**DA261 Machine Learning**

**Fundamentals**

Bachelor of Science (Honours) in Data Science and Artificial Intelligence

**DA261: Machine Learning Fundamentals**

Instructor: Teena Sharma, Ph.D.

Mehta Family School of Data Science and Artificial Intelligence,

Indian Institute of Technology Guwahati, India

*Teena Sharma, Ph.D.*

---

## Learning Objectives

**DA261 Machine Learning**

**Fundamentals**

- Linear regression
- Overfitting
- Underfitting
- Regularization
- Differentiate between regression and classification tasks

*Teena Sharma, Ph.D.*

2

---

## Introduction to Supervised Learning

**DA261 Machine Learning**

**Fundamentals**

Given a set of input patterns $X$ and a corresponding set of labels $Y$, the underlined mapping function $f: X \rightarrow Y$ is discovered. The classifier operates in two distinct phases, i.e., a training phase (model tuning), and an operating phase, in which the model is kept fixed and tested with different and new data.

*Supervised learning is useful when the map from inputs to outputs is unknown, but we have a lot of input-to-output examples.*

*Teena Sharma, Ph.D.*

3

---

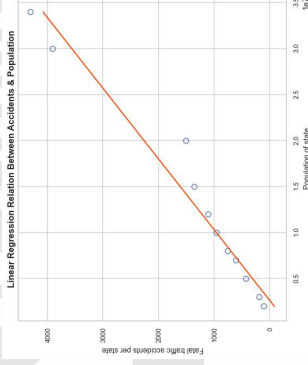## Introduction to Supervised Learning

**DA261 Machine Learning**

**Fundamentals**

$$y_i = f(x_i)$$

- When $y_i$ is continuous, this is a regression problem.
- When $y_i$ is discrete, this is a classification problem.

*Teena Sharma, Ph.D.*

4

## Linear Regression

- Regression analysis is a set of statistical methods used for the estimation of relationships between a dependent variable and one or more independent variables.

- Regression models are used to predict a continuous value.



Linear Regression Relation Between Accidents & Population

---

## Linear Regression

- Let's consider the linear regression problem.

- Feature vector: $\boldsymbol{x_i} = [x_{i_1}, x_{i_2}, \ldots, x_{i_d}]$

- Target: $y_i$

- Prediction: $\hat{y}_i = w_1 x_{i_1} + w_2 x_{i_2} + \cdots w_d x_{i_d}$ (weighted sum of features)

  Weight of feature 1 | feature 1

---

## Linear Regression

- Let's consider the linear regression problem.

- Feature vector: $\boldsymbol{x_i} = [x_{i_1}, x_{i_2}, \ldots, x_{i_d}]$

- Target: $y_i$

- Prediction: $\hat{y}_i = w_1 x_{i_1} + w_2 x_{i_2} + \cdots w_d x_{i_d}$ (weighted sum of features)

  Weight of feature 1 | feature 1

  If $w_1 > 0$, more populations means more murder rate

---

## Linear Regression

- Let's consider the linear regression problem.

- Feature vector: $\boldsymbol{x_i} = [x_{i_1}, x_{i_2}, \ldots, x_{i_d}]$

- Target: $y_i$

- Prediction: $\hat{y}_i = w_1 x_{i_1} + w_2 x_{i_2} + \cdots w_d x_{i_d}$ (weighted sum of features)

  Weight of feature 1 | feature 1

  We want to make $y_i - \hat{y}_i$ small, i.e. we want prediction to be close to the target.

# Linear Regression

Let's consider the linear regression problem.

- Residual: $\quad y_i - \hat{y}_i = y_i - \sum_{j=1}^{d} w_j x_{ij}$

We need to tune the weights to **minimize the prediction error.**

*Teena Sharma, Ph.D.*

---

# Least Square Error

What is the classic way to minimize this error considering we have $n$ training examples?

