

# TỔNG QUAN VỀ MACHINE LEARNING

thipty@vietinbank.vn

Hà Nội, 09/2022

- ❖ Giới thiệu về học máy (machine learning)
- ❖ Các dạng học máy
- ❖ Machine learning workflow
- ❖ Một số thuật toán học có giám sát
- ❖ Một số thuật toán học không giám sát
- ❖ Các độ đo dùng trong đánh giá
- ❖ Các kỹ thuật đánh giá mô hình học máy
- ❖ Các thách thức chính của học máy

# Giới thiệu về Học máy (Machine Learning)

- ❖ *Học máy là khả năng của chương trình máy tính sử dụng kinh nghiệm, quan sát, hoặc dữ liệu trong quá khứ để cải thiện công việc của mình trong tương lai thay vì chỉ thực hiện theo đúng các quy tắc đã được lập trình sẵn. Chẳng hạn, máy tính có thể học cách dự đoán dựa trên các ví dụ, hay học cách tạo ra các hành vi phù hợp dựa trên quan sát trong quá khứ. [GS. Từ Minh Phương, Trí tuệ nhân tạo, HV BCVT]*



- ❖ Biểu diễn một bài toán học máy [Mitchell, 1997]:  
Học máy = Cải thiện hiệu quả một công việc thông qua kinh nghiệm
  - Thực hiện tốt hơn một công việc (nhiệm vụ) T
  - Theo tiêu chí đánh giá hiệu suất P
  - Thông qua (sử dụng) kinh nghiệm E

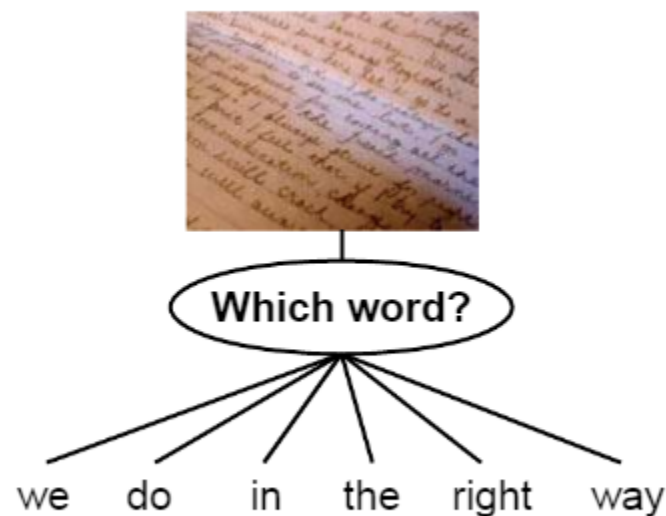
Bài toán lọc các trang Web theo sở thích của một người dùng

- **T**: Dự đoán (để lọc) xem những trang Web nào mà một người dùng cụ thể thích đọc
- **P**: Tỷ lệ (%) các trang Web được dự đoán đúng
- **E**: Một tập các trang Web mà người dùng đã chỉ định là thích đọc và một tập các trang Web mà anh ta đã chỉ định là không thích đọc



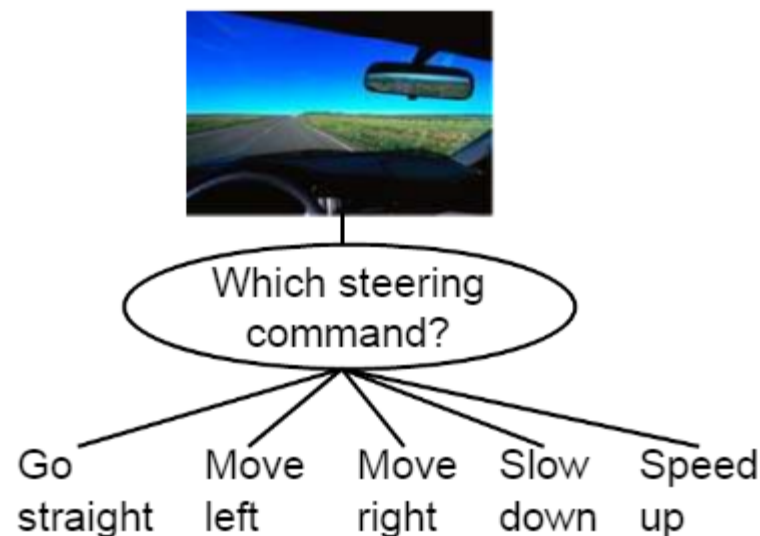
## Bài toán nhận dạng chữ viết tay

- **T**: Nhận dạng và phân loại các từ trong các ảnh chữ viết tay
- **P**: Tỷ lệ (%) các từ được nhận dạng và phân loại đúng
- **E**: Một tập các ảnh chữ viết tay, trong đó mỗi ảnh được gắn với một định danh của một từ



## Bài toán robot lái xe tự động

- **T:** Robot (được trang bị các camera quan sát) lái xe tự động trên đường cao tốc
- **P:** Khoảng cách trung bình mà robot có thể lái xe tự động trước khi xảy ra lỗi (tai nạn)
- **E:** Một tập các ví dụ được ghi lại khi quan sát một người lái xe trên đường cao tốc, trong đó mỗi ví dụ gồm một chuỗi các ảnh và các lệnh điều khiển xe



- Những ứng dụng khó lập trình theo cách thông thường do không tồn tại hoặc khó giải thích kinh nghiệm, kỹ năng của con người. Ví dụ:
  - Nhận dạng chữ viết tay, âm thanh, hình ảnh
  - Lái xe tự động, thám hiểm sao Hỏa
- Chương trình máy tính có khả năng thích nghi: lời giải thay đổi theo thời gian hoặc theo tình huống cụ thể. Ví dụ:
  - Chương trình trợ giúp cá nhân
  - Định tuyến mạng
- Khai phá (phân tích dữ liệu). Ví dụ:
  - Hồ sơ bệnh án --> tri thức y học
  - Dữ liệu bán hàng --> quy luật kinh doanh

