# CIS 678 Machine Learning

**Course Review Week**
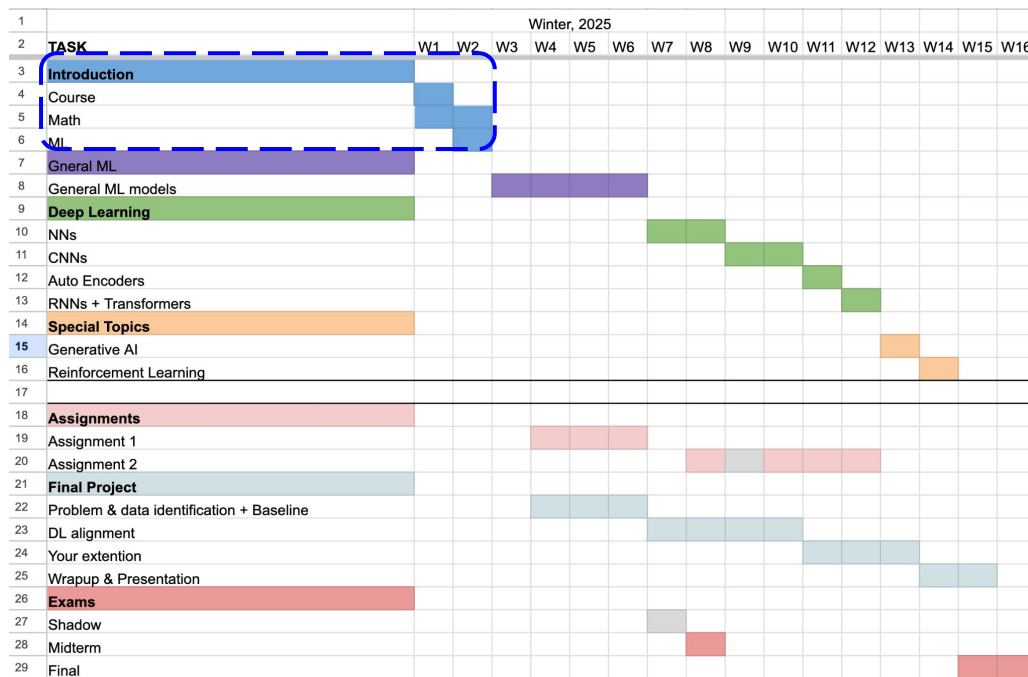
# Our Plan

| # | TASK | Winter, 2025 | | | | | | | | | | | | | | | |
|---|------|----|----|----|----|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|
| | | W1 | W2 | W3 | W4 | W5 | W6 | W7 | W8 | W9 | W10 | W11 | W12 | W13 | W14 | W15 | W16 |
| 3 | **Introduction** | | | | | | | | | | | | | | | | |
| 4 | Course | ● | | | | | | | | | | | | | | | |
| 5 | Math | ● | ● | | | | | | | | | | | | | | |
| 6 | ML | | ● | | | | | | | | | | | | | | |
| 7 | Gneral ML | | | | | | | | | | | | | | | | |
| 8 | General ML models | | | ● | ● | ● | ● | | | | | | | | | | |
| 9 | **Deep Learning** | | | | | | | | | | | | | | | | |
| 10 | NNs | | | | | | | ● | ● | | | | | | | | |
| 11 | CNNs | | | | | | | | ● | ● | | | | | | | |
| 12 | Auto Encoders | | | | | | | | | | ● | ● | | | | | |
| 13 | RNNs + Transformers | | | | | | | | | | | | ● | | | | |
| 14 | **Special Topics** | | | | | | | | | | | | | | | | |
| 15 | Generative AI | | | | | | | | | | | | | ● | | | |
| 16 | Reinforcement Learning | | | | | | | | | | | | | | ● | | |
| 17 | | | | | | | | | | | | | | | | | |
| 18 | **Assignments** | | | | | | | | | | | | | | | | |
| 19 | Assignment 1 | | | | ● | ● | ● | | | | | | | | | | |
| 20 | Assignment 2 | | | | | | | | ● | ● | ● | ● | ● | | | | |
| 21 | **Final Project** | | | | | | | | | | | | | | | | |
| 22 | Problem & data identification + Baseline | | | | ● | ● | ● | | | | | | | | | | |
| 23 | DL alignment | | | | | | | ● | ● | ● | ● | ● | | | | | |
| 24 | Your extention | | | | | | | | | | ● | ● | ● | ● | | | |
| 25 | Wrapup & Presentation | | | | | | | | | | | | | | ● | ● | |
| 26 | **Exams** | | | | | | | | | | | | | | | | |
| 27 | Shadow | | | | | | | ● | | | | | | | | | |
| 28 | Midterm | | | | | | | | ● | | | | | | | | |
| 29 | Final | | | | | | | | | | | | | | | ● | ● |

