



HUST

ĐẠI HỌC BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY

ONE LOVE. ONE FUTURE.



ĐẠI HỌC
BÁCH KHOA HÀ NỘI
HANOI UNIVERSITY
OF SCIENCE AND TECHNOLOGY

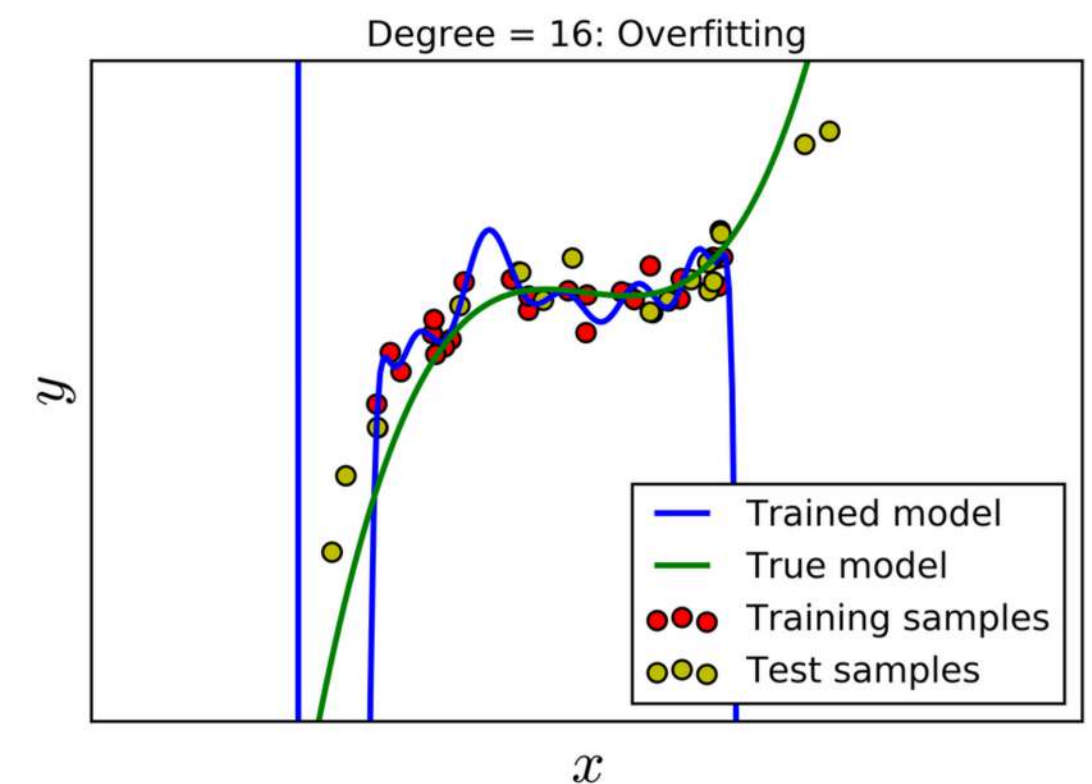