- **Mẫu**, hay ví dụ (samples): là đối tượng cần xử lý (ví dụ xử lý phân loại)
  - Ví dụ: khi lọc thư rác thì mỗi thư là một mẫu, trong chuẩn đoán bệnh thì mẫu là bệnh nhân
- Mẫu thường được mô tả bằng tập thuộc tính hay **đặc trưng** (features)
  - ví dụ: trong chuẩn đoán bệnh, thuộc tính là triệu chứng của người bệnh, và các tham số khác như chiều cao, cân nặng, nhiệt độ, ...
- **Nhãn** phân loại (label): thể hiện loại của đối tượng mà ta cần dự đoán
  - ví dụ: nhãn phân loại thư rác có thể là 'rác' hoặc 'bình thường'



## ❖ Học có giám sát (Supervised learning)

Tập dữ liệu huấn luyện được cung cấp tường minh dưới dạng các mẫu cùng với giá trị đầu ra hay giá trị đích.

- Giá trị đầu ra là rời rạc thì gọi là *phân loại* hay *phân lớp* (classification).
- Giá trị đầu ra nhận giá trị *liên tục*, tức đầu ra là số thực, thì gọi là *hồi quy* (regression).

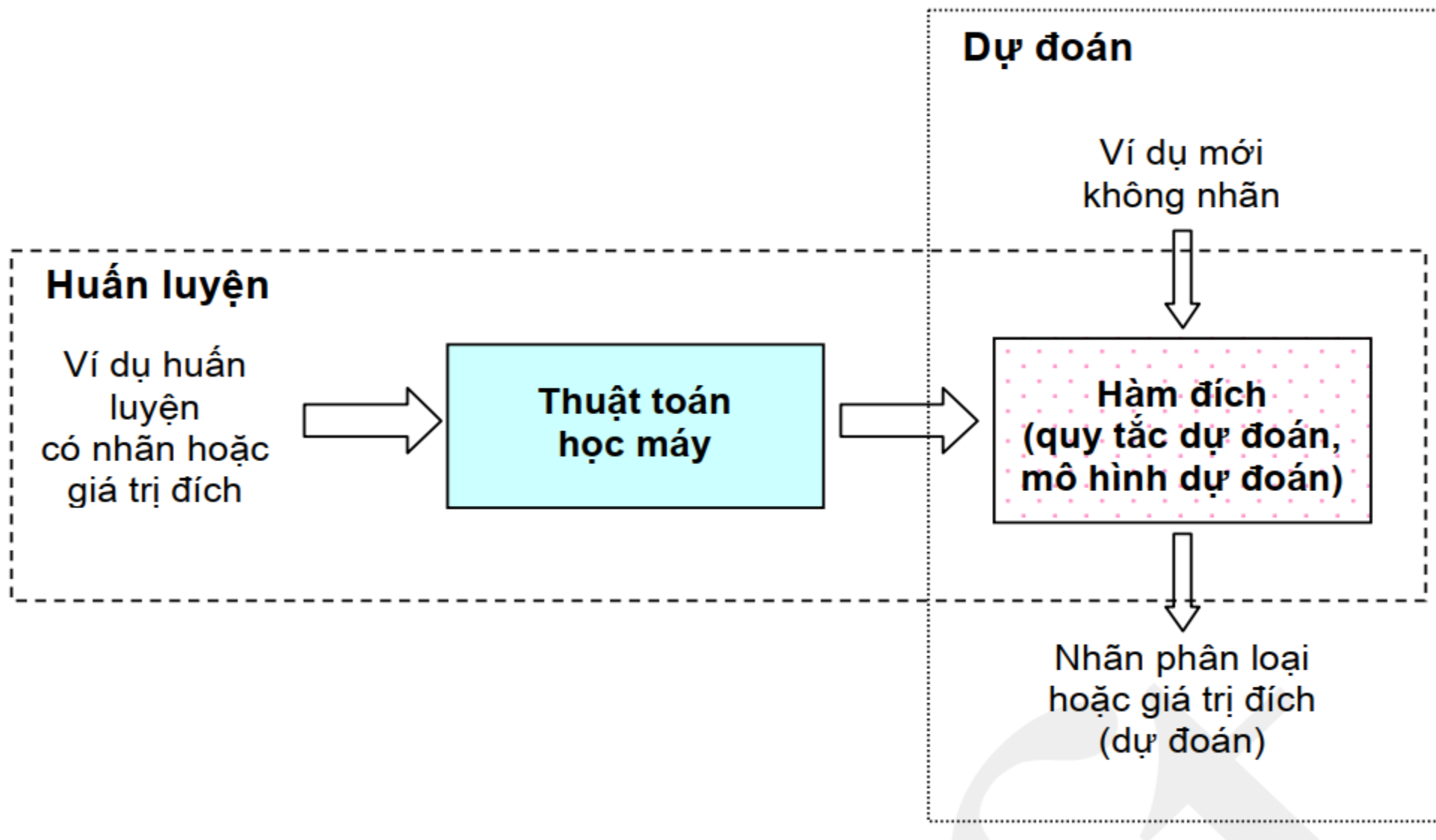
## ❖ Học không giám sát (Un-supervised learning)