# Background diagnostics + ML Introduction

**Vector space**
**Proximity or distance metric**
- L1/Manhattan distance
- L2/Euclidean distance,
- Cosine distances

**kNN model:**
- Distance based
- Can be applied to both Regression and Classification tasks

**Probability (measuring uncertainty)**
**Probability distributions**

**ML introduction:**
Linear to Polynomial Regression

| | TASK | W1 | W2 | W3 | W4 | W5 | W6 | W7 | W8 | W9 | W10 | W11 | W12 | W13 | W14 | W15 | W16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | Winter, 2025 | | | | | | | | | | | | | |
| 2 | **TASK** | | | | | | | | | | | | | | | | |
| 3 | **Introduction** | | | | | | | | | | | | | | | | |
| 4 | Course | | | | | | | | | | | | | | | | |
| 5 | Math | | | | | | | | | | | | | | | | |
| 6 | ML | | | | | | | | | | | | | | | | |
| 7 | Gneral ML | | | | | | | | | | | | | | | | |
| 8 | General ML models | | | | | | | | | | | | | | | | |
| 9 | **Deep Learning** | | | | | | | | | | | | | | | | |
| 10 | NNs | | | | | | | | | | | | | | | | |
| 11 | CNNs | | | | | | | | | | | | | | | | |
| 12 | Auto Encoders | | | | | | | | | | | | | | | | |
| 13 | RNNs + Transformers | | | | | | | | | | | | | | | | |
| 14 | **Special Topics** | | | | | | | | | | | | | | | | |
| 15 | Generative AI | | | | | | | | | | | | | | | | |
| 16 | Reinforcement Learning | | | | | | | | | | | | | | | | |
| 17 | | | | | | | | | | | | | | | | | |
| 18 | **Assignments** | | | | | | | | | | | | | | | | |
| 19 | Assignment 1 | | | | | | | | | | | | | | | | |
| 20 | Assignment 2 | | | | | | | | | | | | | | | | |
| 21 | **Final Project** | | | | | | | | | | | | | | | | |
| 22 | Problem & data identification + Baseline | | | | | | | | | | | | | | | | |
| 23 | DL alignment | | | | | | | | | | | | | | | | |
| 24 | Your extention | | | | | | | | | | | | | | | | |
| 25 | Wrapup & Presentation | | | | | | | | | | | | | | | | |
| 26 | **Exams** | | | | | | | | | | | | | | | | |
| 27 | Shadow | | | | | | | | | | | | | | | | |
| 28 | Midterm | | | | | | | | | | | | | | | | |
| 29 | Final | | | | | | | | | | | | | | | | |

# Linear Regression



y (weight)

X (height)

### Model

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\Theta = \{\beta_0, \beta_1\}$$

### Fitting Error

$$\epsilon = |\hat{y} - y|$$

### Optimization function

$$E_\Theta = \frac{1}{2} \sum_{i=1}^{N} (\hat{y}_i - y_i)^2$$

$$\Theta^* = \mathrm{argmin}_\Theta E\{(x_i, y_i)\}_{i=1,\cdots,N}$$

# Linear to Polynomial Regression



**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

|  | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ |  | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ |  |  | -25.43 | -5321.83 |
| $w_3^\star$ |  |  | 17.37 | 48568.31 |
| $w_4^\star$ |  |  |  | -231639.30 |
| $w_5^\star$ |  |  |  | 640042.26 |
| $w_6^\star$ |  |  |  | -1061800.52 |
| $w_7^\star$ |  |  |  | 1042400.18 |
| $w_8^\star$ |  |  |  | -557682.99 |
| $w_9^\star$ |  |  |  | 125201.43 |

Model generalization

# Linear to Polynomial Regression



**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

|  | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ |  | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ |  |  | -25.43 | -5321.83 |
| $w_3^\star$ |  |  | 17.37 | 48568.31 |
| $w_4^\star$ |  |  |  | -231639.30 |
| $w_5^\star$ |  |  |  | 640042.26 |
| $w_6^\star$ |  |  |  | -1061800.52 |
| $w_7^\star$ |  |  |  | 1042400.18 |
| $w_8^\star$ |  |  |  | -557682.99 |
| $w_9^\star$ |  |  |  | 125201.43 |

Absolute values are increasing

Model generalization

# Linear to Polynomial Regression

Model

$$\hat{y} = \beta_0 + \beta_1 x$$

$$\Theta = \{\beta_0, \beta_1\}$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

$$\epsilon = |\hat{y} - y|$$

Optimization function

Essentially, the same formulation

Generally ML vs Math conventions

$$E_\Theta = \frac{1}{2} \sum_{i=1}^{N} (\hat{y}_i - y_i)^2$$

$$W^* = \operatorname{argmin}_W E\{(x_i, t_i)\}_{i=1,\cdots,N}$$

$$\Theta^* = \operatorname{argmin}_\Theta E\{(x_i, y_i)\}_{i=1,\cdots,N}$$

# Linear to Polynomial Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

Regularizer

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{1}{2} \|\mathbf{w}\|^2$$

Model generalization: Regularization

**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

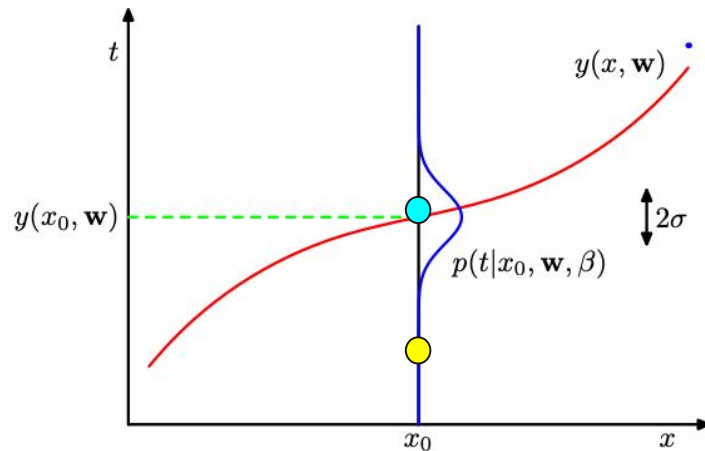| | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ | | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ | | | -25.43 | -5321.83 |
| $w_3^\star$ | | | 17.37 | 48568.31 |
| $w_4^\star$ | | | | -231639.30 |
| $w_5^\star$ | | | | 640042.26 |
| $w_6^\star$ | | | | -1061800.52 |
| $w_7^\star$ | | | | 1042400.18 |
| $w_8^\star$ | | | | -557682.99 |
| $w_9^\star$ | | | | 125201.43 |

