



データマイニング

Data Mining

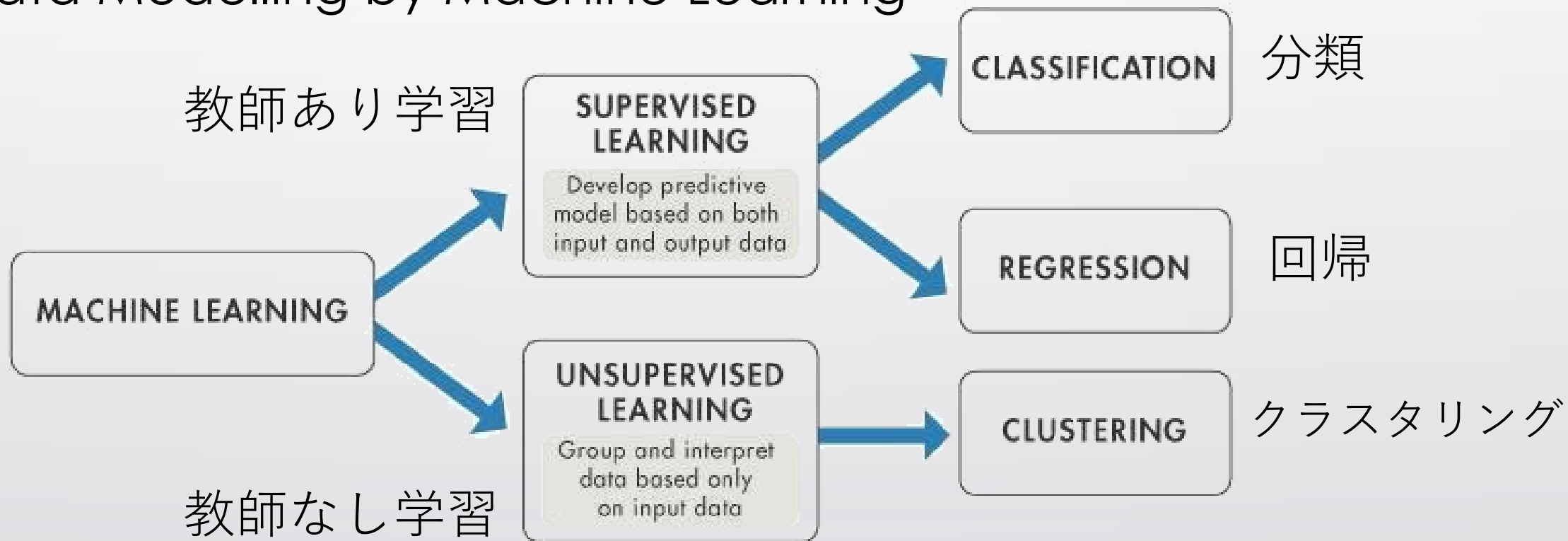
6: 分類① Classification

土居 裕和 Hirokazu Doi

長岡技術科学大学 Nagaoka University of Technology

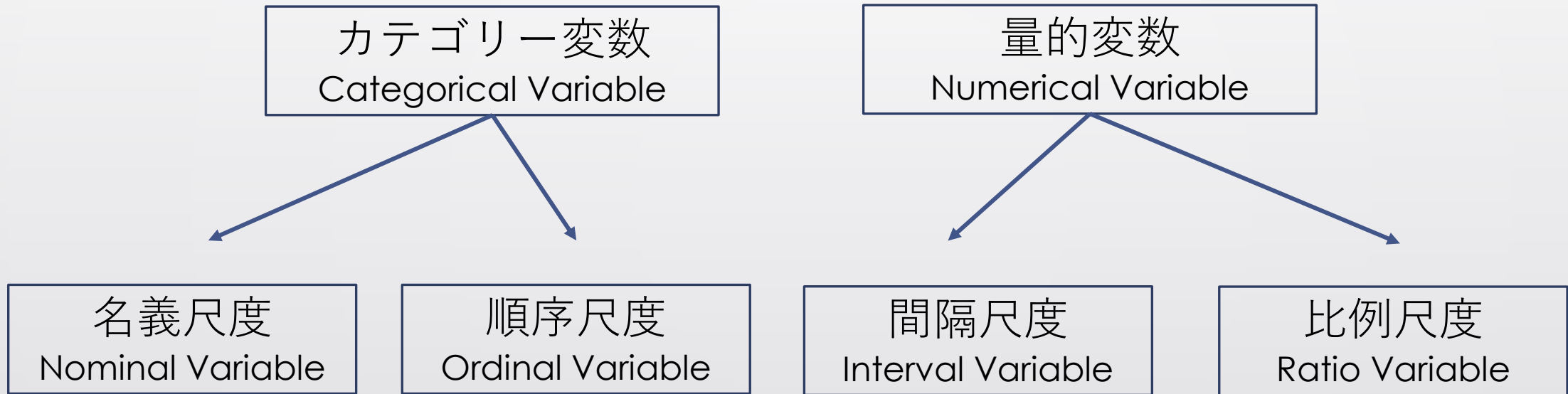
機械学習によるモデル化

Data Modelling by Machine Learning



Supervised Learning versus Unsupervised Learning (Mathworks, n.d.)