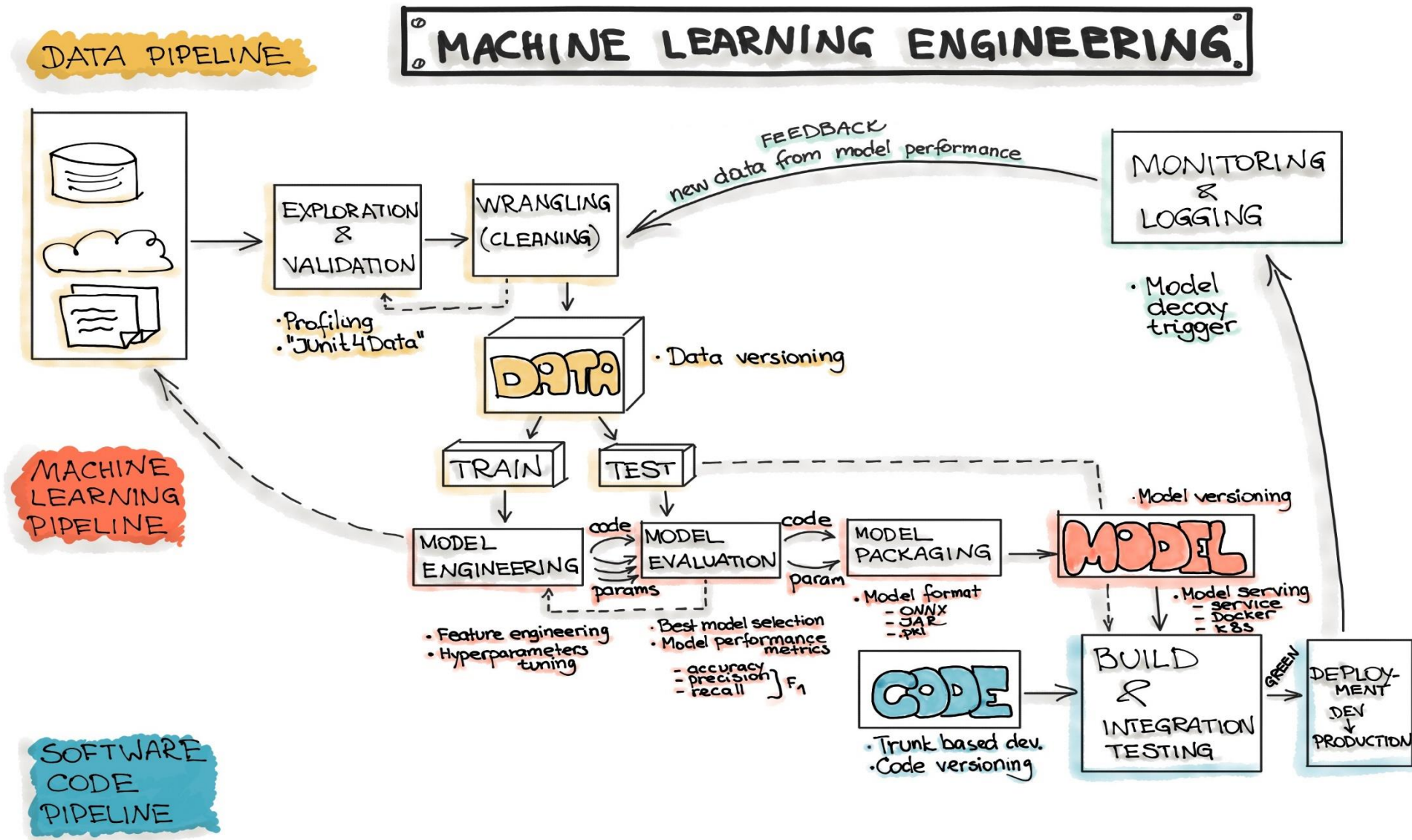
Các mẫu được cung cấp nhưng không có giá trị đầu ra hay giá trị đích

- Phân cụm (clustering): Dựa trên độ tương tự giữa các ví dụ để xếp chúng thành các nhóm, mỗi nhóm gồm các ví dụ tương tự nhau, khác với các ví dụ ở nhóm khác.
- Phát hiện luật kết hợp (association rule): Luật kết hợp có dạng  $P(A|B)$ , cho thấy xác suất hai tính chất A và B xuất hiện cùng với nhau.

## ❖ Học tăng cường (Reinforcement learning)

Kinh nghiệm không được cho trực tiếp dưới dạng đầu vào/đầu ra cho mỗi trạng thái hoặc mỗi hành động. Hệ thống nhận được một giá trị thưởng (reward) là kết quả cho một chuỗi hành động nào đó. Thuật toán cần học cách hành động để cực đại hóa giá trị thưởng.





# Một số thuật toán học có giám sát - Linear Regression

Mô hình hồi quy tuyến tính có dạng sau:

$$h(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

$$h(x) = \sum_{i=0}^n \beta_i x_i$$

Hàm Lỗi có dạng sau:

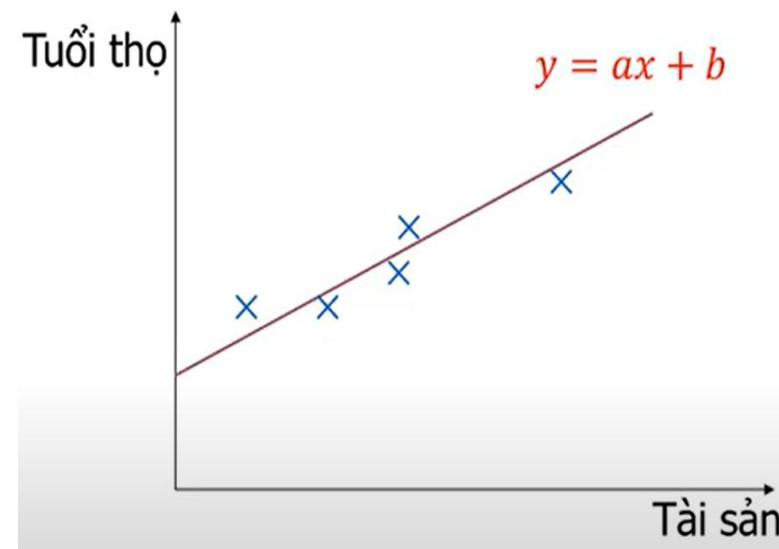
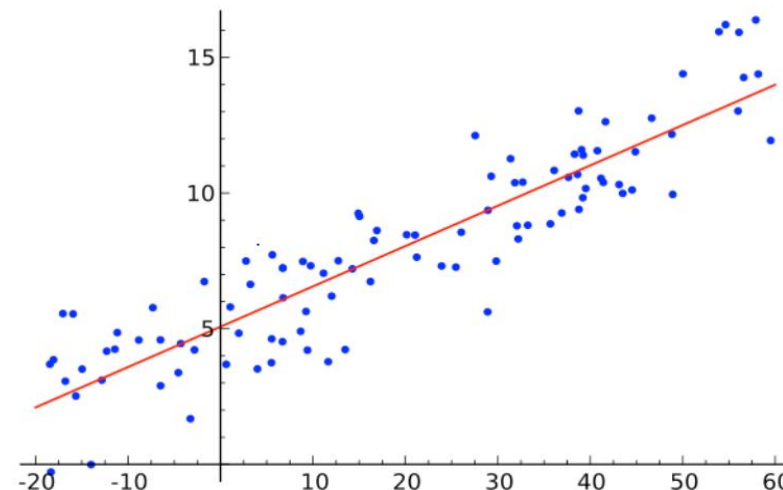
$$J(\beta) = \frac{1}{2} \sum_{i=1}^m (h(x_i) - y_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=0}^n \beta_j x_{ij} - y_i \right)^2$$

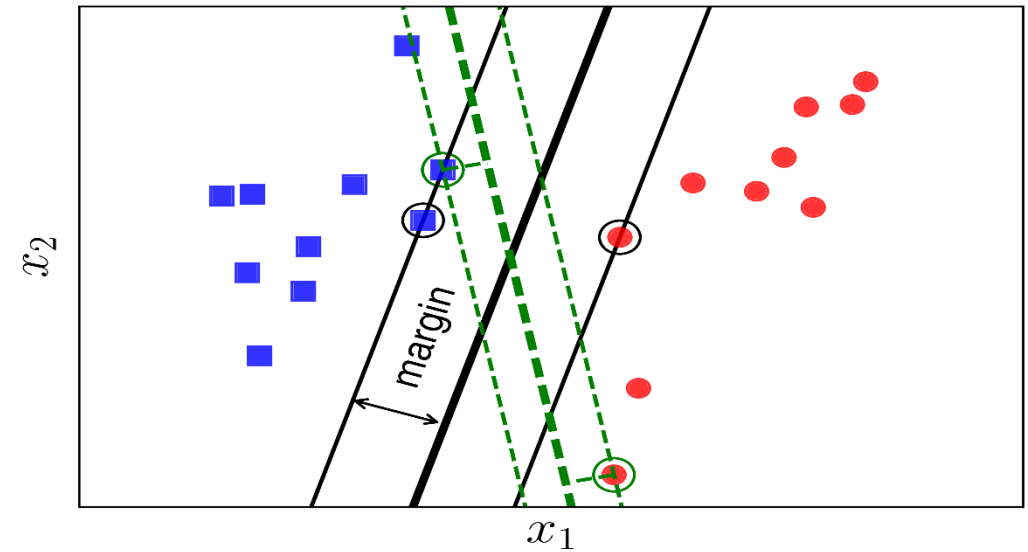
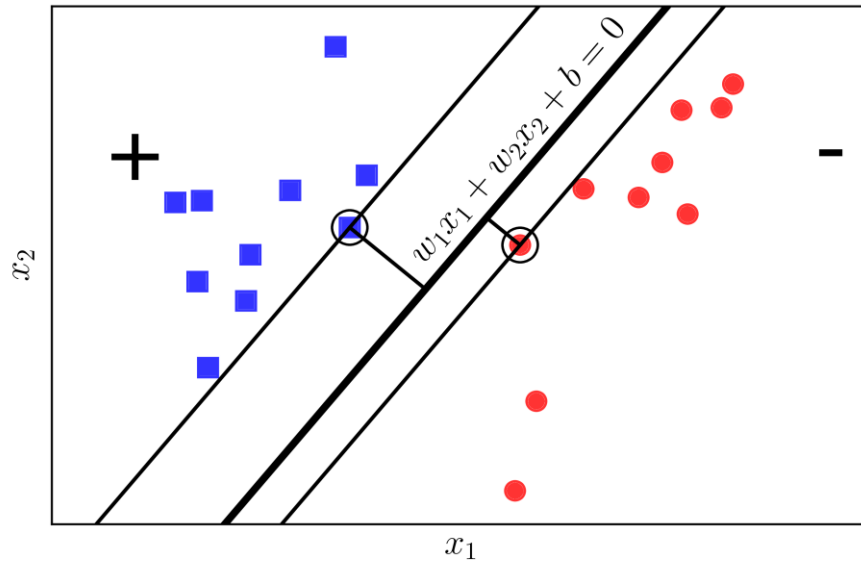
$$J(\beta) = \frac{1}{2} (\mathbf{X}\beta - \mathbf{y})^T (\mathbf{X}\beta - \mathbf{y})$$

Tính đạo hàm  $J$  theo  $\beta$  và cho đạo hàm bằng không:  $\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) = 0$