Abs values
Are increasing

# Linear to Polynomial Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{1}{2} \|\mathbf{w}\|^2$$

$\|\mathbf{w}\|^2$

decreasing

**Model generalization: Regularization**

**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

| | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ | | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ | | | -25.43 | -5321.83 |
| $w_3^\star$ | | | 17.37 | 48568.31 |
| $w_4^\star$ | | | | -231639.30 |
| $w_5^\star$ | | | | 640042.26 |
| $w_6^\star$ | | | | -1061800.52 |
| $w_7^\star$ | | | | 1042400.18 |
| $w_8^\star$ | | | | -557682.99 |
| $w_9^\star$ | | | | 125201.43 |

Abs values
Are increasing

# Linear to Polynomial Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

| | $M=0$ | $M=1$ | $M=6$ | $M=9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ | | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ | | | -25.43 | -5321.83 |
| $w_3^\star$ | | | 17.37 | 48568.31 |
| $w_4^\star$ | | | | -231639.30 |
| $w_5^\star$ | | | | 640042.26 |
| $w_6^\star$ | | | | -1061800.52 |
| $w_7^\star$ | | | | 1042400.18 |
| $w_8^\star$ | | | | -557682.99 |
| $w_9^\star$ | | | | 125201.43 |

$$\|\mathbf{w}\|^2$$

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{1}{2} \|\mathbf{w}\|^2$$

*How to control this?*

decreasing

Model generalization: Regularization

# Linear to Polynomial Regression

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

**Table 1.1** Table of the coefficients $\mathbf{w}^\star$ for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

| | $M = 0$ | $M = 1$ | $M = 6$ | $M = 9$ |
|---|---|---|---|---|
| $w_0^\star$ | 0.19 | 0.82 | 0.31 | 0.35 |
| $w_1^\star$ | | -1.27 | 7.99 | 232.37 |
| $w_2^\star$ | | | -25.43 | -5321.83 |
| $w_3^\star$ | | | 17.37 | 48568.31 |
| $w_4^\star$ | | | | -231639.30 |
| $w_5^\star$ | | | | 640042.26 |
| $w_6^\star$ | | | | -1061800.52 |
| $w_7^\star$ | | | | 1042400.18 |
| $w_8^\star$ | | | | -557682.99 |
| $w_9^\star$ | | | | 125201.43 |

**Hyperparameter**

$$\|\mathbf{W}\|^2$$

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

*Who to control this?*

**Model generalization: Regularization**

# Probabilistic equivalent

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

# Probabilistic equivalent

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

# Probabilistic equivalent

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2$$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}\left(t_n | y(x_n, \mathbf{w}), \beta^{-1}\right)$$

# General ML models, Assignment 1 & Midterm

# General ML models

**Regression**

- kNN
- Linear Regression
- Decision Tree
- Meta learners
  - Random Forest Regressor
  - Boosting Regressor
- Support Vector Regressor (SVRs)

**Classification**

- kNN
- Logistic Regression
- Decision Tree
- Meta Learners
  - Random Forest Classifier
  - Boosting Classifiers
- Support Vector Classifiers (SVCs)
- Naive Bayes

# General ML models



**Supervised**

Regression:
- kNN
- Linear Regression
- Decision Tree
- Meta learners
  - Random Forest Regressor
  - Boosting Regressor
- Support Vector Regressor (SVRs)

Classification:
- kNN
- Logistic Regression
- Decision Tree
- Meta Learners
  - Random Forest Classifier
  - Boosting Classifiers
- Support Vector Classifiers (SVCs)
- Naive Bayes

Regression | Parametric vs nonparametric | Classification

# General ML models

## Unsupervised

### Clustering

- Clustering algorithms
  - **k-means**: Centroid Based
  - k-modes:  Mode Based (categorical)
  - **Hierarchical clustering**: Distance connectivity based
  - **GMM**: Distribution based
  - **DBSCAN**: Density Based
- How to choose the optimal number of clusters.

### Linear Dimensionality Reduction

- Principal Component Analysis (PCA)
- Singular Value Decomposition (SVD)

# General ML models

**Model generalization**

- Universal concepts (applies to all models)
  - Cross validation
  - HP optimization

- Overfitting
- Underfitting

Universal concepts

Overfitting vs Under fitting

# General ML models

- Training set, Validation set, Test set
- iid data

| Training | Test |
|---|---|

| Training | Validation |
|---|---|

**Data splits**

- Overfitting
- Underfitting

**Overfitting vs Under fitting**

# General ML models

**Model generalization**

- Training set, Validation set, Test set
- iid data

| Training | Hold out set |

| Training | Validation |

**Data splits**

- Overfitting
- Underfitting

**Overfitting vs Under fitting**

# K-fold-cross validation

Data

| D1 | D2 | D3 |

3-fold-cv

Train    validate

*What HP gives the best validation score?*

# Assignment 1

### Vector space



Figure 1: A read and a blue box containing a set of orange and green balls.