変数の種類 Types of Variables



変数の種類 Types of Variables

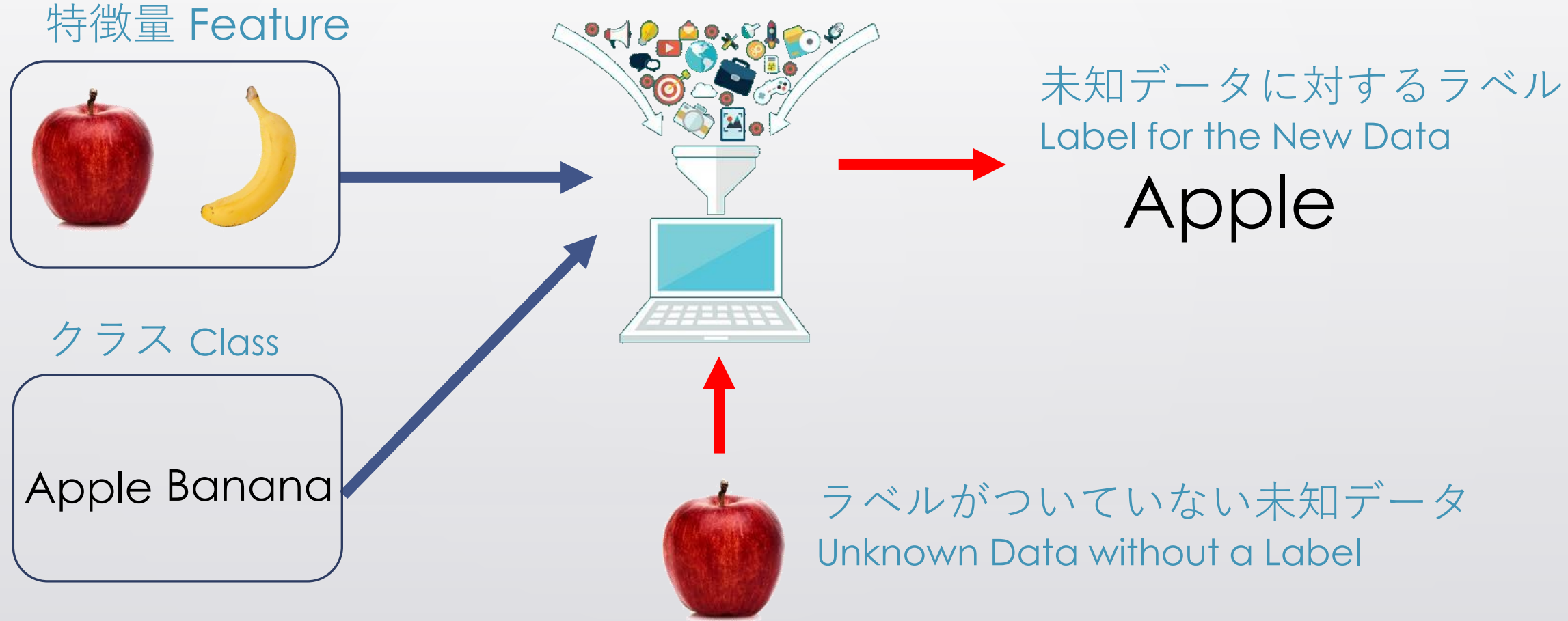
名義尺度
Nominal Variable

あるカテゴリーを、別のカテゴリーと区別するために用いられる、数値自体には意味がない変数

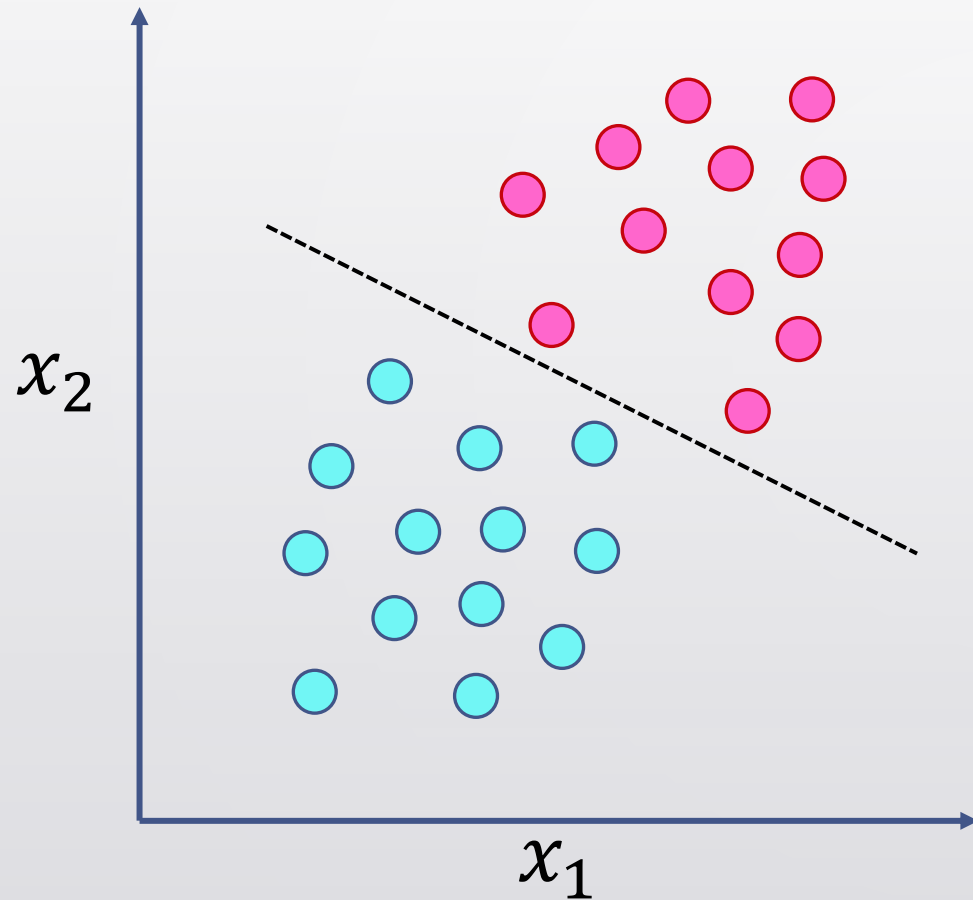
Variables, whose number has no numerical value, often used to discriminate multiple categories

ユーザーID	性別	年齢	年収
sanapon	男	26	411万円
oggi1985	女	33	536万円
murachan	女	39	681万円
shozan.s	男	23	309万円

教師あり学習 Supervised Learning



線型判別分析 Linear Discriminant Analysis (LDA)