Lời giải như sau:  $\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$



# Một số thuật toán học có giám sát - SVM

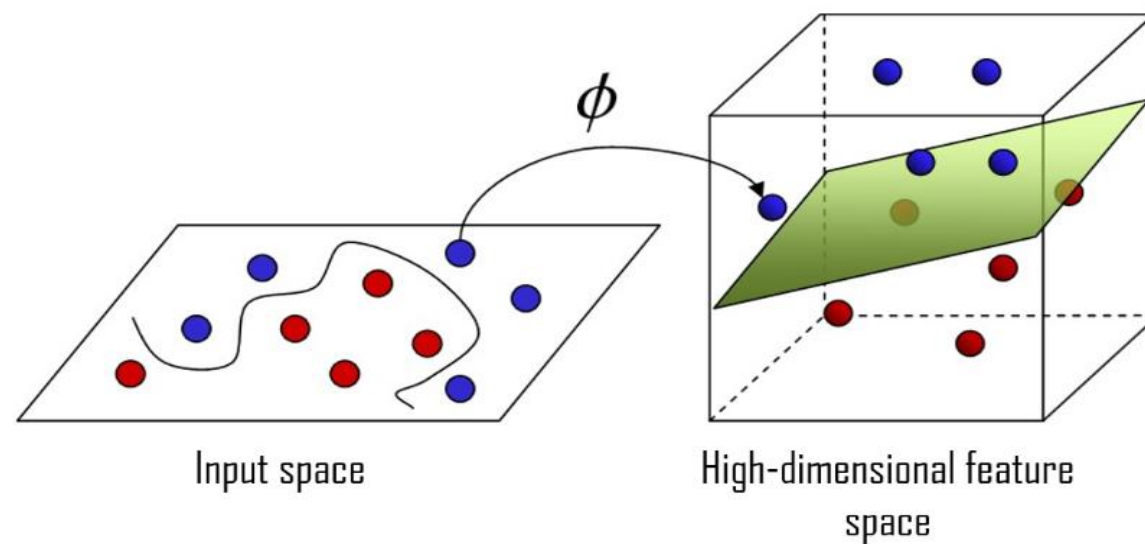
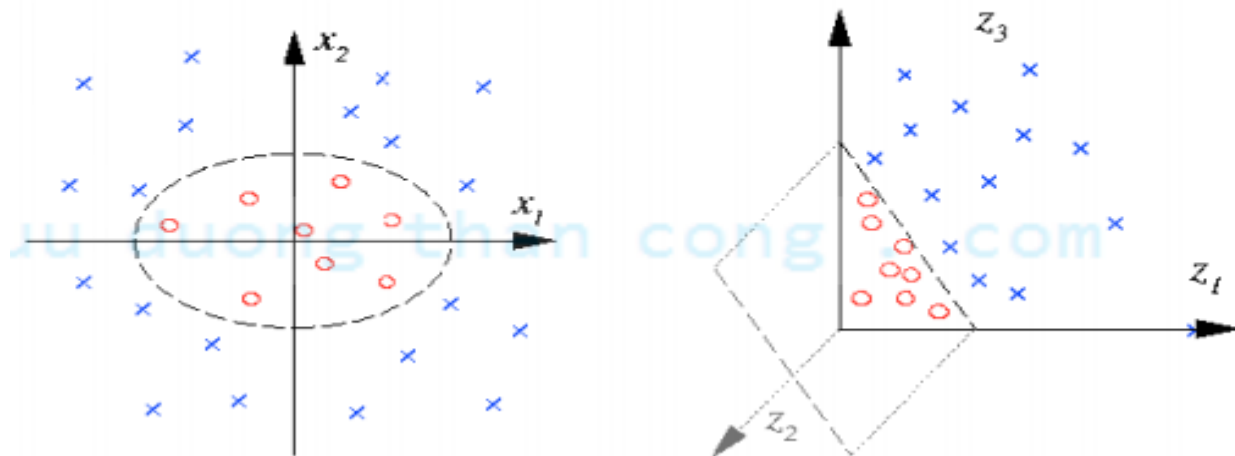


Hàm phân loại có dạng: 
$$f(\mathbf{x}) = \sum_{i=1}^d w_i x_i + b = \mathbf{w}^T \mathbf{x} + b$$

Nhãn phân loại  $y$  được xác định như sau:  $y = \text{sgn}(f(\mathbf{x}))$

# Một số thuật toán toán học có giám sát - SVM

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$
$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$