### Probability & distribution

1. [**2x2 points**] For each dimension $d \in [2^1, 2^3, 2^5, ..., 2^{11}]$, sample 100 random points from corresponding vector spaces *(sample code to generate random samples is provided below)*, and

   - Record the $l2$ and the *cosine* distances between all pairs (of points); then
   - Fit two normal/Gaussian distributions, one for each distance metric. Share the mean ($\mu$) and the standard distribution ($\sigma$) parameters of each distribution that you have learned.
   - Plot these normal/Gaussian distributions using your preferred visualization package(s).

   **normal/Gaussian distribution**: $p(x) \sim \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

   **Sample code** (to generate $n = 100$ random samples from a $d = 4$ dimensional vector space):

   ```
   import numpy as np
   d, n = 4, 100
   sample_data = np.random.randn(n, d)
   ```

2. [**2x2 points**]
   There are some orange and green balls in a red and a blue box as illustrated in Figure 1. Someone (blinded by tying a piece of cloth around his/her eyes) is asked to pick up a ball twice (with replacement, i.e. the ball is placed back in the same box from where it was picked up).

   (a) What is the **log-probability** of the first ball to be **orange** and second ball to be **green**?
   (b) What is is **probability** that both the time the ball came from the **red** box?

   **Note**: We assume, the selection of the red and the blue boxes follow an **uniform** distribution (given that the person is not able to see the color of those boxes).

# Assignment 1

Loss function

Building your own model from scratch

3. [**2 points**] Given a Linear Regression model, $\Theta = \{\beta_0 = 0.1, \beta_1 = 0.9, \beta_2 = -3.5\}$ (where $\beta_0$ is the bias and $\beta_1$ and $\beta_2$ are the parameters associated with two input features/variables), the regularizer parameter $\lambda = 1.5$ and the following data set (with $y$ being the target variable), estimate the

- Quadratic error or loss with l2 regularizer as defined below

$$E_\Theta = \frac{1}{2}\sum_{i=0}^{N-1}(y_i - \hat{y}_i)^2 + \frac{\lambda}{2}\|\boldsymbol{\beta}\|^2$$

4. [**4 points**] You are asked to fit a **second order/degree** polynomial regression model, $y = \beta_0 + \beta_1 x + \beta_2 x^2$ with parameter $\Theta = \{\beta_0, \beta_1, \beta_2\}$ (where $\beta_0$ is the bias of your model) on the following dataset.

**data file:**
https://raw.githubusercontent.com/mdkamrulhasan/data-public/refs/heads/main/miscellaneous/second_degree_polynomial_regression_data.csv

For the setup below, we ask you to find out the updated version of $\Theta$ after two(2) iterations of any gradient descent algorithm (you can use the algorithm that we shared as a part of our linear regression model illustration in the class):

<u>Setup details</u>:

- Use the quadratic error/loss function
- Initialize, $\Theta_0 = \{\beta_0 = 0.0, \beta_1 = 0.0, \beta_2 = 1.0\}$
- Use learning rate (parameter), $L = 0.001$

**Note**: You have to define your error/loss function and also will need to estimate partial derivatives of your loss function.

# Assignment 1

Regularization:
model generalization

5. [**3x2 points**] For the given dataset below, we ask you develop, test and compare following models (follow the instructions under the **Detailed specification** section below).

   (a) Linear Regression

   (b) Linear Regression with $l1$ regularizer

   (c) Linear Regession with $l2$ regularizer

**data file:**

https://raw.githubusercontent.com/mdkamrulhasan/data_mining_kdd/main/data/medical-cost/insurance.csv

**Detailed specification**:

- Use a 50%-50% test/train setup.

- Use Mean Squared Error (MSE) as your evaluation metric.

- Visualize (using Bar charts) the model parameter **absolute** values (covert any negative values to positives before plotting).

- You can use python packages such as **sklearn** for your solution.

# Shadow test followed by Mid-term!

# Sequential cross validation



*Invalid configuration:*

*Training on the last two and validation on the* **first fold**.

# Sequential cross validation



*Always follow!*

# Sequential cross validation



*Always follow!*

# Sequential cross validation



*Always follow!*

# Sequential cross validation



*Always follow!*

# Deep Learning

Assignment 2

# Logistic Regression

- Probabilistic classifier

- Sigmoid function

$$p(x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}}$$

# Neural Networks (Node)



$z = w_0 + w_1 x$

$\dfrac{1}{1+e^{-z}}$

1

$w_0$

x

$w_1$

$z$

y

**Linear**
(aggregation)

**Activation
function**

$y \in \{0, 1\}$

**Logistic Regression**

# Neural Networks (Node)

A NN **node** with 2 sequential steps:
- Aggregation
- Activation

1

$w_0$

$x$

$w_1$

$z$

$y$

$y$

**Linear**
(aggregation)

**+**

**Activation function**

# Neural Networks (Node)



| Name | Plot | Function, $g(x)$ |
|---|---|---|
| Identity | | $x$ |
| Binary step | | $\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| Logistic, sigmoid, or soft step | | $\sigma(x) \doteq \dfrac{1}{1+e^{-x}}$ |
| Hyperbolic tangent (tanh) | | $\tanh(x) \doteq \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$ |
| Soboleva modified hyperbolic tangent (smht) | | $\text{smht}(x) \doteq \dfrac{e^{ax} - e^{-bx}}{e^{cx} + e^{-dx}}$ |
| Rectified linear unit (ReLU)[13] | | $(x)^+ \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max(0, x) = x\mathbf{1}_{x>0}$ |

1

$w_0$

x

$w_1$

$z = w_0 + w_1 x$

**Linear**
(aggregation)

$z$

???