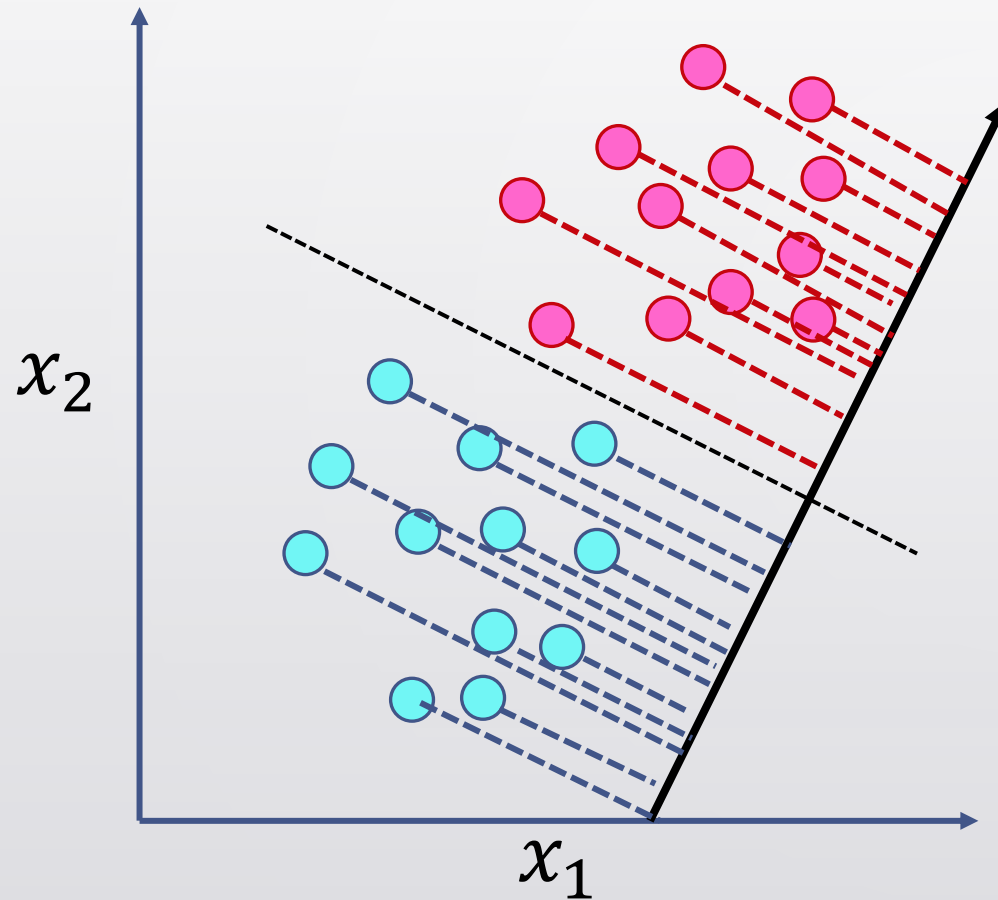


クラス1とクラス2を分離できる決定境界の引き方を見つける

Find how to draw a decision boundary that separates Class 1 and Class 2

決定境界
Decision Boundary

線型判別分析 Linear Discriminant Analysis (LDA)