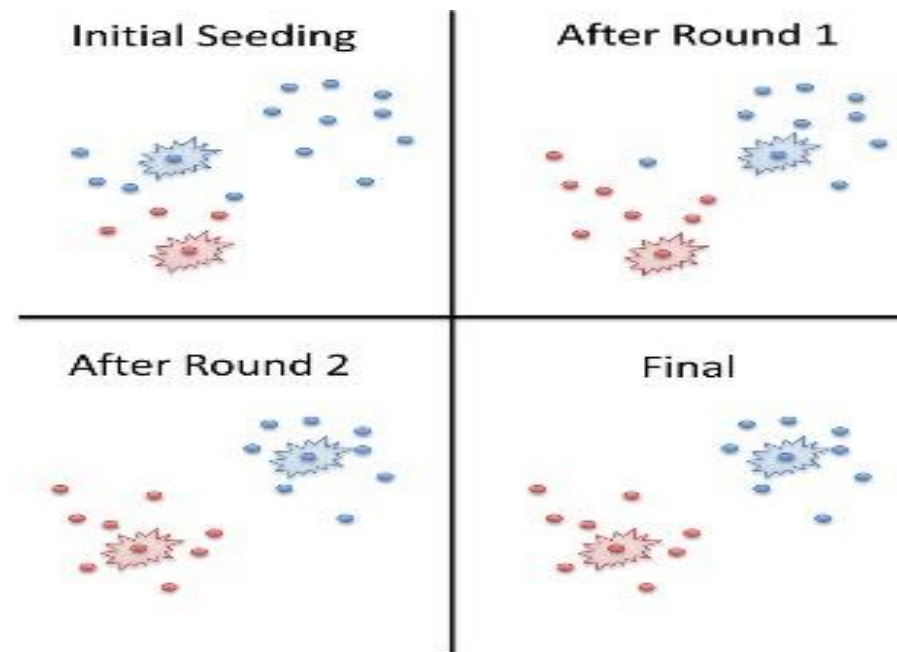


# Một số thuật toán học không giám sát - Kmean

**Đầu vào:** Dữ liệu X và số lượng cluster K.

**Đầu ra:** Các center M và label vector cho từng điểm dữ liệu Y.

1. Chọn K điểm bất kỳ làm các center ban đầu.
2. Phân mỗi điểm dữ liệu vào cluster có center gần nó nhất.
3. Nếu việc gán dữ liệu vào từng cluster ở bước 2 không thay đổi so với vòng lặp trước nó thì ta dừng thuật toán.
4. Cập nhật center cho từng cluster bằng cách lấy trung bình cộng của tất cả các điểm dữ liệu đã được gán vào cluster đó sau bước 2.
5. Quay lại bước 2.



# Các độ đo đánh giá - Độ đo dùng trong classification

Nhấn thật	Nhấn dự đoán	
	Positive	Negative
Positive	True Positive (TP)	False Negative (FN)
Negative	False Negative (FN)	True Negative (TN)

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$F - \text{measure} = \frac{2 * \text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$



**Xét ví dụ thực tế :** Vào ngày 16/04/2021 tại bệnh viện A có 100 bệnh nhân đến khám một loại bệnh, giả sử biết trước trong 100 bệnh nhân có 60 người mắc bệnh, 40 người không có bệnh. Sau khi thăm khám, bệnh viện đưa ra kết quả:

- Trong 60 người bệnh thật thì có 45 người chuẩn đoán có bệnh, 15 người chuẩn đoán không mắc bệnh.
- Trong 40 người không mắc bệnh thì có 30 người chuẩn đoán không mắc bệnh, 10 người chuẩn đoán là mắc bệnh.

→ TP = ?

→ FN = ?

→ TN = ?

→ FP = ?

→ Precision = ?

→ Recall = ?

→ F-measure = ?

→ Accuracy = ?

**Xét ví dụ thực tế :** Vào ngày 16/04/2021 tại bệnh viện A có 100 bệnh nhân đến khám một loại bệnh, giả sử biết trước trong 100 bệnh nhân có 60 người mắc bệnh, 40 người không có bệnh. Sau khi thăm khám, bệnh viện đưa ra kết quả:

- Trong 60 người bệnh thật thì có 45 người chuẩn đoán có bệnh, 15 người chuẩn đoán không mắc bệnh.
- Trong 40 người không mắc bệnh thì có 30 người chuẩn đoán không mắc bệnh, 10 người chuẩn đoán là mắc bệnh.

$$\rightarrow TP = 45$$

$$\rightarrow FN = 15$$

$$\rightarrow TN = 30$$

$$\rightarrow FP = 10$$

$$\rightarrow \text{Precision} = TP / (TP + FP) = 45 / (45 + 10) = 45 / 55 = 0.82$$

$$\rightarrow \text{Recall} = TP / (TP + FN) = 45 / (45 + 15) = 45 / 60 = 0.75$$

$$\rightarrow \text{F-measure} = 2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall}) = 2 * 0.82 * 0.75 / (0.82 + 0.75) = 0.78$$

$$\rightarrow \text{Accuracy} = (TP + TN) / (TP + TN + FP + FN) = (45 + 30) / (45 + 30 + 10 + 15) = 75 / 100 = 0.75$$

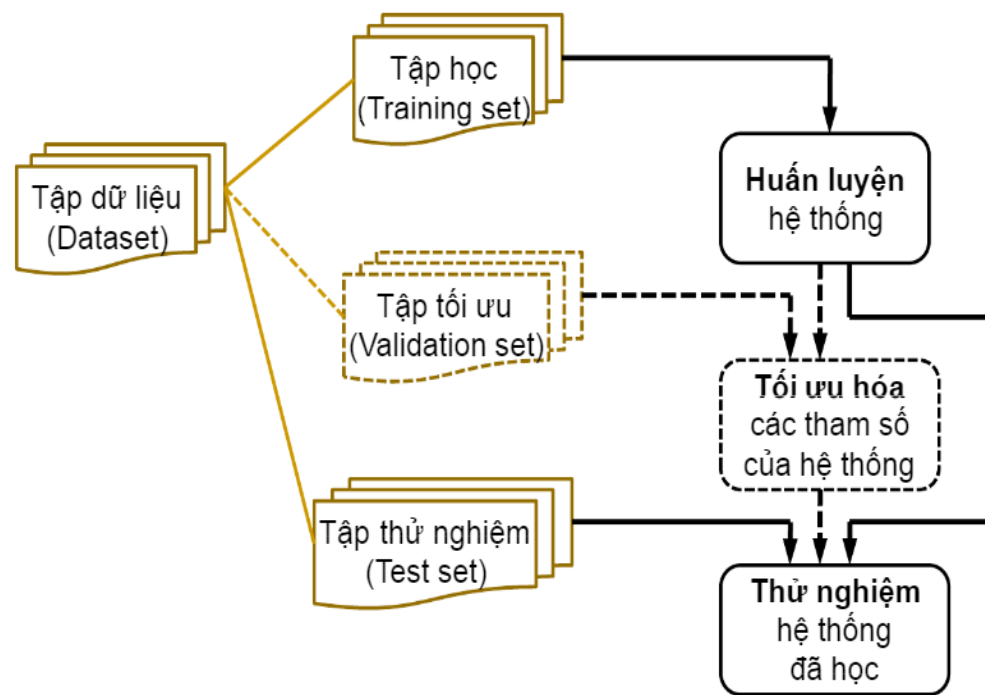
Hai độ đo thường dùng trong bài toán Regression là *lỗi trung bình tuyệt đối*  $MAE$  (Mean absolute error), và *lỗi trung bình bình phương*  $MSE$  (Mean squared error)