**Activation**
**function**

y

**A NN with ???**
**Activation function**

# Feed-forward (FF) neural networks



- Feedforward NN
- Perceptron

$y = f(\mathbf{X}, \mathbf{Z})$

Input layer          Hidden layer(s)          Output layer

# Error Back propagation

# Gradient Descent (Error Back Propagation)

**The Backwards Pass**

Let's focus on $\frac{\partial E_{total}}{\partial w_5}$

What would be the gradient update for w5?



Adapted from

# Gradient Descent (Error Back Propagation)

**The Backwards Pass**



$$\frac{\partial net_{o1}}{\partial w_5} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial E_{total}}{\partial out_{o1}} = \frac{\partial E_{total}}{\partial w_5}$$

output h1

w5

output h2   w6   net $_{o1}$  out $_{o1}$

b2

1

$E_{o1}$ = ½(target $_{o1}$ - out $_{o1}$)²

$E_{total}$ = $E_{o1}$ + $E_{o2}$

$$E_{total} = \tfrac{1}{2}(target_{o1} - out_{o1})^2 + \tfrac{1}{2}(target_{o2} - out_{o2})^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \tfrac{1}{2}(target_{o1} - out_{o1})^{2-1} * -1 + 0$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = -(target_{o1} - out_{o1}) = -(0.01 - 0.75136507) = 0.74136507$$

[Adapted from](#)

# Gradient Descent (Error Back Propagation)

**The Backwards Pass**

$$\frac{\partial net_{o1}}{\partial w_5} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial E_{total}}{\partial out_{o1}} = \frac{\partial E_{total}}{\partial w_5}$$

output h1

**w5**

output h2   **w6**   **net**$_{o1}$ **out**$_{o1}$

**b2**

1

$E_{o1}$ = ½(target$_{o1}$ - **out**$_{o1}$)²

$E_{total}$ = $E_{o1}$ + $E_{o2}$

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}}$$

$$\frac{\partial out_{o1}}{\partial net_{o1}} = out_{o1}(1 - out_{o1})$$

$$= 0.75136507(1 - 0.75136507)$$

$$= 0.186815602$$

Adapted from

# Gradient Descent (Error Back Propagation)

**The Backwards Pass**



$$\frac{\partial net_{o1}}{\partial w_5} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial E_{total}}{\partial out_{o1}} = \frac{\partial E_{total}}{\partial w_5}$$

$$E_{o1} = \tfrac{1}{2}(target_{o1} - out_{o1})^2$$

$$E_{total} = E_{o1} + E_{o2}$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial w_5} = 1 * out_{h1} * w_5^{(1-1)} + 0 + 0 = out_{h1} = 0.593269992$$

[Adapted from](#)

# Gradient Descent (Error Back Propagation)

**Hidden Layer**



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

$$E_{total} = E_{o1} + E_{o2}$$

$E_{o1} + E_{o2}$

- While changing w5 affects only O1, a change in w1 will change both O1 and O2

# Gradient Descent (Error Back Propagation)

**Hidden Layer**

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

$E_{total} = E_{o1} + E_{o2}$

$E_{o1} + E_{o2}$

- Can you think of an arbitrary node in a giant and complex NN? What challenges we may encounter?

# Convolutional NNs

- State of the Art for CV and some other problems
- Filters/Convolutional Kernels

Examples:

- Alexnet
- VGG
- ResNet
- GoogLeNet
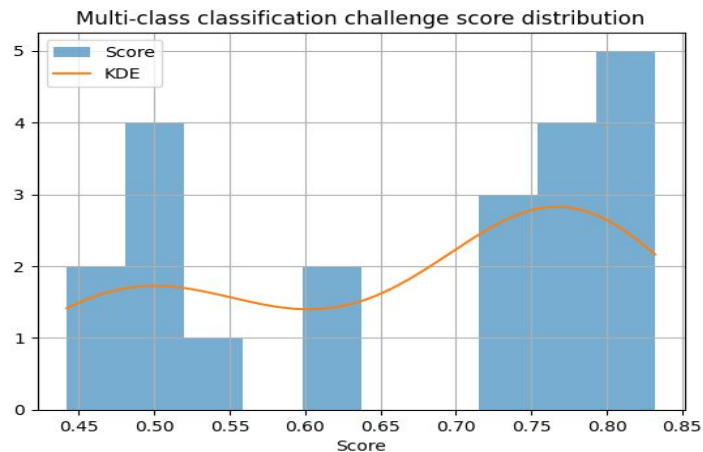- ..