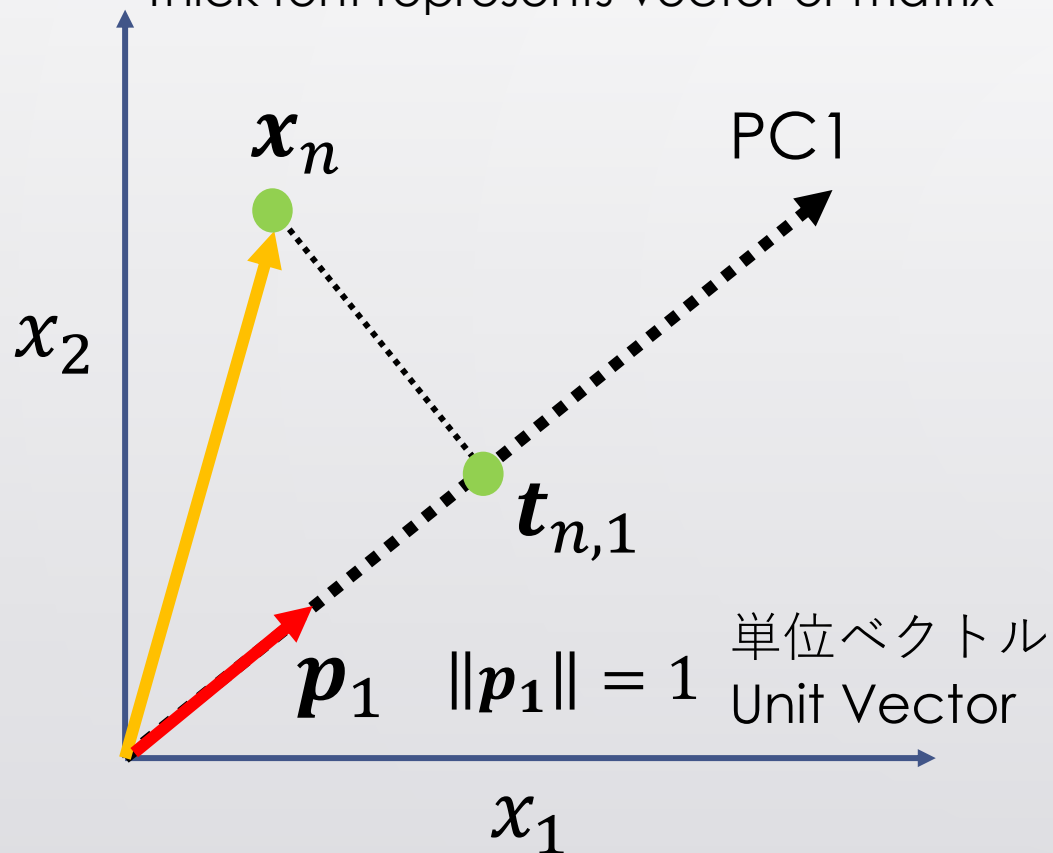
決定境界に直交する軸へのデータの射影を計算する

Consider the projection of data onto axis orthogonal to the decision boundary

第1主成分の計算 Computation of PC1

太字はベクトルか行列

Thick font represents vector or matrix



変数を標準化しておく

Normalize the variables

観測データ x_n の第1主成分軸方向への

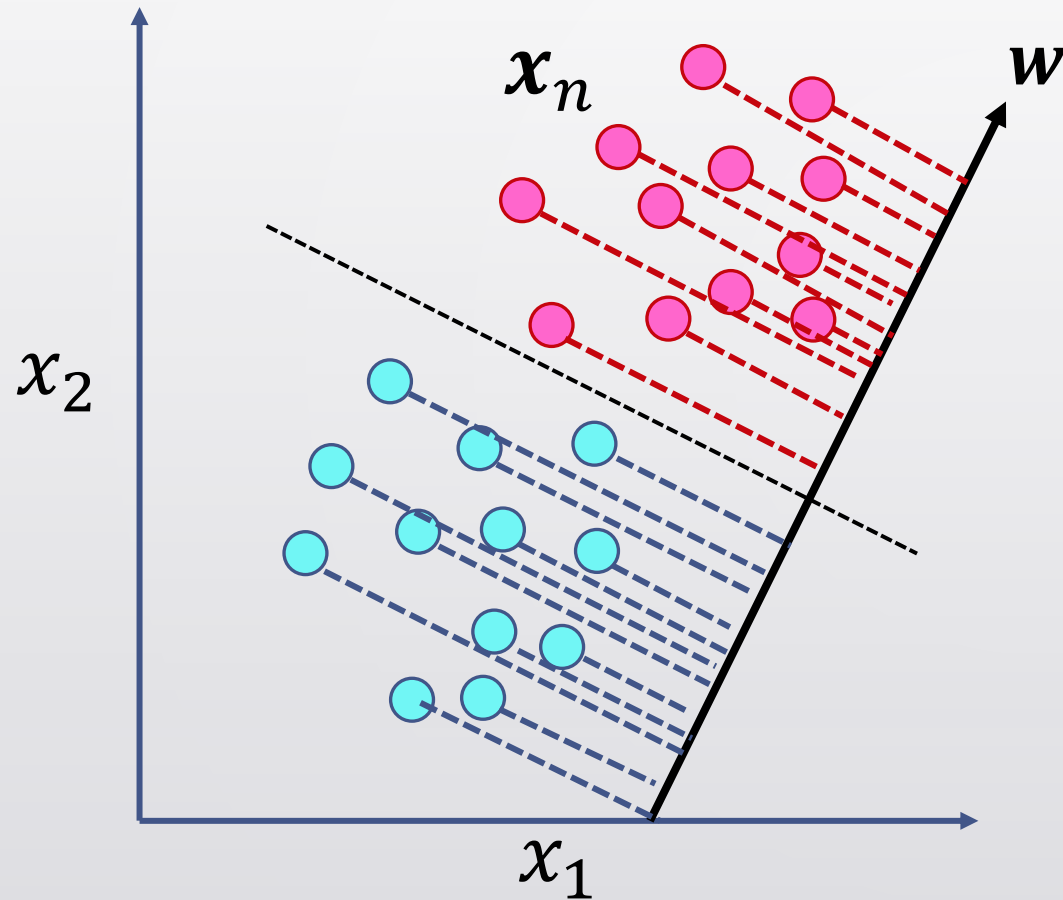
射影 $t_{n,1}$ を計算する

Compute the projection $t_{n,1}$ of the observed data x_n onto the first PC axis

$t_{n,1}$ は x_n と p_1 の内積

$t_{n,1}$ is dot(inner) product of x_n and p_1

線型判別分析 Linear Discriminant Analysis (LDA)

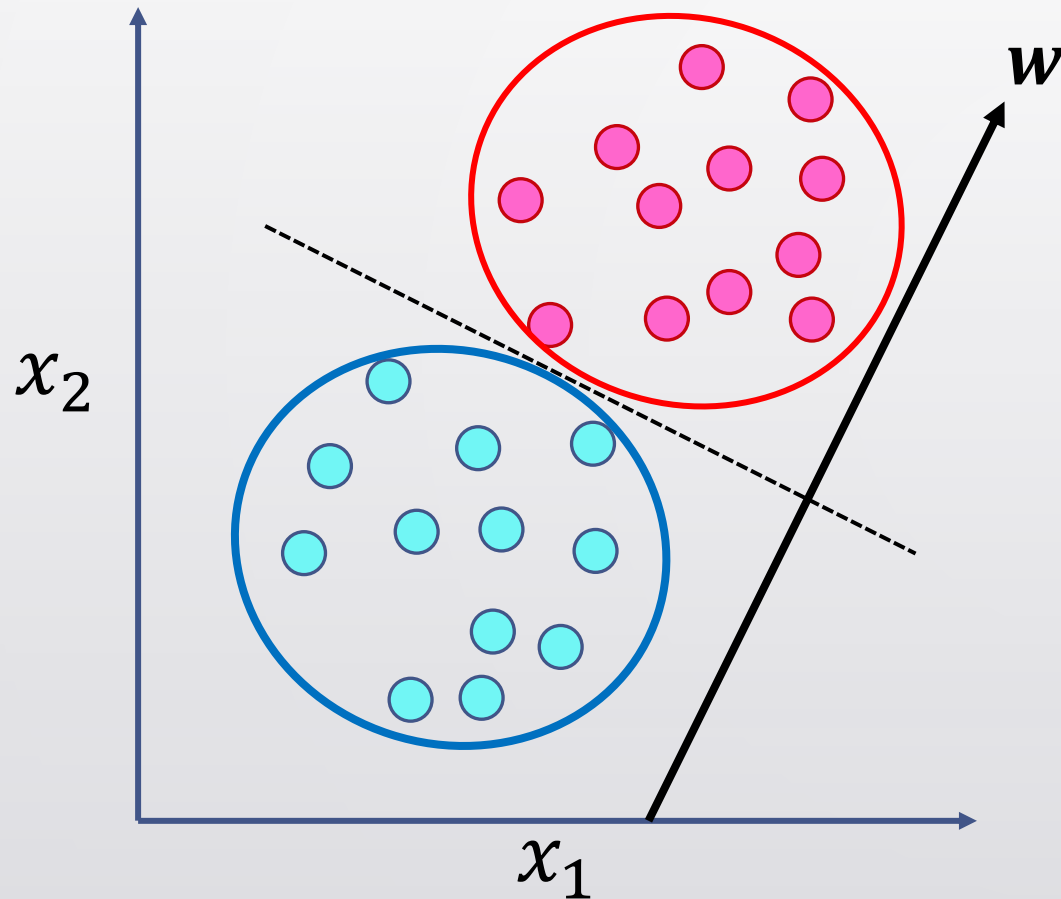


$$\mathbf{x}_n = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_M \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} \quad \|\mathbf{w}\| = 1$$