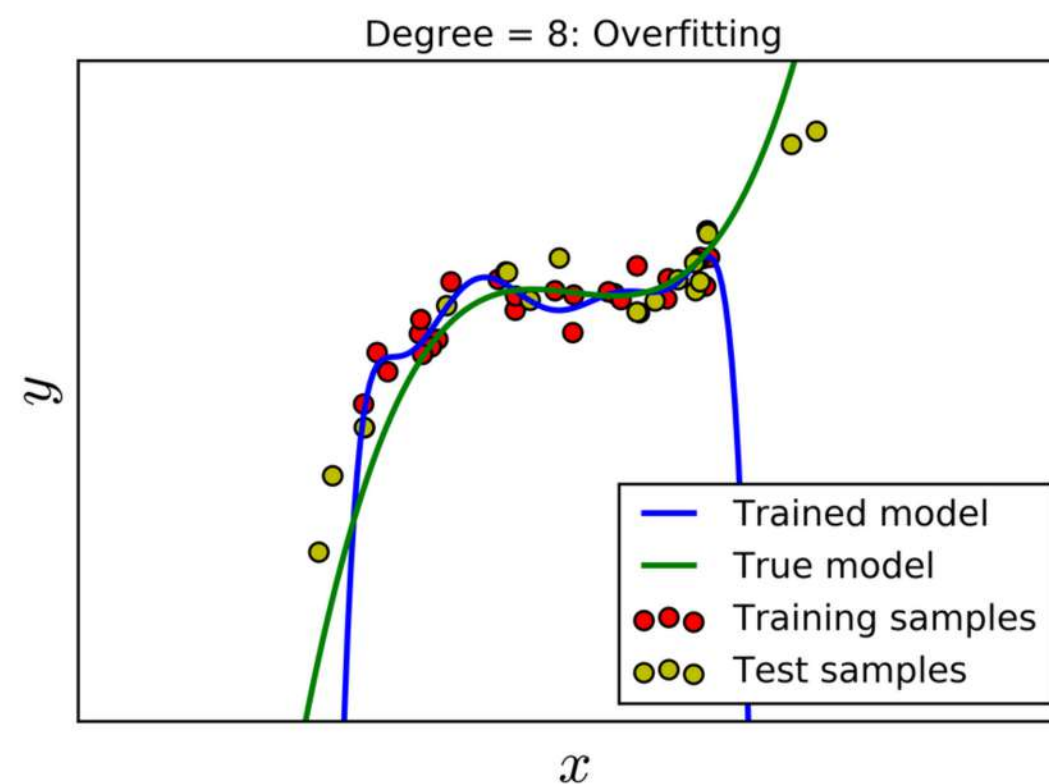
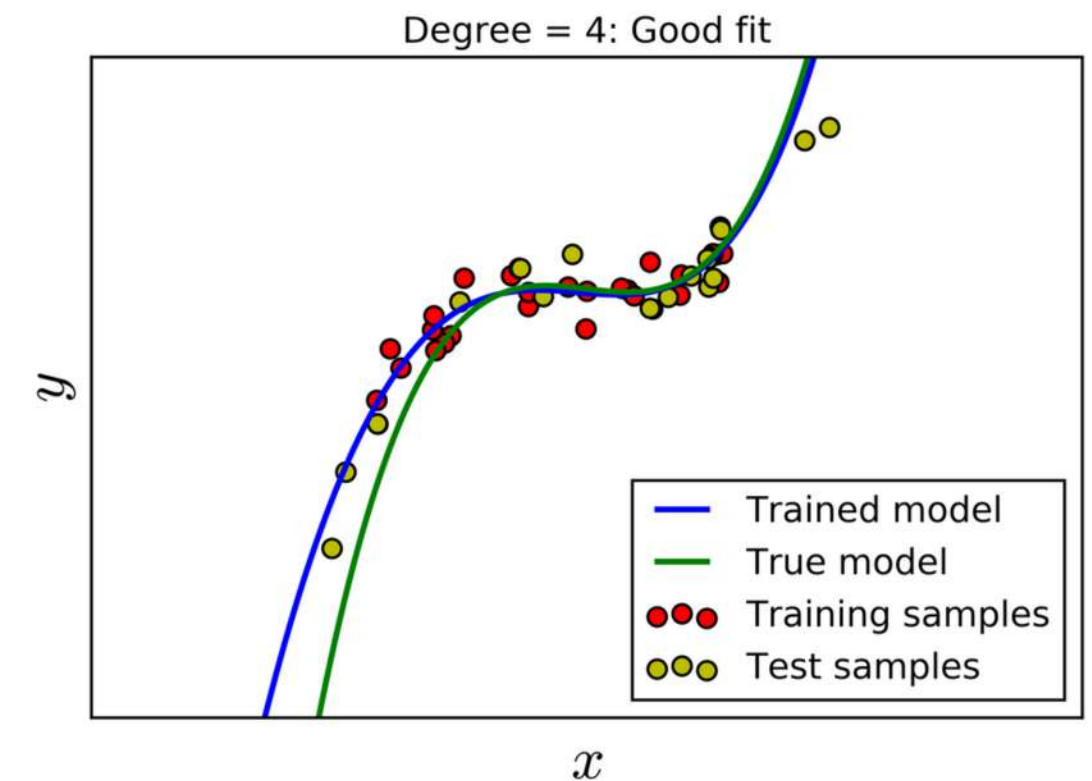
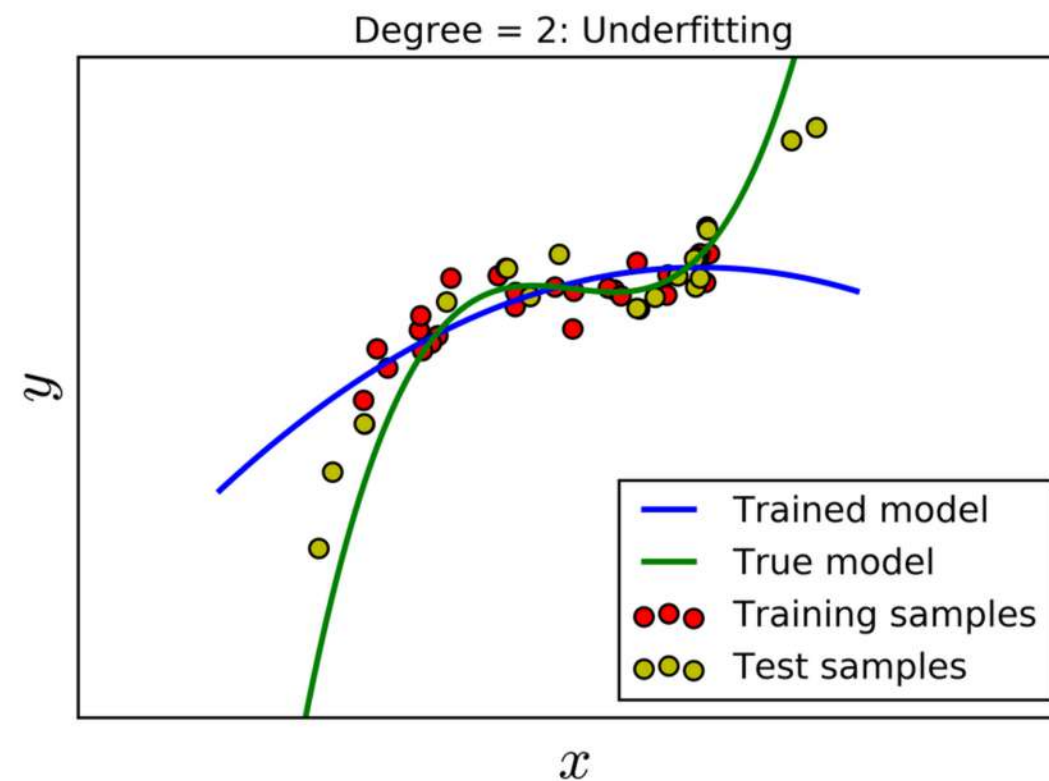
Các phương pháp chống Overfitting trong học máy