ref

# Convolutional NNs

# Convolutional NNs

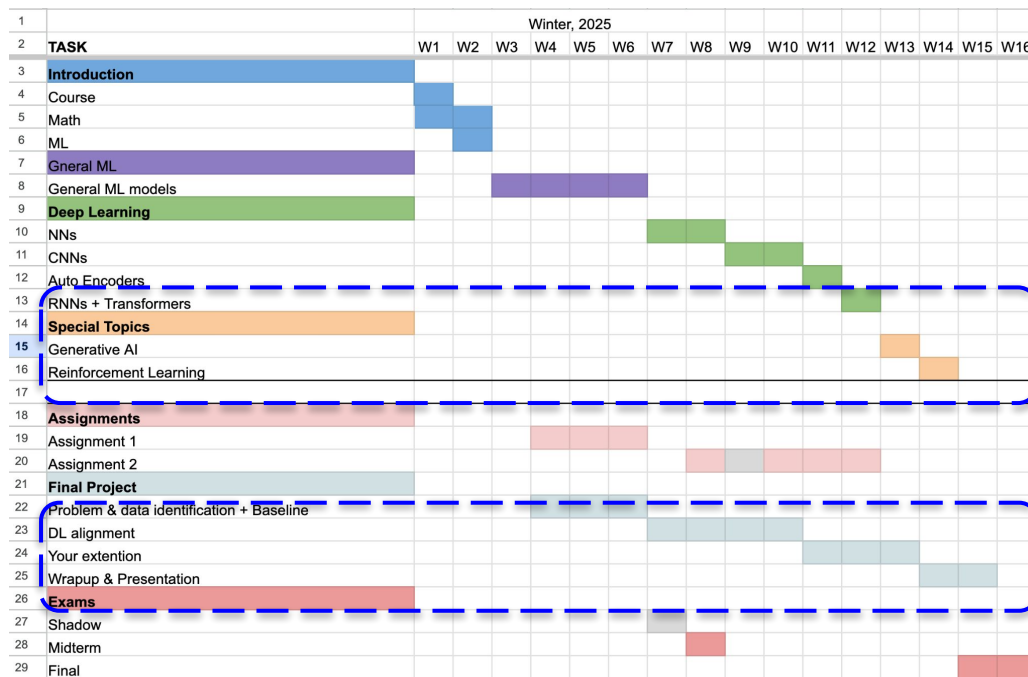# Assignment 1 (classification challenge)

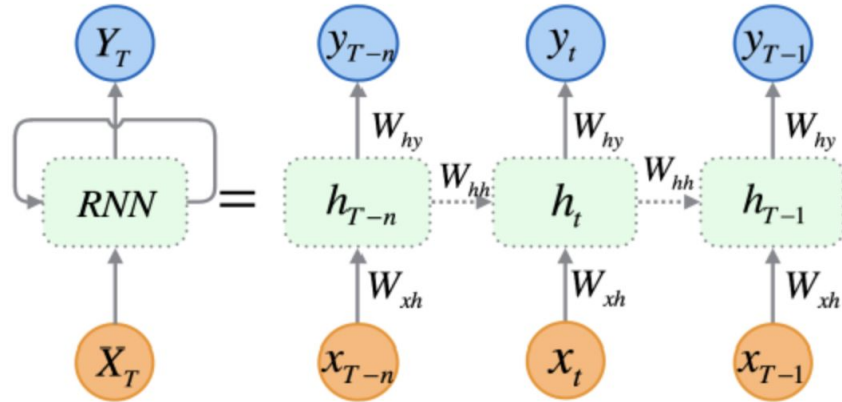- Deep learning models doubled the performance to general ML models

# DL continuation + Project



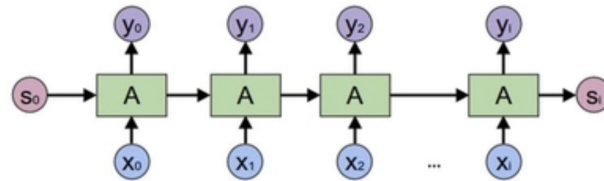| | TASK | W1 | W2 | W3 | W4 | W5 | W6 | W7 | W8 | W9 | W10 | W11 | W12 | W13 | W14 | W15 | W16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | | | | Winter, 2025 | | | | | | | | | | | | |
| 2 | **TASK** | W1 | W2 | W3 | W4 | W5 | W6 | W7 | W8 | W9 | W10 | W11 | W12 | W13 | W14 | W15 | W16 |
| 3 | **Introduction** | | | | | | | | | | | | | | | | |
| 4 | Course | ▓ | | | | | | | | | | | | | | | |
| 5 | Math | ▓ | ▓ | | | | | | | | | | | | | | |
| 6 | ML | | ▓ | | | | | | | | | | | | | | |
| 7 | Gneral ML | | | | | | | | | | | | | | | | |
| 8 | General ML models | | | ▓ | ▓ | ▓ | | | | | | | | | | | |
| 9 | **Deep Learning** | | | | | | | | | | | | | | | | |
| 10 | NNs | | | | | | | ▓ | | | | | | | | | |
| 11 | CNNs | | | | | | | | ▓ | | ▓ | | | | | | |
| 12 | Auto Encoders | | | | | | | | | | | ▓ | | | | | |
| 13 | RNNs + Transformers | | | | | | | | | | | | ▓ | | | | |
| 14 | **Special Topics** | | | | | | | | | | | | | | | | |
| 15 | Generative AI | | | | | | | | | | | | | ▓ | | | |
| 16 | Reinforcement Learning | | | | | | | | | | | | | | ▓ | | |
| 17 | | | | | | | | | | | | | | | | | |
| 18 | **Assignments** | | | | | | | | | | | | | | | | |
| 19 | Assignment 1 | | | | ▓ | ▓ | | | | | | | | | | | |
| 20 | Assignment 2 | | | | | | | ▓ | ▓ | ▓ | | | | | | | |
| 21 | **Final Project** | | | | | | | | | | | | | | | | |
| 22 | Problem & data identification + Baseline | | | | | | | ▓ | ▓ | | | | | | | | |
| 23 | DL alignment | | | | | | | | | ▓ | ▓ | | | | | | |
| 24 | Your extention | | | | | | | | | | | ▓ | ▓ | | | | |
| 25 | Wrapup & Presentation | | | | | | | | | | | | | ▓ | ▓ | | |
| 26 | **Exams** | | | | | | | | | | | | | | | | |
| 27 | Shadow | | | | | | | ▓ | | | | | | | | | |
| 28 | Midterm | | | | | | | | ▓ | | | | | | | | |
| 29 | Final | | | | | | | | | | | | | | | ▓ | ▓ |