\mathbf{y}_n は \mathbf{x}_n の軸 \mathbf{w} への射影
 \mathbf{y}_n is projection of \mathbf{x}_n onto axis \mathbf{w}

$$\mathbf{y}_n = \mathbf{w}^T \mathbf{x}_n$$

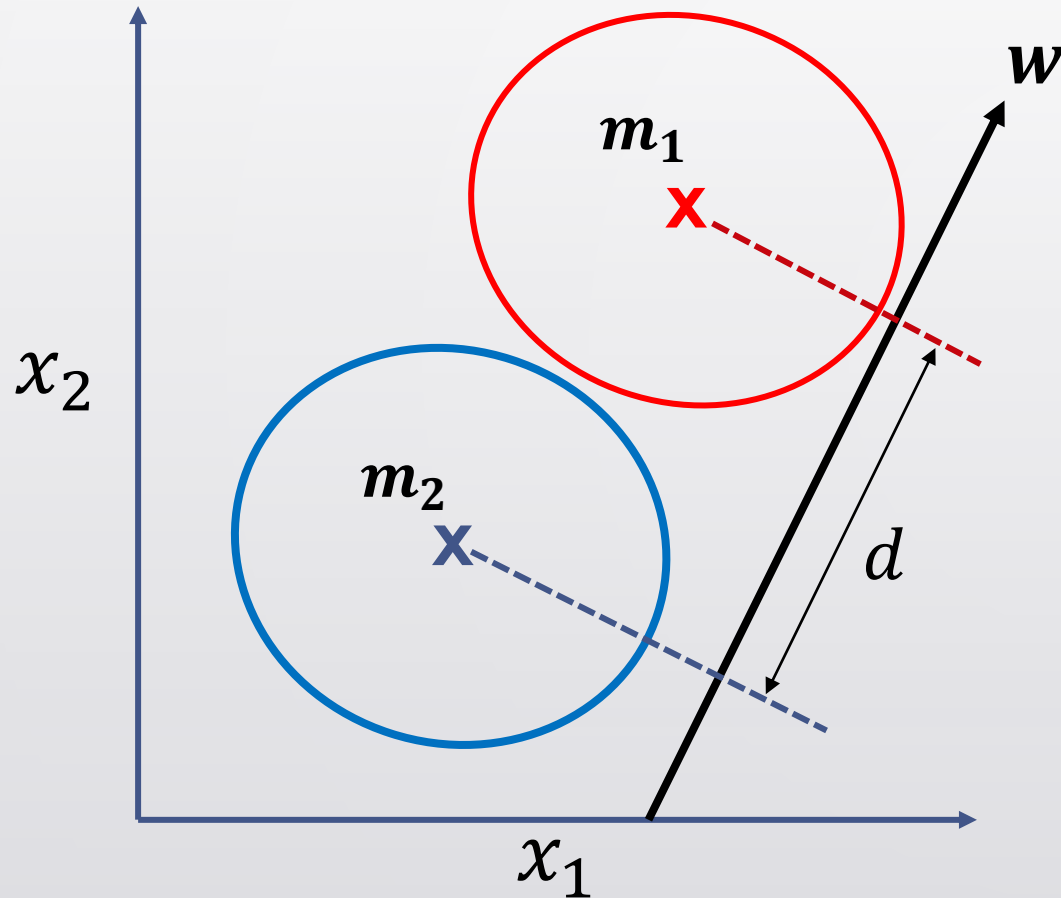
線型判別分析 Linear Discriminant Analysis (LDA)



よい決定境界は、下の二つの条件を満たす
A good decision boundary meets the two conditions below

1. 2 クラスの中心が離れている
Centers of the two classes are distant from each other
2. 各クラスのクラス内分散が小さい
Within-class variance of each class is small

線型判別分析 Linear Discriminant Analysis (LDA)



1. 2 クラスの中心が離れている

Centers of the two classes are distant from each other

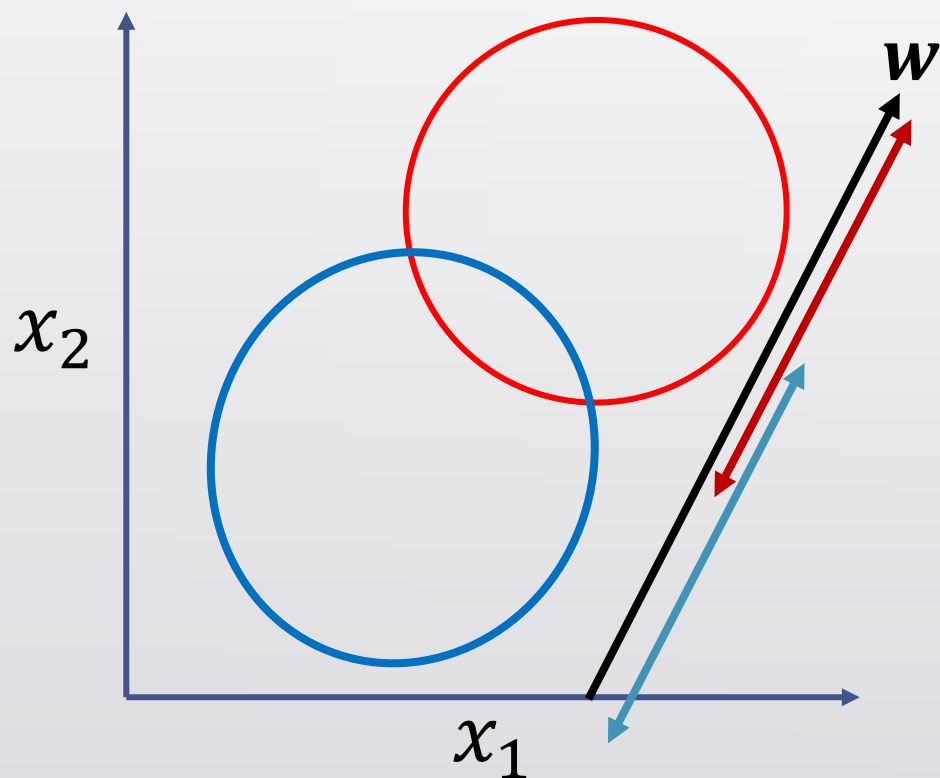
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{\mathbf{x}_k \in C_1} \mathbf{x}_k \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{\mathbf{x}_k \in C_2} \mathbf{x}_k$$