ONE LOVE. ONE FUTURE.

- Overfitting
- Các phương pháp chống Overfitting
 - Validation
 - Regularization
- Ứng dụng trong bài toán thực tế

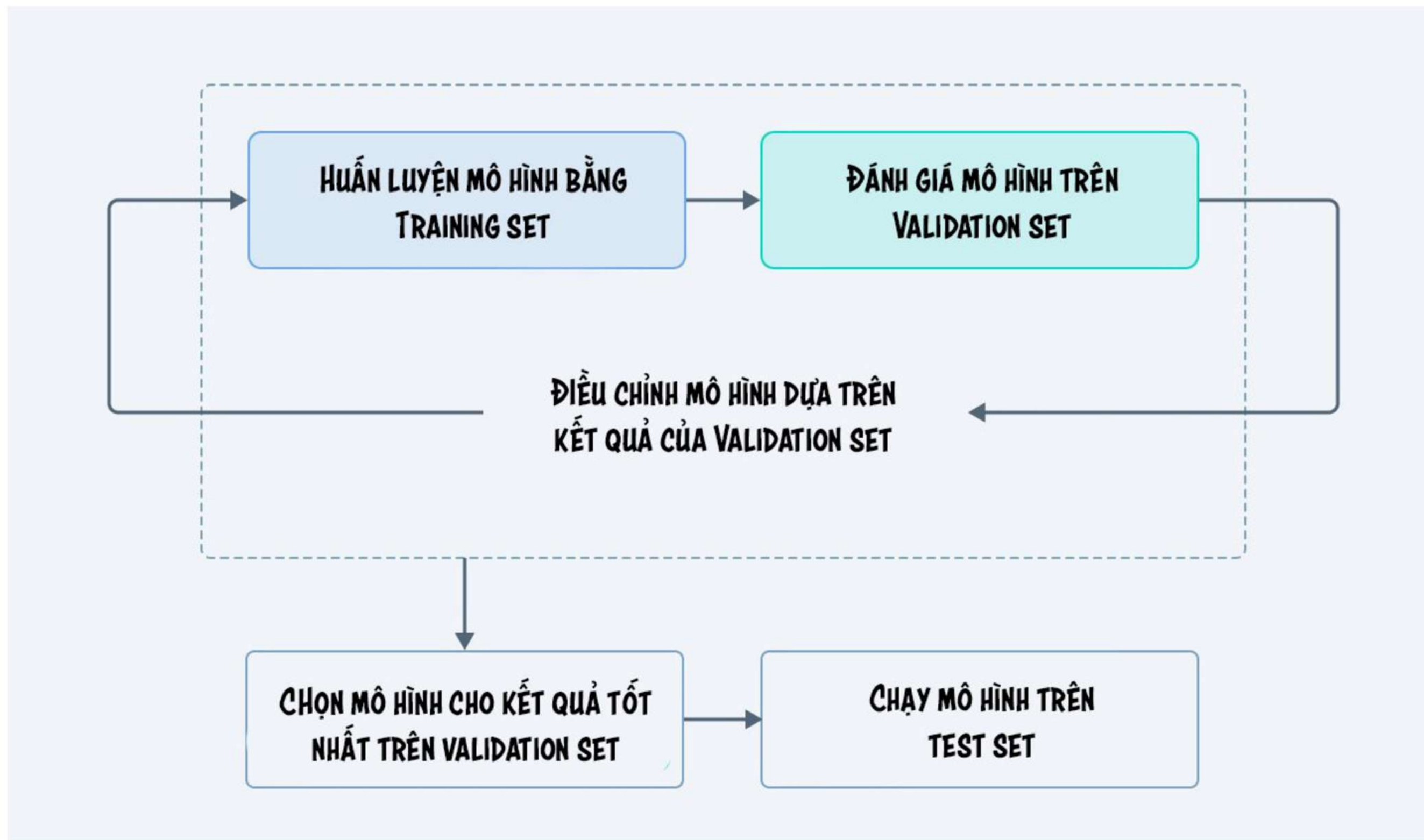
Overfitting

- Overfitting xảy ra khi mô hình học quá kỹ dữ liệu huấn luyện, đến mức ghi nhớ cả nhiễu và các đặc điểm không quan trọng. Ngược lại của overfitting là underfitting



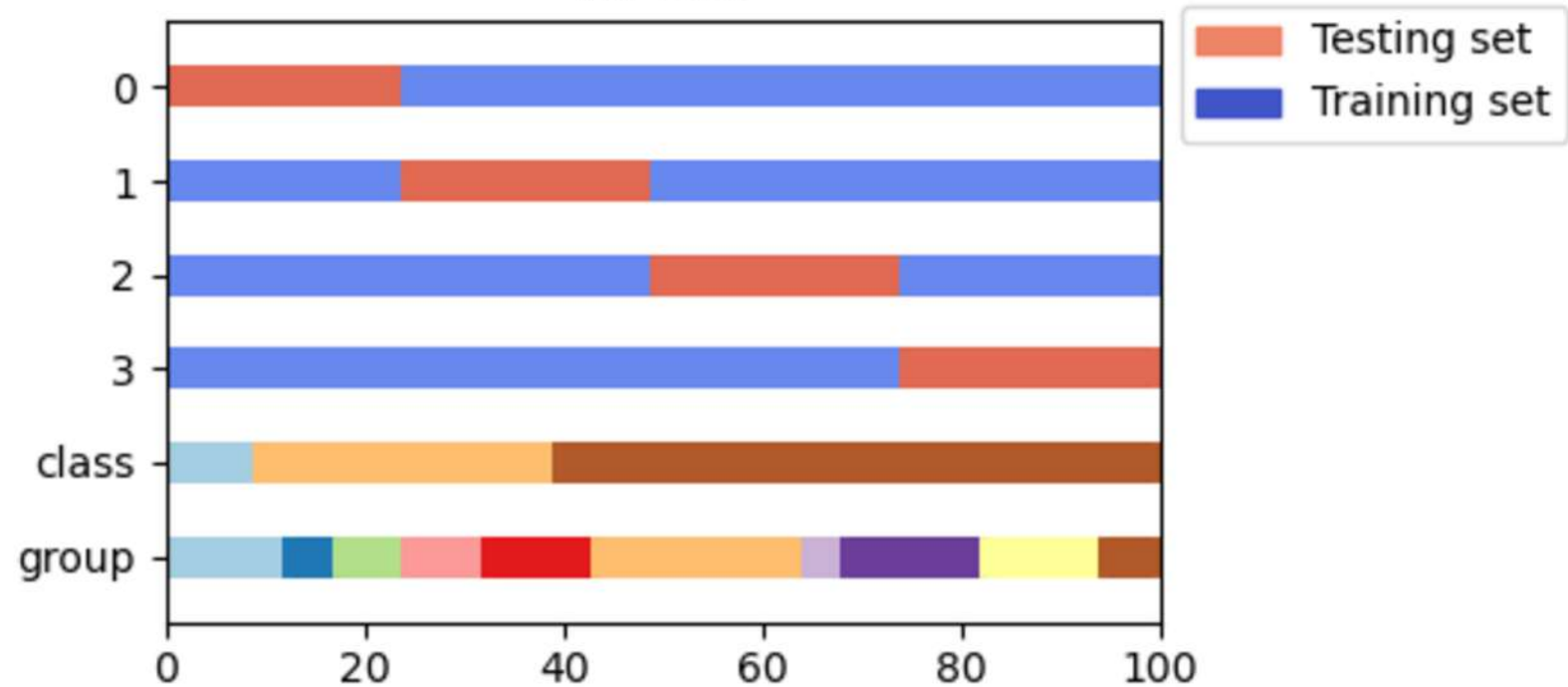
- Một số nguyên nhân dẫn đến overfitting:
 - Mô hình quá phức tạp
 - Dữ liệu huấn luyện ít hoặc không đủ bao quát
 - Huấn luyện quá lâu
 - Dữ liệu chưa qua xử lý, chứa nhiều nhiễu