# Recurrent Neural Networks

Examples:

- LSTMs
- GRU
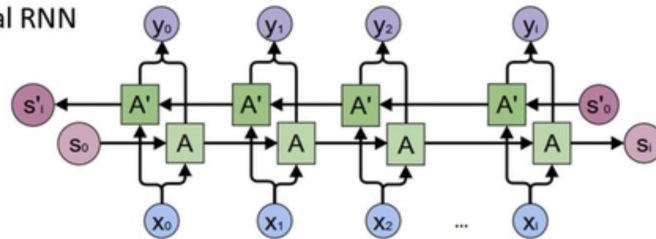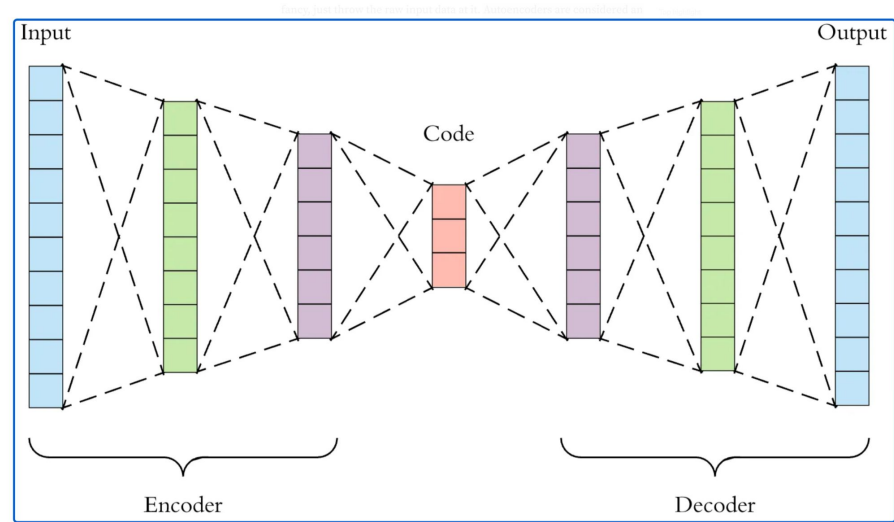


ref

# RNNs

# Unsupervised learning (nonlinear)

- Auto Encoders
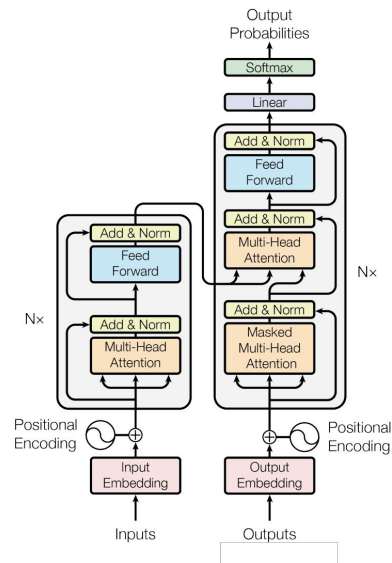- Restricted Boltzmann Machines (RBMs)

# Transformers

Examples:

- Encoder decoder pair
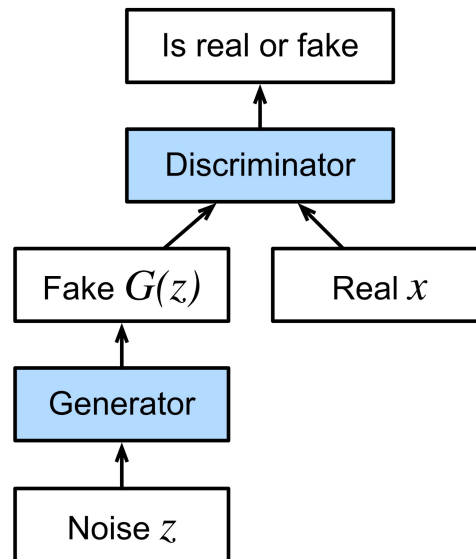- GPT
- BERT

**BERT**

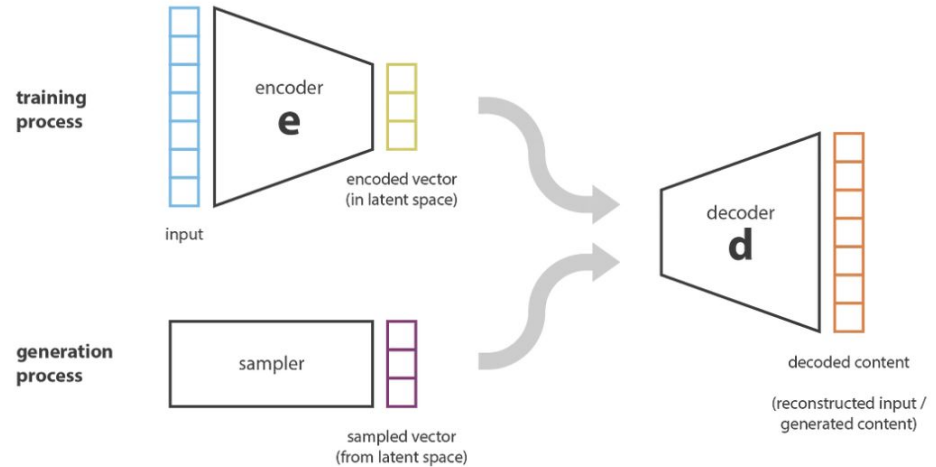Encoder

**GPT**

Decoder

# Generative AI

Examples:

- Generative Adversarial Networks (GANs)

ref

# Generative AI

Examples:

- Generative Adversarial Networks (G
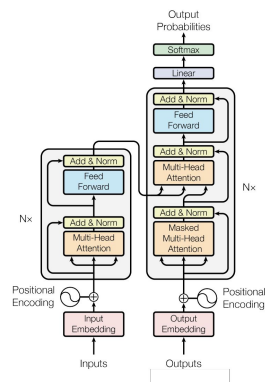- **Variational Autoencoders (VAEs)**



training process

encoder
**e**

input

encoded vector
(in latent space)

generation process

sampler

sampled vector
(from latent space)

decoder
**d**

decoded content

(reconstructed input /
generated content)

ref

# Foundational models

- **We don't train independent models explicitly (say machine translation)**

BERT

Encoder

GPT

Decoder

Machine translation

# Foundational models

- We don't train independent models explicitly (say machine translation)
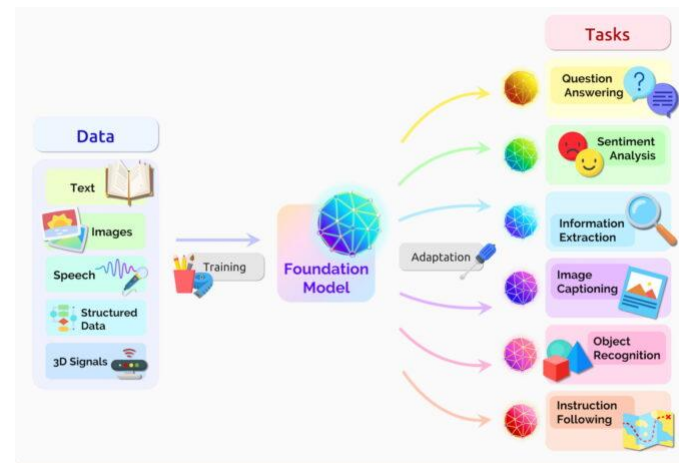- Or document summarization

BERT

Encoder

GPT

Decoder



Document summarization

# Foundational models

- We don't train independent models explicitly (say machine translation)
- Or document summarization
- **We train a Base model for multiple tasks jointly, and then fine tune for specific tasks**

# Foundational models

- We don't train independent models explicitly (say machine translation)
- Or document summarization
- **We train a Base model for multiple tasks jointly, and then fine tune for specific tasks**



Training compute of notable machine learning models by domain, 2012–23
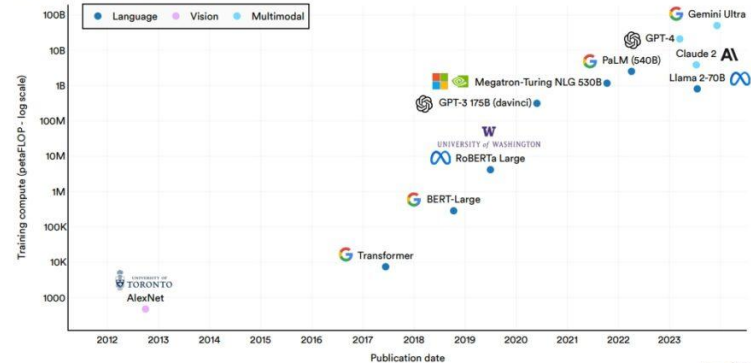Source: Epoch, 2023 | Chart: 2024 AI Index report

Figure 1.3.7

# Training Neural Network Challenges

- Intractable gradients
  - Vanishing, and
  - Exploding gradients

# Training Neural Network Challenges

- Intractable gradients
  - Vanishing, and
  - Exploding gradients
- Various normalizations
  - Input normalization (standard scalar)
  - Batch normalization
  - Layer normalization

# Training Neural Network Challenges

- Intractable gradients
    - Vanishing, and
    - Exploding gradients
- Various normalizations
    - Input normalization (standard scalar)
    - Batch normalization
    - Layer normalization
- Controlling overfitting
    - Regularization
    - Early stopping
    - Drop out

# QA