$$d = \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2 = \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)$$

線型判別分析 Linear Discriminant Analysis (LDA)

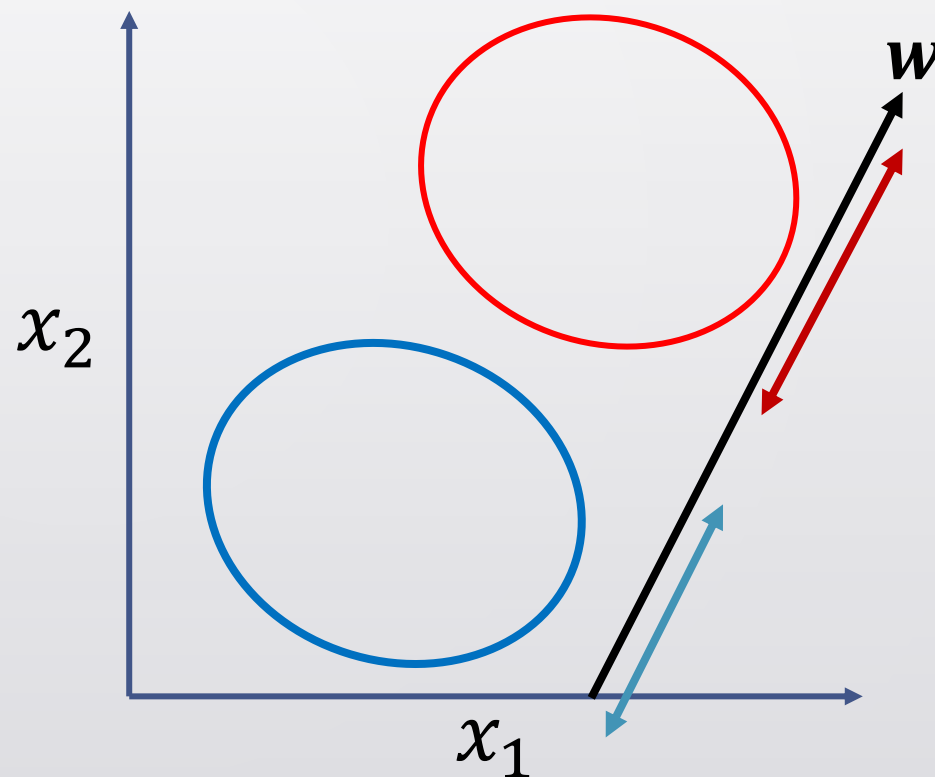
クラス内分散 大

Large within-class variance



クラス内分散 小

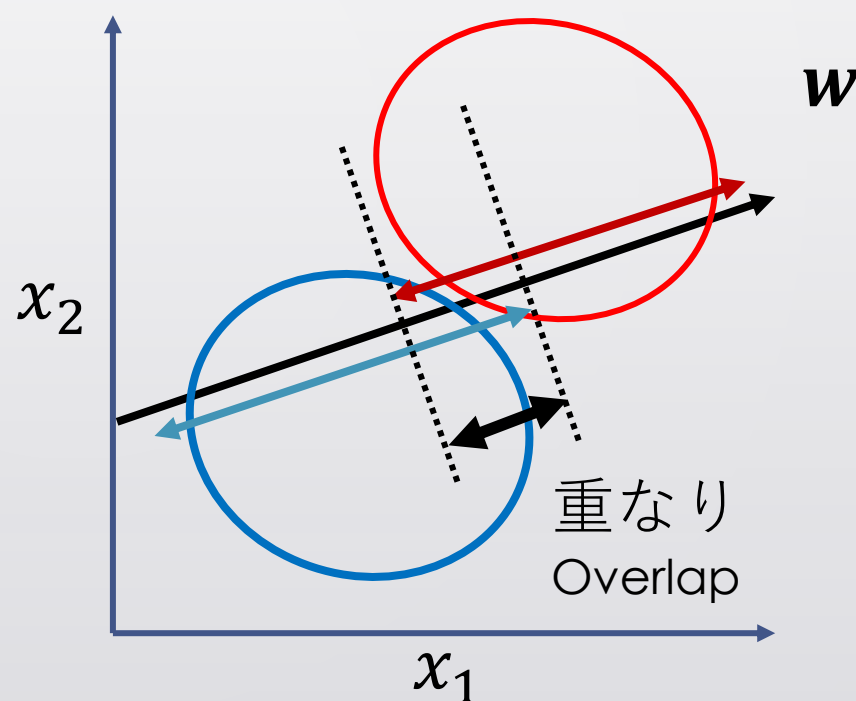
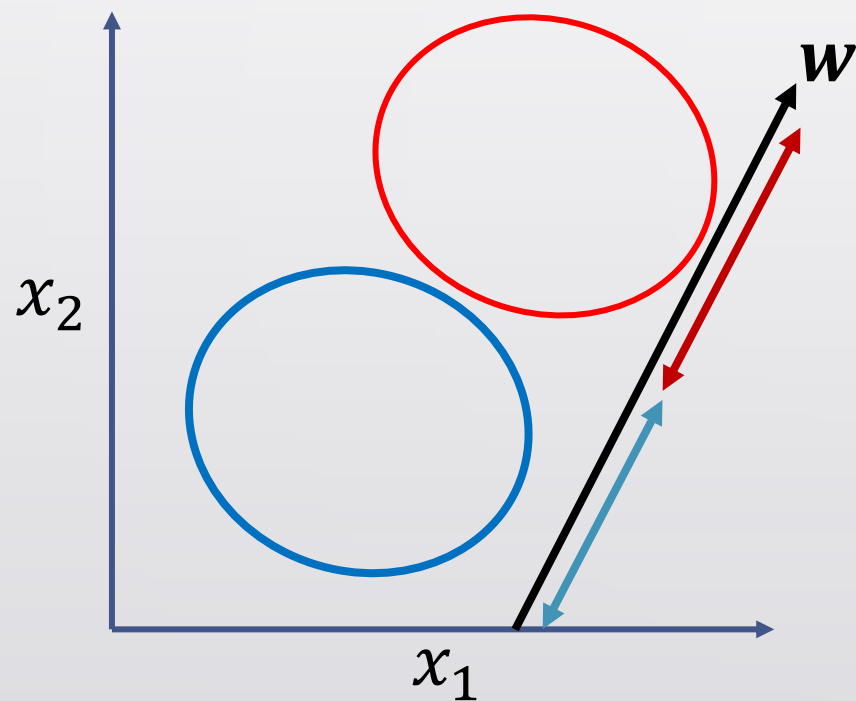
Small within-class variance