Validation

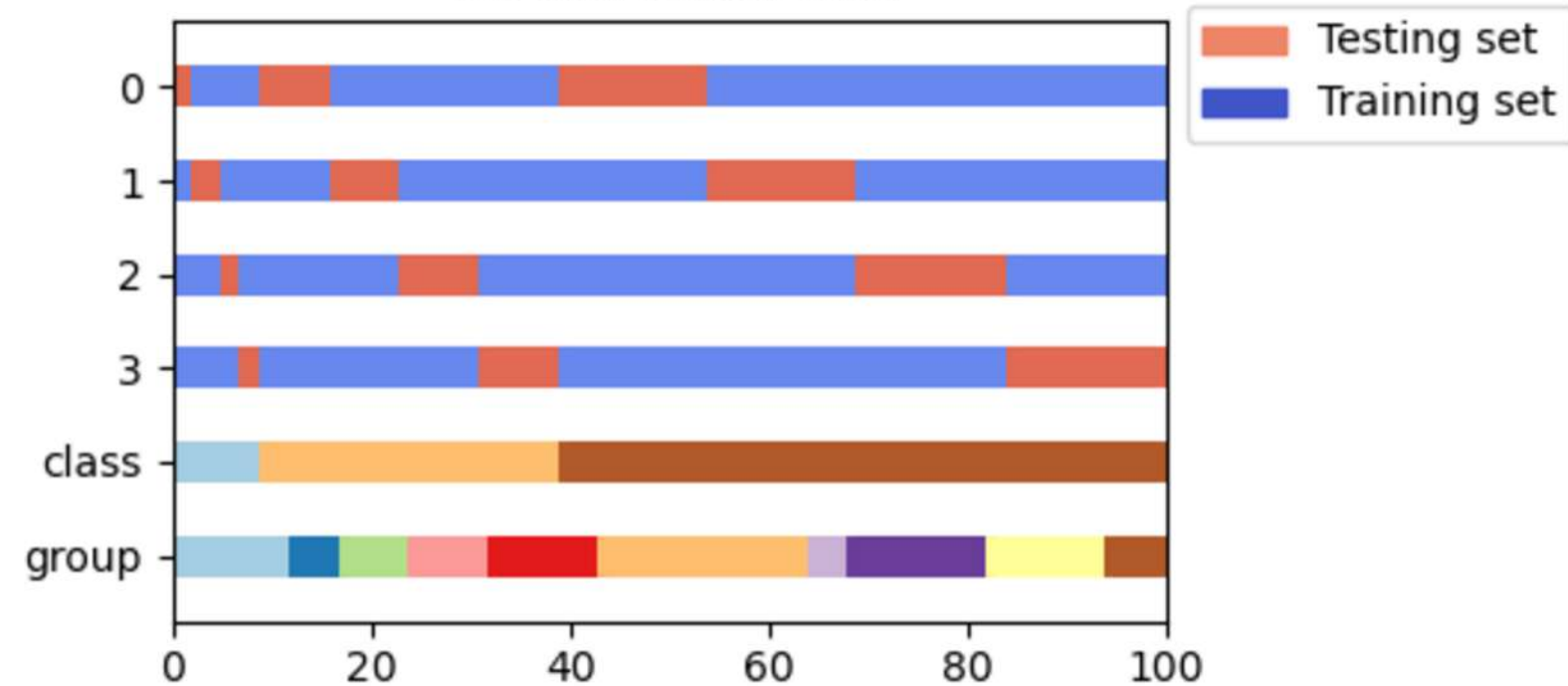


Cross-Validation

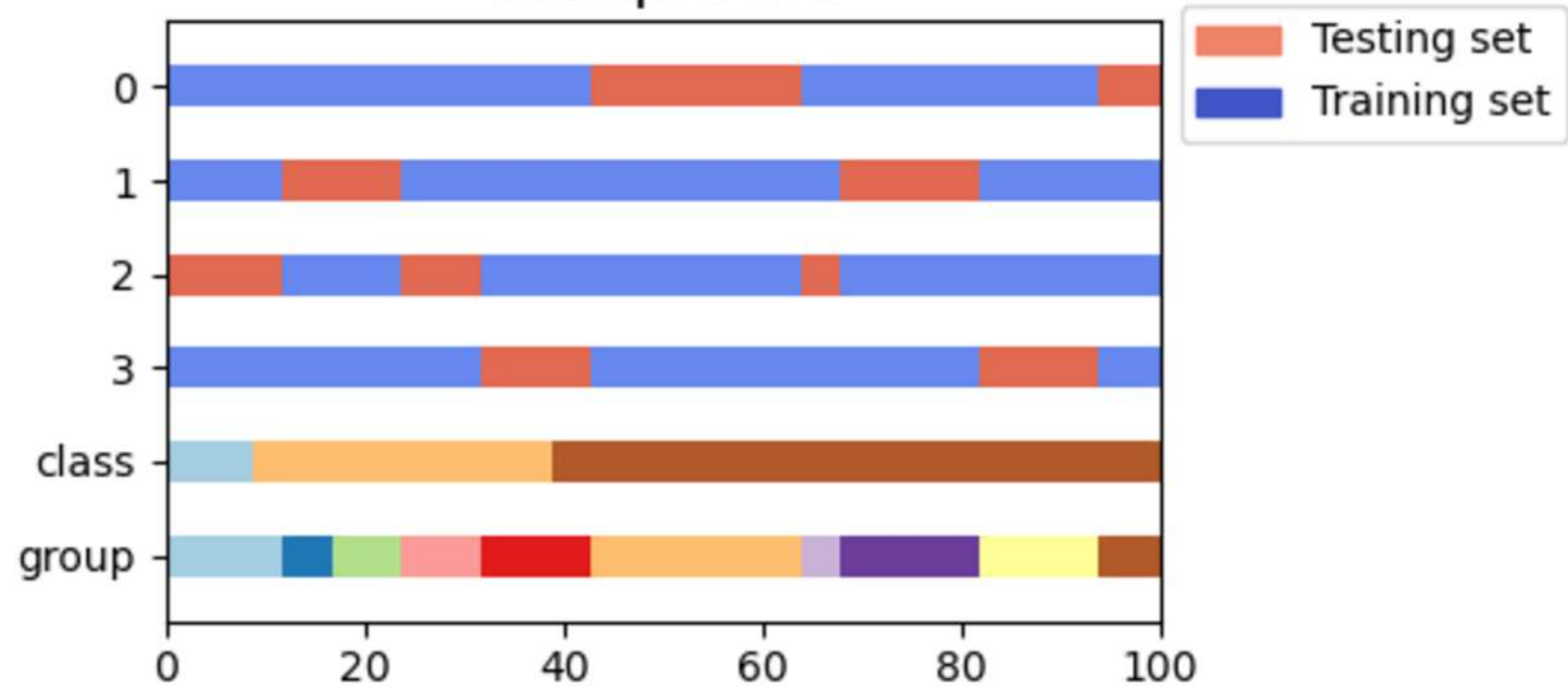
KFold



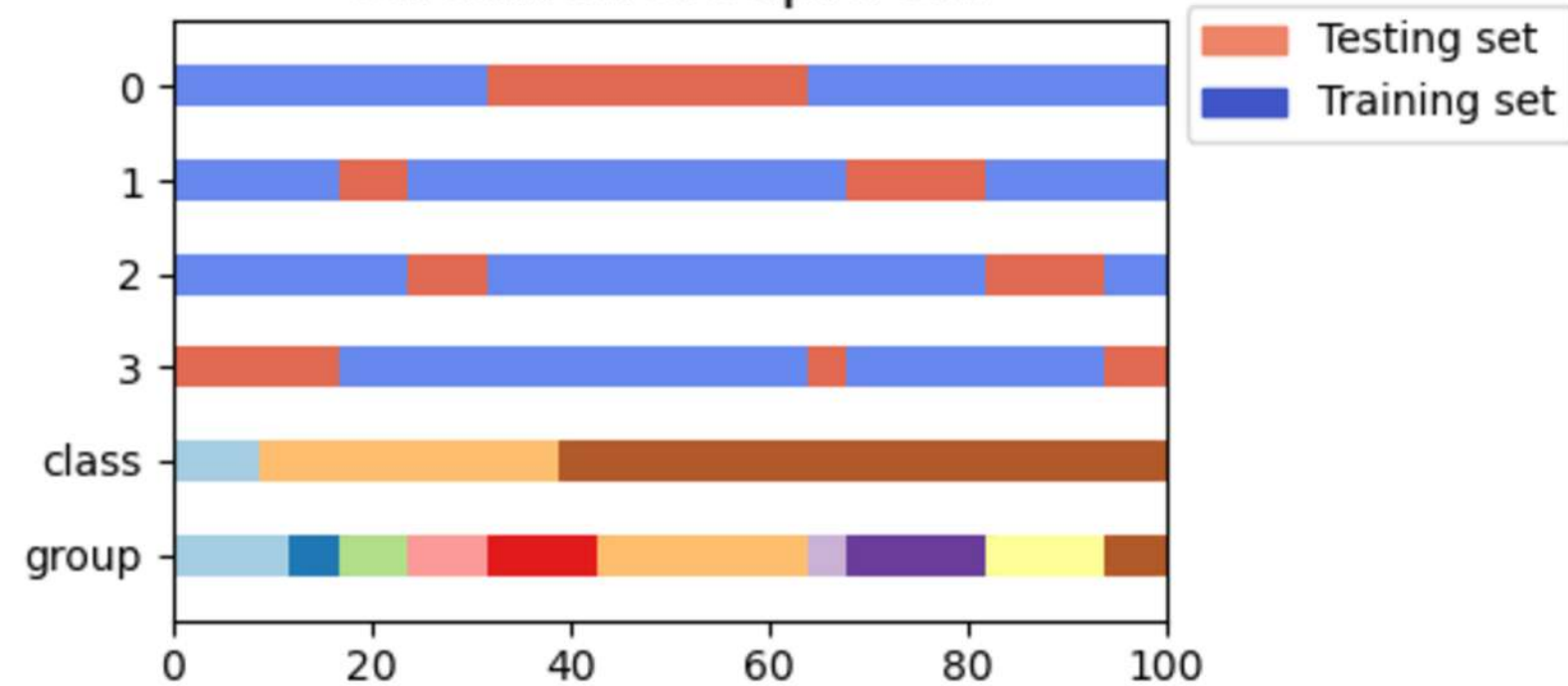
StratifiedKFold



GroupKFold



StratifiedGroupKFold

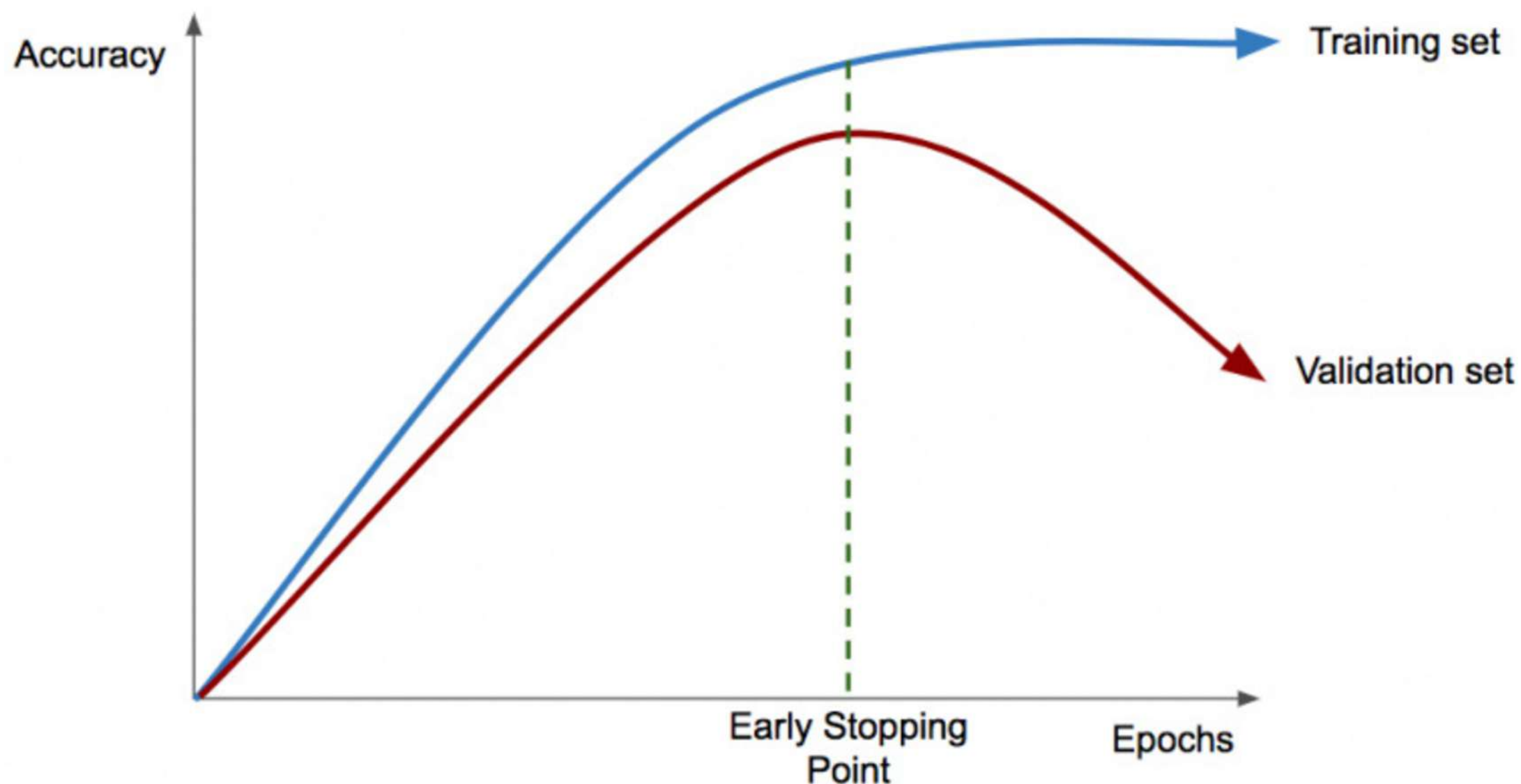


- Lợi ích của Validation/Cross-Validation:
 - Cho phép đánh giá mô hình trên dữ liệu chưa từng thấy
 - Giúp điều chỉnh siêu tham số, chọn kiến trúc mô hình phù hợp
 - Cross-Validation phù hợp khi dữ liệu ít, giúp tận dụng tối đa
- Hạn chế của Validation/Cross-Validation:
 - Nếu chỉ chia 1 lần, kết quả phụ thuộc vào cách chia. Nếu chia nhiều lần thì tốn thời gian tính toán
 - Cross-Validation không cẩn thận sẽ khiến rò rỉ dữ liệu

- Regularization là kỹ thuật “điều chuẩn” mô hình nhằm ngăn overfitting
- Regularization là bất kỳ sự điều chỉnh nào chúng ta thực hiện nhằm mục đích giảm sai số khái quát hóa (generalization error) nhưng không làm giảm sai số huấn luyện (training error)

Early Stopping

- Early stopping tức là dừng thuật toán trước khi hàm mất mát trên tập train giảm nhưng trên tập validation thì bắt đầu tăng, hoặc độ chính xác trên tập train tăng nhưng trên tập validation bắt đầu giảm



L1, L2 Regularization

- Trong huấn luyện mô hình, ta thêm một thuật ngữ phạt (penalty) vào hàm mất mát để giữ trọng số nhỏ và đơn giản, giúp mô hình tổng quát tốt hơn, tránh overfitting