**Minimizing Least Square Error!**

We need to find $w$ such that it minimizes: $\displaystyle\sum_{j=1}^{n}(y_i - \hat{y}_i)^2 \;=\; \sum_{j=1}^{n}(y_i - w^T x_i)^2$

$$= ||Xw - y||^2$$

where, $w = [w_1, w_2, \dots, w_d]$

*Teena Sharma, Ph.D.*

---

# Minimizing Least Square Error

$$||Xw - y||^2$$

$X = [x_1, x_2, \dots, x_d]$

$n * d$ — Feature matrix

$w = [w_1, w_2, \dots, w_d]^T$

$d * 1$ — Weight vector

$y = [y_1, y_2, \dots, y_n]^T$

$n * 1$ — Target vector

*Teena Sharma, Ph.D.*

---

# Minimizing Least Square Error



$$||Xw - y||^2$$

$X = [x_1, x_2, \dots, x_d]$

$n * d$ — Feature matrix

This is unknown

$w = [w_1, w_2, \dots, w_d]^T$

$d * 1$ — Weight vector

$y = [y_1, y_2, \dots, y_n]^T$

$n * 1$ — Target vector

*Teena Sharma, Ph.D.*

## Minimizing Least Square Error

Finding solution to get $w$:

$$\text{minimize } \|Xw - y\|^2 \quad \text{over } w \in R^d$$

Set gradient equal to zero and solve for $w$

Gradient is a vector of partial derivatives:

$$\nabla f(w) = \left[\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \ldots, \frac{\partial f}{\partial w_d}\right]^T$$

---

## Finding Gradient

Finding the gradient of linear functions:

Linear function: $f(w) = \alpha^T w + \beta$ —— Scalar

Vector

Finding gradient:

Step-1: $f(w) = \sum_{i=1}^d \alpha_i w_i + \beta$

Step-2: $\dfrac{\partial f(w)}{\partial w_i} = \alpha_i$

Step-3: $\nabla f(w) = [\alpha_1, \alpha_2, \ldots, \alpha_d]^T = \alpha$

---

## Finding Gradient

Finding the gradient of Quadratic functions:

Quadratic function: $f(w) = w^T A w$

Finding gradient:

Step-1: Convert to summation notation:

$$f(w) = \sum_{i=1}^d \sum_{j=1}^d w_i a_{ij} w_j$$

where, $a_{ij}$ is the element in row $i$ and column $j$ of $A$. To help with computing the partial derivatives, it helps to re-write it in the form

$$f(w) = \sum_{i=1}^d a_{ii} w_i^2 + \sum_{j\neq i} w_i a_{ij} w_j$$

---

## Finding Gradient

Finding the gradient of Quadratic functions:

Quadratic function: $f(w) = w^T A w$

Finding gradient:

Step-2: Take the partial derivative with respect to an element $k$:

$$\frac{\partial}{\partial w_k}\left[\sum_{i=1}^d a_{ii} w_i^2 + \sum_{j\neq i} w_i a_{ij} w_j\right] = 2a_{kk} w_k + \sum_{j\neq k} w_j a_{jk} + \sum_{j\neq k} a_{kj} w_j$$

where, $a_{ij}$ is the element in row $i$ and column $j$ of $A$. To help with computing the partial derivatives, it helps to re-write it in the form

$$\frac{\partial}{\partial w_k}\left[\sum_{i=1}^d a_{ii} w_i^2 + \sum_{j\neq i} w_i a_{ij} w_j\right] = \sum_{j=1}^d w_j a_{jk} + \sum_{j=1}^d a_{kj} w_j$$

## Finding Gradient

$$\frac{\partial}{\partial w_k}\left[\sum_{\{i=1\}}^{d} a_{ii}w_i^2 + \sum_{\{j\neq i\}} w_i a_{ij}w_j\right] = 2a_{kk}w_k + \sum_{\{j\neq k\}} w_j a_{jk} + \sum_{\{j\neq k\}} a_{kj}w_j$$

---

## Finding Gradient

Finding the gradient of Quadratic functions:

Quadratic function:   $f(w) = w^T A w$

Finding gradient:

Step-3: Assemble the partial derivatives into a vector:

$$\frac{\partial}{\partial w_k}\left[\sum_{\{i=1\}}^{d} a_{ii}w_i^2 + \sum_{\{j\neq i\}} w_i a_{ij}w_j\right] = \sum_{\{j=1\}}^{d} w_j a_{jk} + \sum_{\{j=1\}}^{d} a_{kj}w_j$$

$$\nabla f(w) = \begin{bmatrix} \frac{\partial}{\partial w_1} \\ \frac{\partial}{\partial w_2} \\ \cdots \\ \frac{\partial}{\partial w_d} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{d} w_j a_{j1} + \sum_{j=1}^{d} a_{1j}w_j \\ \sum_{j=1}^{d} w_j a_{j2} + \sum_{j=1}^{d} a_{2j}w_j \\ \cdots \\ \sum_{j=1}^{d} w_j a_{jd} + \sum_{j=1}^{d} a_{dj}w_j \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{d} w_j a_{j1} \\ \sum_{j=1}^{d} w_j a_{j2} \\ \cdots \\ \sum_{j=1}^{d} w_j a_{jd} \end{bmatrix} + \begin{bmatrix} \sum_{j=1}^{d} a_{1j}w_j \\ \sum_{j=1}^{d} a_{2j}w_j \\ \cdots \\ \sum_{j=1}^{d} a_{dj}w_j \end{bmatrix}$$

---

## Finding Gradient

Finding the gradient of Quadratic functions:

Quadratic function:   $f(w) = w^T A w$

Finding gradient:

Step-4: Convert to matrix notation:

$$\frac{\partial}{\partial w_k}\left[\sum_{\{i=1\}}^{d} a_{ii}w_i^2 + \sum_{\{j\neq i\}} w_i a_{ij}w_j\right] = \sum_{\{j=1\}}^{d} w_j a_{jk} + \sum_{\{j=1\}}^{d} a_{kj}w_j$$

$$\nabla f(w) = \begin{bmatrix} \sum_{j=1}^{d} w_j a_{j1} \\ \sum_{j=1}^{d} w_j a_{j2} \\ \cdots \\ \sum_{j=1}^{d} w_j a_{jd} \end{bmatrix} + \begin{bmatrix} \sum_{j=1}^{d} a_{1j}w_j \\ \sum_{j=1}^{d} a_{2j}w_j \\ \cdots \\ \sum_{j=1}^{d} a_{dj}w_j \end{bmatrix} = A^T w + A w = (A^T + A)w$$

---

## Finding Gradient

Finding the gradient of Quadratic functions:

Quadratic function:   $f(w) = w^T A w$

Finding gradient:

Final result:   $\nabla f(w) = (A^T + A)w$

## Least Square Problem

Our least square problem:

$$\underset{w \in R^d}{\text{Minimize}} \ \ ||Xw - y||^2 \ \sim \ \underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}||Xw - y||^2$$

$$f(w) = \frac{1}{2}||Xw - y||^2$$

---

## Least Square Problem

Our least square problem:

$$\underset{w \in R^d}{\text{Minimize}} \ \ ||Xw - y||^2 \ \sim \ \underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}||Xw - y||^2$$

$$f(w) = \frac{1}{2}||Xw - y||^2 = \frac{1}{2}(Xw - y)^T(Xw - y)$$

$$= \frac{1}{2}(w^TX^T - y^T)(Xw - y)$$

$$= \frac{1}{2}(w^TX^T(Xw - y) - y^T(Xw - y))$$

$$= \frac{1}{2}(w^TX^TXw - w^TX^Ty - y^TXw + y^Ty)$$

$$= \frac{1}{2}(w^TX^TXw - (X^Ty)^Tw - (X^Ty)^Tw + y^Ty)$$

$$= \frac{1}{2}w^TX^TXw - (X^Ty)^Tw + \frac{1}{2}y^Ty$$

$[||a||^2 = a^Ta]$

$[(ab)^T = b^Ta^T]$

$w^TX^Ty = (y^TXw)^T = y^TXw$
$= (X^Ty)^Tw$
Since both are scalars, so they are equal

$[w^Tv = v^Tw]$

---

## Least Square Problem

Objective function: $f(w) = \frac{1}{2}\underbrace{w^TX^TXw}_{\text{Quadratic}} - \underbrace{(X^Ty)^Tw}_{\text{Linear}} + \underbrace{\frac{1}{2}y^Ty}_{\text{Constant}}$

Gradient computation: $\nabla f(w) = X^TXw - X^Ty$

Setting gradient to zero:

---

## Least Square Problem

Objective function: $f(w) = \frac{1}{2}\underbrace{w^TX^TXw}_{\text{Quadratic}} - \underbrace{(X^Ty)^Tw}_{\text{Linear}} + \underbrace{\frac{1}{2}y^Ty}_{\text{Constant}}$

Gradient computation: $\nabla f(w) = X^TXw - X^Ty$

Setting gradient to zero: $\nabla f(w) = X^TXw - X^Ty = 0 \rightarrow (X^TX)w = X^Ty$

(if $(X^TX)$ is invertible $\rightarrow (X^TX)^{-1}(X^TX)w = (X^TX)^{-1}X^Ty$

$\rightarrow w = (X^TX)^{-1}X^Ty$

$[A^{-1}A = I]$

Least Square Solution

# Cost of Computing Least Squares:

- Getting $X^T y$ costs $O(nd)$
  - Getting $X^T X$ costs $O(nd^2)$
  - Solving a $d * d$ linear system costs $O(d^3)$: Using Gaussian Elimination

> **Total cost = $O(nd^2 + d^3)$**

---

# Is Least Squares any good?

Issues with least squares model:

- It assumes a linear relationship between $x$ and $y$.
- It might predict poorly for new values of $x$.
- $X^T X$ might not be invertible.
- It is sensitive to outliers.
- It might predict outside known range of $y$ values.
- It always uses all features.
- Number of dimensions $d$ might be so big we can't store $X^T X$.

---

# Problem with Basic Linear Models?

This is our model: $\quad \hat{y}_i = \sum_{j=1}^{d} w_j x_{ij} = w^T x_i$



- This is always satisfied: $x_i = 0 \rightarrow \hat{y}_i = 0$
- The fitted line always pass through origin. *Undesirable!*

---

# Problem with Basic Linear Models?

This is our model: $\quad \hat{y}_i = \sum_{j=1}^{d} w_j x_{ij} = w^T x_i$



- Now our updated model: $\quad y_i = w^T x_i + \beta$
  
  Bias
- Now: $x_i = 0 \rightarrow \hat{y}_i = \beta$

## How to include the Bias variable?

To include the bias variable in the objective function of the least square problem:

1. Modify the weight Vector

$$w = [w_1, w_2, ..., w_d]^T \rightarrow \bar{w} = [w_0, w_1, w_2, ..., w_d]^T, \ where\ w_0 = \beta$$

$d * 1$

$(d + 1) * 1$

---

## Solution with Bias variable

$$(\bar{X}^T\bar{X})\bar{w} = \bar{X}^T\bar{y} \rightarrow \bar{w} = (\bar{X}^T\bar{X})^{-1}\bar{X}^T\bar{y}$$

Let's take simple example:

$X \rightarrow 11 * 1$
$w \rightarrow 1 * 1$
$y \rightarrow 11 * 1$



$n = 11, \ d = 1$

| x | y |
|---|---|
| -1.0000 | 1.0882 |
| -0.8000 | 0.6600 |
| -0.6000 | 0.4089 |
| -0.4000 | 0.2720 |
| -0.2000 | 0.1334 |
| 0.0000 | -0.0489 |
| 0.2000 | 0.0875 |
| 0.4000 | 0.1524 |
| 0.6000 | 0.3548 |
| 0.8000 | 0.6605 |
| 1.0000 | 1.0072 |

---

## Solution with Bias variable

Linear Regression with bias:   $\hat{y}_i = \beta + w_1 x_i$

We can clearly see that our data is not linear!

Now let's assume a quadratic basis:

$$\hat{y}_i = \beta + w_1 x_i + w_2 x_i^2$$



*Fitting only linear + bias*

---

## Solution with Quadratic variable

Model with Quadratic bias:   $\hat{y}_i = \beta + w_1 x_i + w_2 x_i^2$



$$w = [\beta, w_1, w_2]^T$$

$$y = [y_1, y_2, ..., y_n]^T$$

Final solution:   $w^* = (X_{poly}^T X_{poly})^{-1} X_{poly}^T y$

## Solution with Quadratic variable

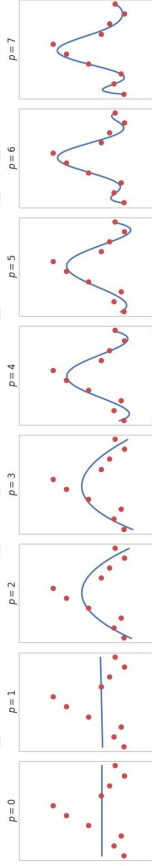We can also have polynomial of degree $p$:    $\hat{y}_i = \beta + w_1 x_i + w_2 x_i^2 + \cdots + w_p x_i^p$