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}|$$

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

# Các kỹ thuật đánh giá mô hình (cross validation)

- Kiểm tra chéo đơn giản với tập thử nghiệm tách riêng (hold-out cross validation)
- Kiểm tra chéo k-fold (k-fold cross validation)



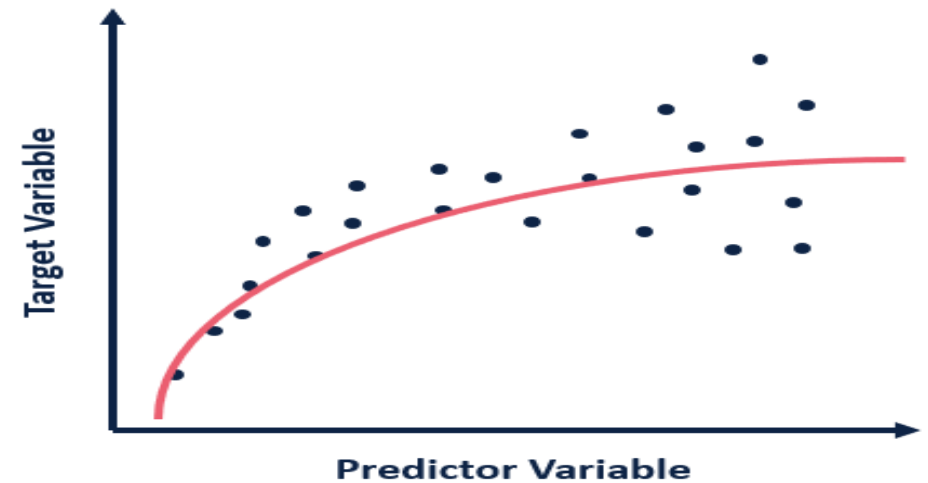
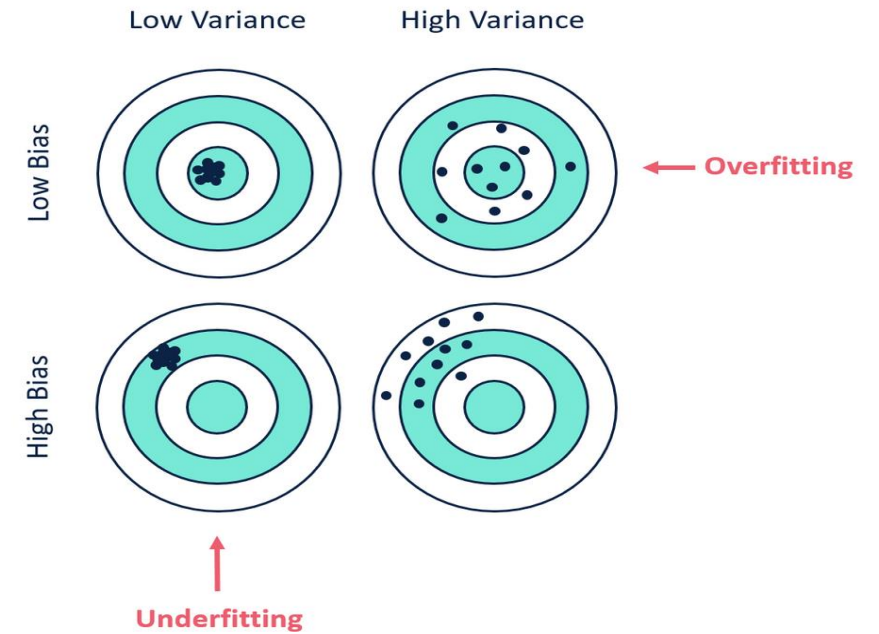
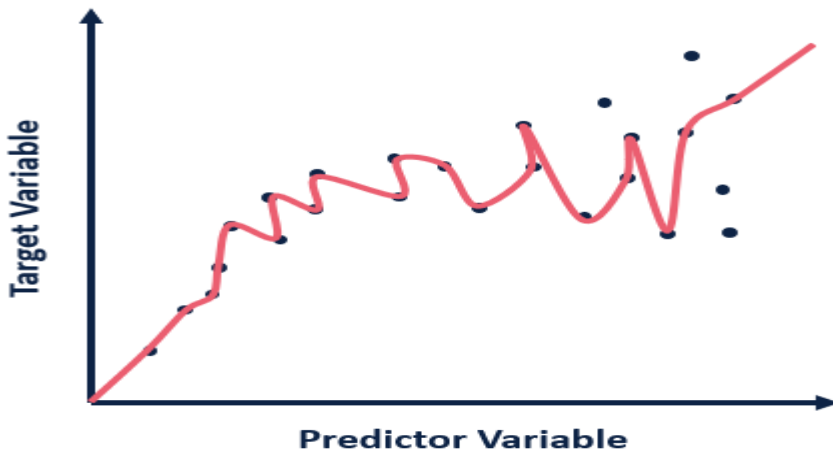
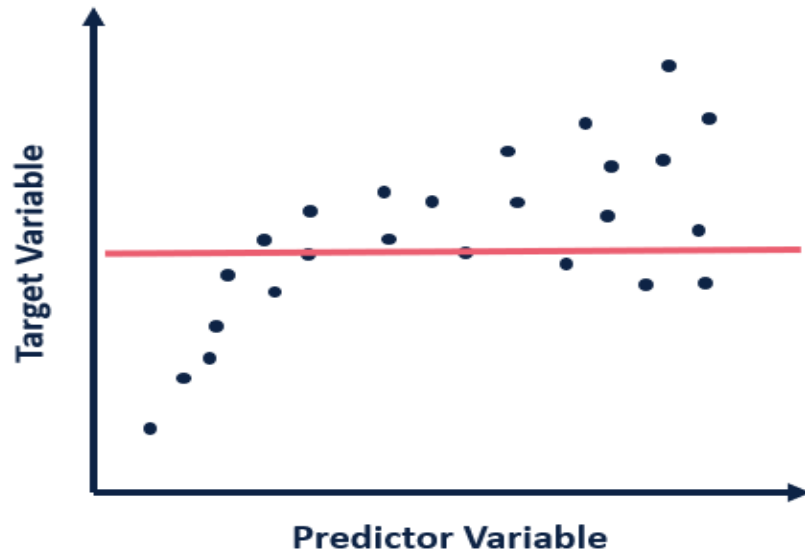
## ❖ Dữ liệu chưa tốt