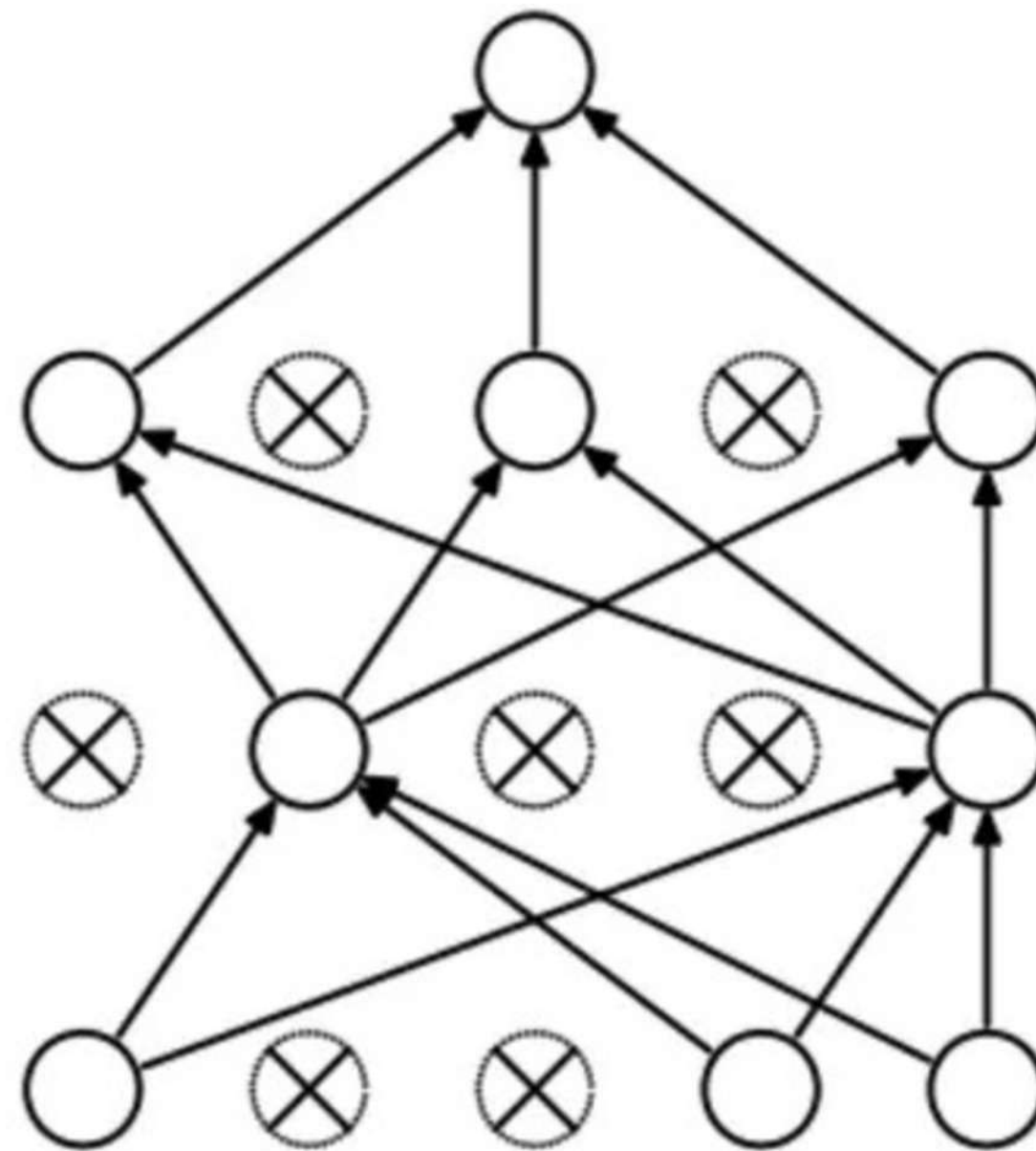
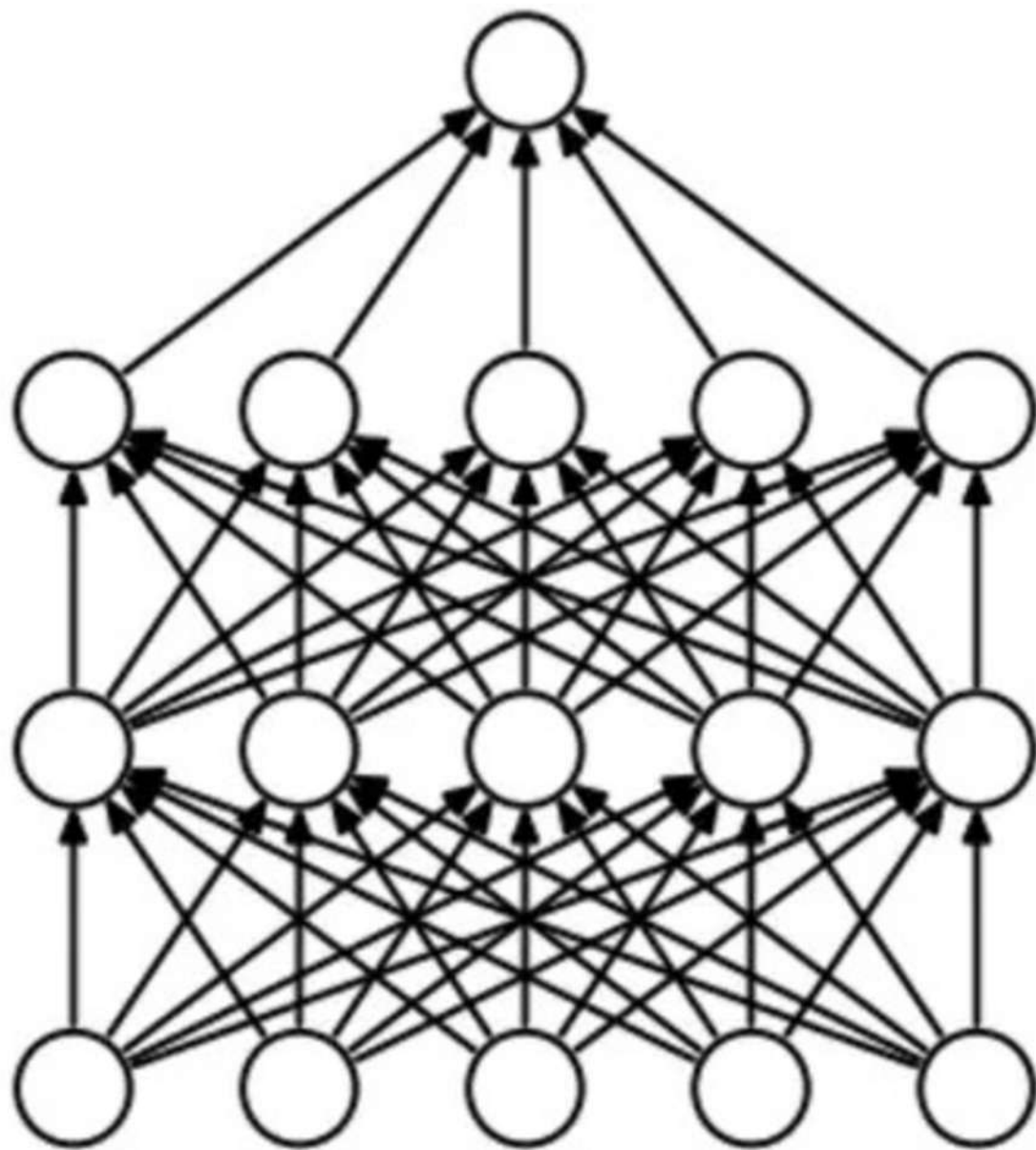
- L1 Regularization (Lasso): $J(\theta) = L(\theta) + \lambda \sum_i^n |\theta_i|$

- L2 Regularization (Ridge): $J(\theta) = L(\theta) + \lambda \sum_i^n \theta_i^2$

- L1 + L2 (Elastic Net): $J(\theta) = L(\theta) + \lambda_1 \sum_i^n |\theta_i| + \lambda_2 \sum_i^n \theta_i^2$