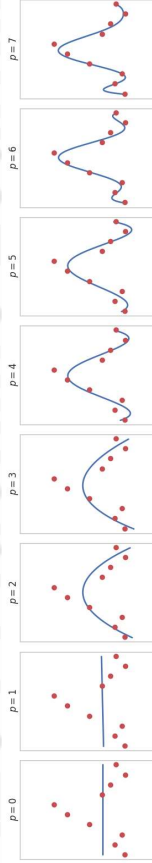$$w = [\beta, w_1, w_2 \ldots w_p]^T$$
$$y = [y_1, y_2, \ldots, y_n]^T$$

*Teena Sharma, Ph.D.*

---

## True or False?

**Higher $p$ means better model (Think and try to answer)?**



*Teena Sharma, Ph.D.*

---

## True or False?

**Higher $p$ means better model (Think and try to answer)?**



**Low degree**
- Less likely to fit data well
- Model does not change much with change in data

*Teena Sharma, Ph.D.*

---

## True or False?

**Higher $p$ means better model (Think and try to answer)?**



**Low degree**
- Less likely to fit data well
- Model does not change much with change in data

**High degree**
- Very likely to fit data well
- Model change a lot with change in data

*Teena Sharma, Ph.D.*

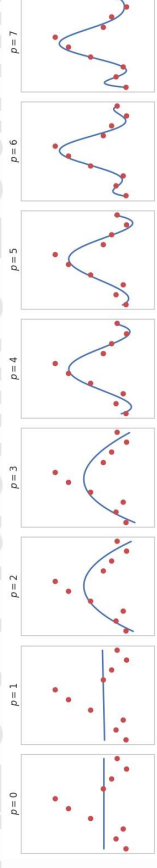## True or False?

**Higher *p* means better model (Think and try to answer)?**



p=0, p=1, p=2, p=3, p=4, p=5, p=6, p=7

**Low degree**
- Less likely to fit data well
- Model does not change much with change in data

**High Bias, Low Variance**
- High expected error due to wrong model

**High degree**
- Very likely to fit data well
- Model change a lot with change in data

**High Variance, Low Bias**
- High sensitive the model is to particular training set

---

## Overfitting



p = 7

- Fit all the data, training error is very low
- The model fit is sensitive to the training data.
- Does not generalize well for new test data (not in training example).

**High Variance, Low Bias**
- High sensitive the model is to particular training set

---

## True or False?

**Higher *p* means better model (Think and try to answer)?**



p=0, p=1, p=2, p=3, p=4, p=5, p=6, p=7

**Low degree**
- Less likely to fit data well
- Model does not change much with change in data

**High Bias, Low Variance**
- High expected error due to wrong model

**High degree**
- Very likely to fit data well
- Model change a lot with change in data

---

## Underfitting



p = 1

- Model does not fit the data well.
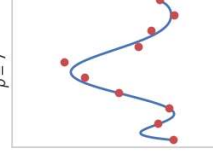- The model fit is not very sensitive to the training data

**High Bias, Low Variance**
- High expected error due to wrong model

# Determining Model Overfitting or Underfitting

We can determine model underfitting or overfitting by observing the error during model training and testing.

For that we need to understand following things:

➢ What is model training?

➢ What is model testing?

➢ Can we use same dataset for both training and testing?

---

# Training and Testing

First we have to divide our dataset into two parts:

➢ Training Data

➢ Testing Data

You can select the dataset in 70% train and 30% test data.