- Không đủ dữ liệu để huấn luyện mô hình (Insufficient Quantity of Training Data)
- Dữ liệu dùng cho huấn luyện không mang tính biểu diễn (Nonrepresentative Training Data)
- Chất lượng dữ liệu thấp (Poor-Quality Data)
- Các đặc trưng không liên quan (Irrelevant Features)

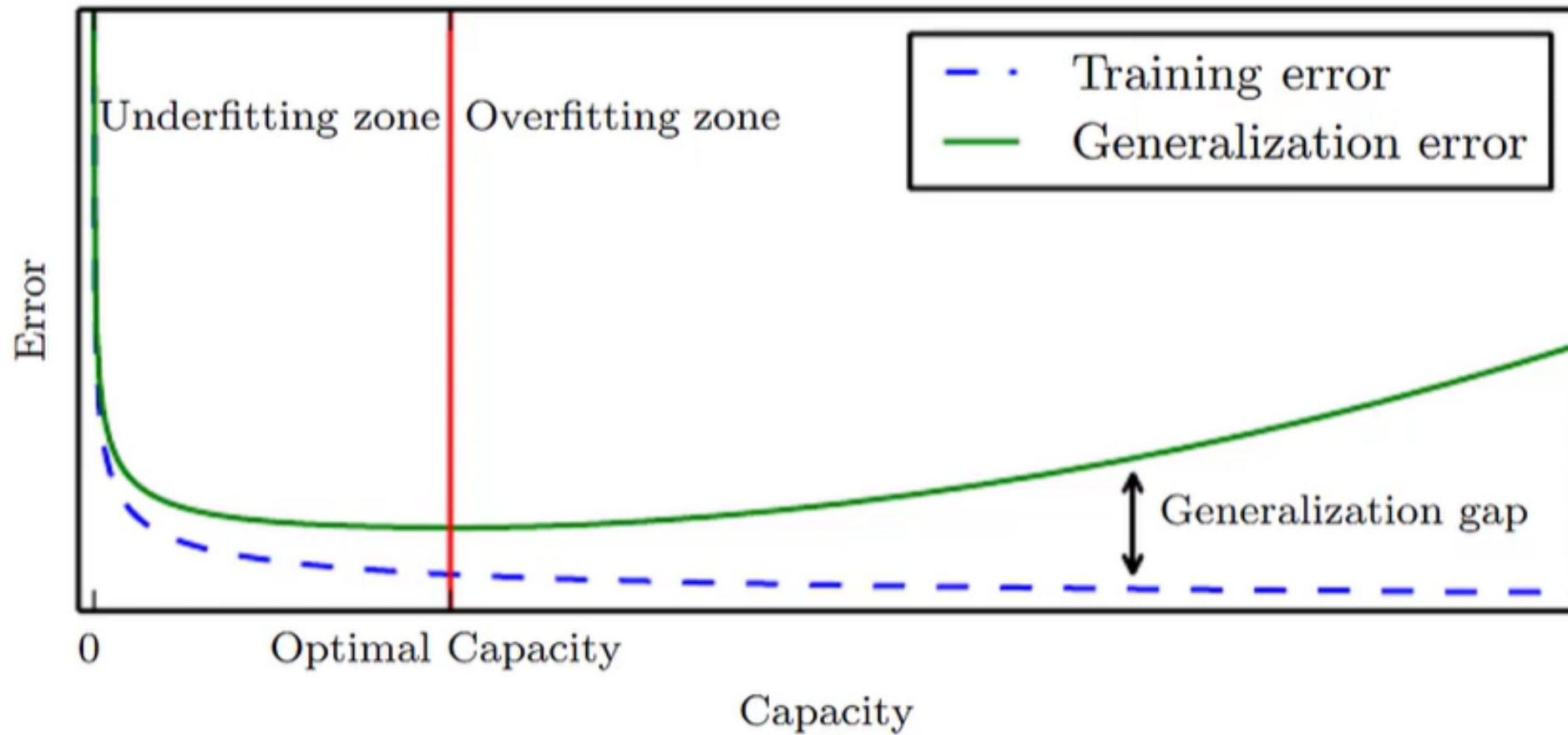
## ❖ Thuật toán chưa tốt

- Underfitting
- Overfitting

# Các thách thức chính của học máy – underfitting và Overfitting



# Các thách thức chính của học máy – underfitting và Overfitting



Từ Minh Phương, Nhập môn trí tuệ nhân tạo, Học viện Bưu chính Viễn thông

<https://machinelearningcoban.com/>

<https://www.researchgate.net/>

<https://developers.google.com/machine-learning/crash-course>

<https://ml-ops.org/>

<https://community.alteryx.com/datascience>



**THANK YOU!**