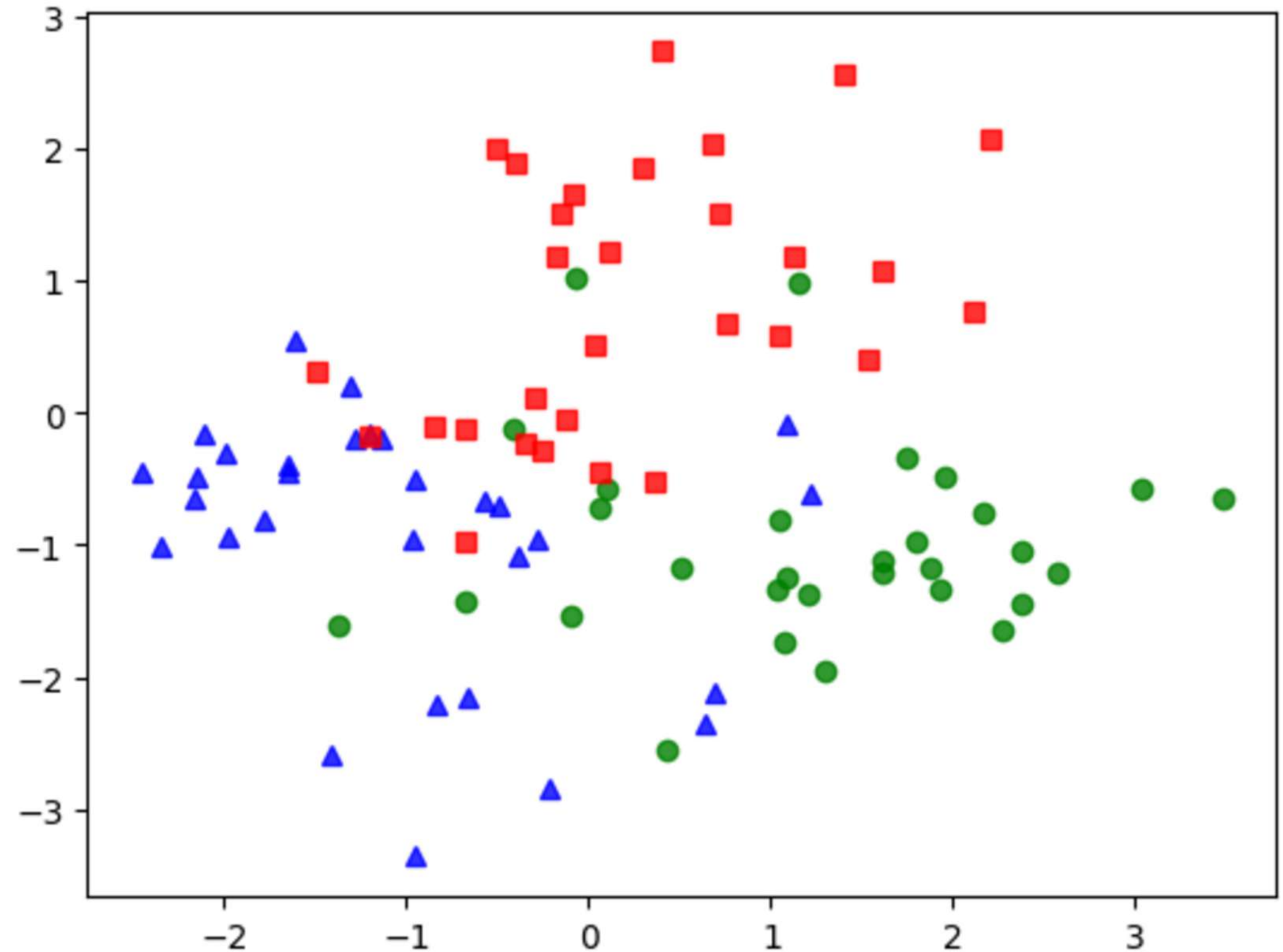
Dropout

- Dropout tức là tắt ngẫu nhiên một số node ở các lớp trong lúc huấn luyện để mạng học biểu diễn tổng quát hơn

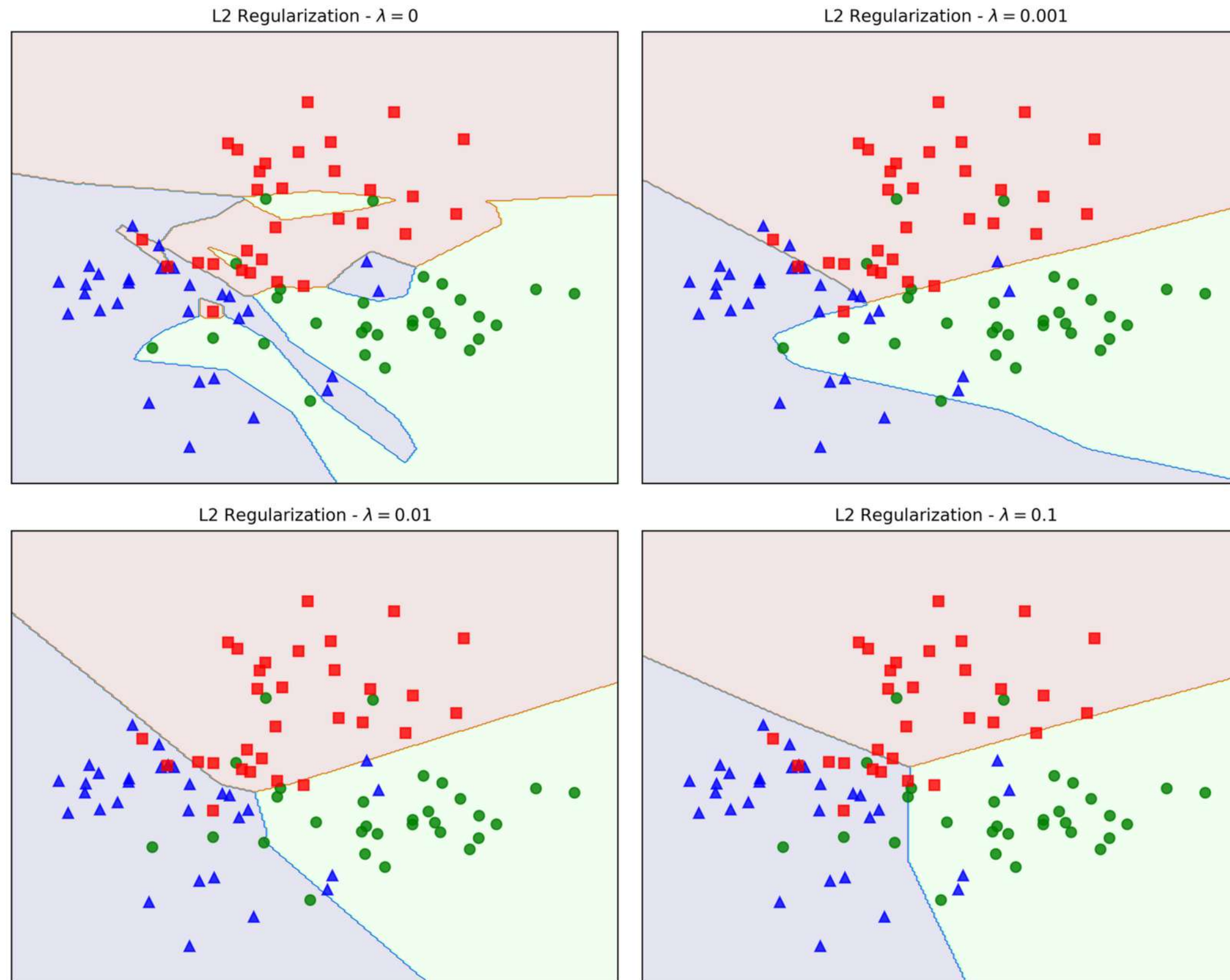


Ứng dụng trong bài toán thực tế

- Bộ dữ liệu: Customers Clustering
- Bộ dữ liệu gồm 3 nhóm khách hàng, mỗi nhóm gồm 30 mẫu bao gồm độ tuổi và số tiền chi tiêu



Ứng dụng trong bài toán thực tế



A large graphic on the left side of the slide. It features a dark blue background with a pattern of red dots of varying sizes arranged in concentric, slightly irregular circles, creating a sense of depth and movement. The word "HUST" is centered within this graphic.

HUST

THANK YOU !



hust.edu.vn



fb.com/dhbkhn