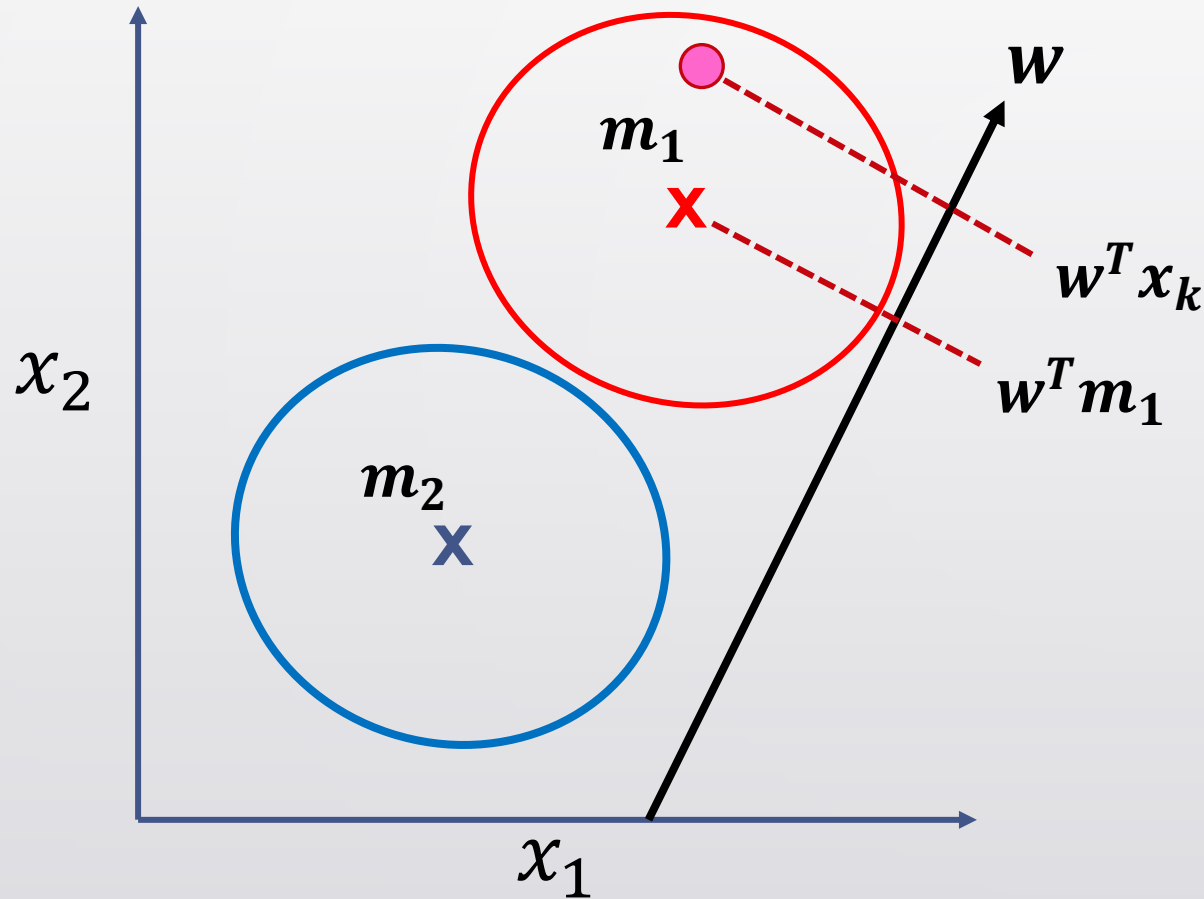
線型判別分析 Linear Discriminant Analysis (LDA)

射影のクラス内分散は軸 \mathbf{w} の向きにより変化する

Within-class variance of projection is dependent on the direction of axis \mathbf{w}



線型判別分析 Linear Discriminant Analysis (LDA)



$$s_1^2 = \frac{1}{N_1} \sum_{x_k \in C_1} (w^T x_k - w^T m_1)^2$$

$$s_2^2 = \frac{1}{N_2} \sum_{x_k \in C_2} (w^T x_k - w^T m_2)^2$$

$$s^2 = s_1^2 + s_2^2$$

線型判別分析 Linear Discriminant Analysis (LDA)

1. 2 クラスの中心が離れている Centers of the two classes are distant from each other



d を最大化する Maximize d

$$d = \mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2 = \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)$$

