch/fashion-mnist`

# Review: K-Means Clustering

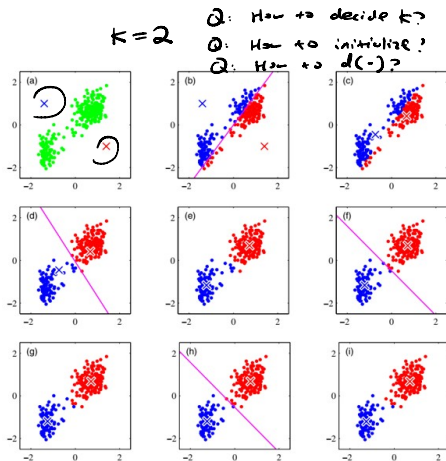- Clustering (unsupervised learning): given a set of $D$-dimensional data $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N\}$, partition them into $K$ clusters such that each data point is similar to the data in the same cluster and dissimilar to the data in different clusters.

- Denote cluster $j$ by $\mathcal{C}_j$ and let $\boldsymbol{\mu}_j$ be the centroid of $\mathcal{C}_j$. K-means clustering minimizes

$$J(\mu) = \sum_{j=1}^{K} \sum_{\mathbf{x} \in \mathcal{C}_j} \left\| \mathbf{x} - \boldsymbol{\mu}_j \right\|^2 \tag{11}$$

- Algorithm (alternate)
    1. Assign each data point to the closest center
    2. Update the cluster center

$K=2$

Q: How to decide k?
Q: How to initialize?
Q: How to $d(\cdot)$?

# Review: Gaussian Mixture Models

Idea: consider using conditional Gaussians to represent data distribution.

- Multivariate Gaussian distribution

$$p(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

- Gaussian mixture distribution

$$p(\mathbf{x}) = \sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \Sigma_j)$$

- $K$ different Gaussian distributions
- $\{\pi_j\}$: mixing coefficients
- $\sum_{j=1}^{K} \pi_j = 1, \quad 0 \leq \pi_j \leq 1$

- Algorithm: Expectation-Maximization
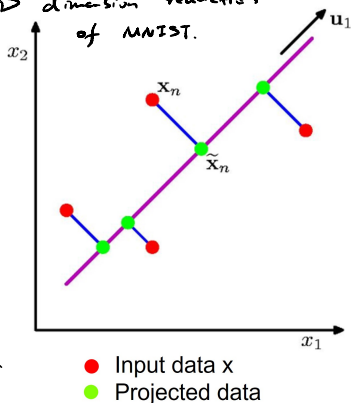
# Review: Principal Component Analysis

- PCA: find the orthogonal projection of data onto a lower-dimensional subspace that
  - maximizes the variance of projected data
  - or minimizes the reconstruction error, i.e.,

$$J = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^{\top}\mathbf{x}_i\|^2 \quad (12)$$

≤ VD.

  * $\mathbf{x} \in \mathbb{R}^D, \ \mathbf{U} \in \mathbb{R}^{D \times d}$

- Solution of PCA: eigenvalue decomposition or singular value decomposition

- Doesn't work well for 2D dimension reduction of MNIST.



$x_2$

$\mathbf{u}_1$

$\mathbf{x}_n$

$\tilde{\mathbf{x}}_n$

$x_1$

- ● Input data x
- ● Projected data

# Review: More Topics (optional)

- Bayes' theorem, maximum likelihood estimation (MLE), maximum a posteriori estimation (MAP)
- Classification evaluation metrics
  - Precision, recall, accuracy, F1-score, AUC (TPR/FPR)
- Cross-validation
- Over-/under-fitting and bias-variance trade-off
- Expectation maximization
- Kernel density estimation
- Clustering evaluation metrics