➢ Train model using training data

➢ Test model using testing data

**Good training performance is useless, if it does not perform well with test data.**

*Example:*

**Train data**

| Inhabitants | Below income | Percentage Unemployed | Murders |
|---|---|---|---|
| 587000.0 | 16.5 | 6.2 | 11.2 |
| 643000.0 | 20.5 | 5.4 | 13.4 |
| 635000.0 | 20.3 | 5.3 | 40.7 |
| 692000.0 | 19.6 | 5.2 | 22.1 |
| 1248000.0 | 16.2 | 5.6 | 13.4 |
| 643000.0 | 21.3 | 6.5 | 25.9 |
| 1964000.0 | 17.3 | 6.4 | 12.6 |
| 1531000.0 | 13.3 | 6.4 | 9.6 |
| 713000.0 | 13.1 | 6.5 | 9.5 |
| 749000.0 | 12.7 | 5.6 | 10.1 |
| 7895000.0 | 12.7 | 5.8 | 12.7 |
| 762000.0 | 12.0 | 5.7 | 12.3 |
| 2793000.0 | 16.4 | 6.3 | 8.5 |
| 834000.0 | 24.9 | 8.3 | 28.9 |

**Test data**

| Inhabitants | Below income | Percentage Unemployed | Murders |
|---|---|---|---|
| 7895000.0 | 17.9 | 6.7 | 14.8 |
| 762000.0 | 22.4 | 8.6 | 25.3 |
| 2793000.0 | 20.2 | 8.4 | 22.7 |
| 3350000.0 | 16.9 | 6.7 | 25.7 |

---

# Training and Testing

There are two phases of supervised learning:

➢ **Training Phase:** Model fitting based on the training data $(X_{train}, y_{train})$.

➢ **Testing Phase:** Model evaluation on test data $(X_{test}, y_{test})$, that was not used for training the model.

We also need to use evaluation metric to find the **testing error** e.g. $||y_{pred} - y_{test}||^2$

By observing only testing error we can say how **good** our model is.

**Test data should never be used for training the model.**

---

# Managing Model Complexity

➢ In practical scenarios, the relationship between features and the target variable can be **quite complex.**

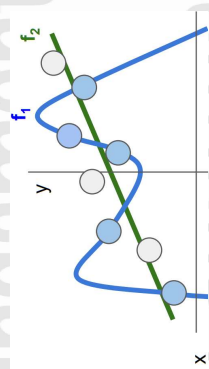➢ However, **highly complex models tend to overfit the training data.**

So, how can we systematically control model complexity to achieve better generalization?

## Regularization

One most popular technique is **Regularization**

➢ We can add the **penalty** on the complexity of the model by using **regularization.**



Regularization pushes against fitting the data too well so we don't fit noise in the data

---

## Regularization

The optimization problem for regression:  $\underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}\|Xw - y\|^2$

Instead we need to solve this!

Standard L2-regularization strategy is to add a penalty on the L2-norm:

$$\underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2$$

L2-norm:  $\|X\|_2 = \left(\sum x_i^2\right)^{1/2} = \sqrt{\{x_1^2 + x_2^2 + \cdots + x_N^2\}}$

---

## Regularization

The optimization problem for regression:  $\underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}\|Xw - y\|^2$

Instead we need to solve this!

Standard L2-regularization strategy is to add a penalty on the L2-norm:

$$\underset{w \in R^d}{\text{Minimize}} \ \frac{1}{2}\|Xw - y\|^2 + \boxed{\frac{\lambda}{2}\|w\|^2}$$  Regularization term

L2-norm:  $\|X\|_2 = \left(\sum x_i^2\right)^{1/2} = \sqrt{\{x_1^2 + x_2^2 + \cdots + x_N^2\}}$

---

## Try this to find the solution

Objective function:  $f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2$

Find gradient: ?

Set gradient to zero: ?

Solution:?

# Thank you!

DA261 Machine Learning

Teena Sharma, Ph.D.

---

## Try this to find the solution

Objective function: $f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{\lambda}{2}\|w\|^2$

Find gradient: ?

Set gradient to zero: ?

Solution: $\boxed{w^* = (X^TX + \lambda I)^{-1}X^Ty}$

DA261 Machine Learning

Teena Sharma, Ph.D.

Source: https://cs231n.stanford.edu/slides/2024/lecture_3.pdf