2. 各クラスのクラス内分散が小さい Within-class variance of each class is small



s^2 を最小化する Minimize s^2

$$s^2 = s_1^2 + s_2^2 \quad s_j^2 = \frac{1}{N_j} \sum_{\mathbf{x}_k \in C_j} (\mathbf{w}^T \mathbf{x}_k - \mathbf{w}^T \mathbf{m}_j)^2$$

線型判別分析 Linear Discriminant Analysis (LDA)

1. 2 クラスの中心が離れている Centers of the two classes are distant from each other

➡ d を最大化する ➡ d^2 を最大化する

$$d^2 = \{w^T(m_1 - m_2)\}^2 = \{w^T(m_1 - m_2)\} \{w^T(m_1 - m_2)\}^T$$

$$= w^T(m_1 - m_2)(m_1 - m_2)^T w = \mathbf{w^T S_B w}$$

$$S_B = (m_1 - m_2)(m_1 - m_2)^T$$

線型判別分析 Linear Discriminant Analysis (LDA)

2. 各クラスのクラス内分散が小さい Within-class variance of each class is small

➡ s^2 を最小化する Minimize s^2

$$s^2 = s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_w \mathbf{w}$$

$$\mathbf{S}_w = \sum_{x_k \in C_1} (x_k - \mathbf{m}_1)(x_k - \mathbf{m}_1)^T + \sum_{x_k \in C_2} (x_k - \mathbf{m}_2)(x_k - \mathbf{m}_2)^T$$

線型判別分析 Linear Discriminant Analysis (LDA)

1. 2 クラスの中心が離れている Centers of the two classes are distant from each other
2. 各クラスのクラス内分散が小さい Within-class variance of each class is small

➡ $J(w)$ を最大化する Maximize $J(w)$

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

線型判別分析 Linear Discriminant Analysis (LDA)

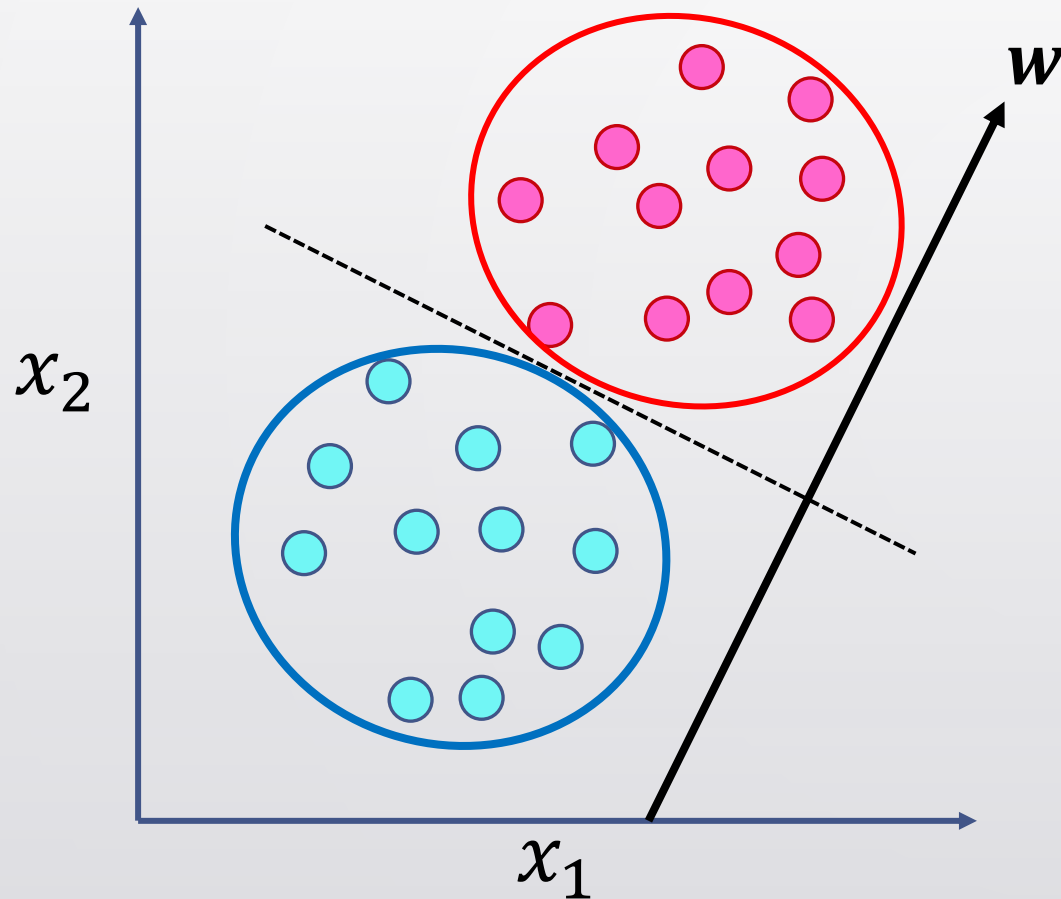
$$J(w) \text{を最大化する Maximize } J(w) \quad J(w) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

$J(w)$ を最大化する \mathbf{w} は下の固有方程式を満たす
 \mathbf{w} that maximizes $J(w)$ satisfies the eigen equation below

$$\mathbf{S}_w^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w} \qquad \mathbf{S}_B \mathbf{w} = \mathbf{S}_w \mathbf{S}_w^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

$$J(w) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} = \frac{\lambda \mathbf{w}^T \mathbf{S}_w \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} = \lambda$$

線型判別分析 Linear Discriminant Analysis (LDA)

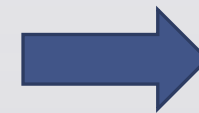


$$S_w^{-1} S_B w = \lambda w$$

$$J(w) = \frac{w^T S_B w}{w^T S_w w} = \frac{\lambda w^T S_w w}{w^T S_w w} = \lambda$$

軸 w は最大の固有値に対応する固有ベクトルと並行

Axis w is in parallel with eigen vector corresponding to the largest eigen value

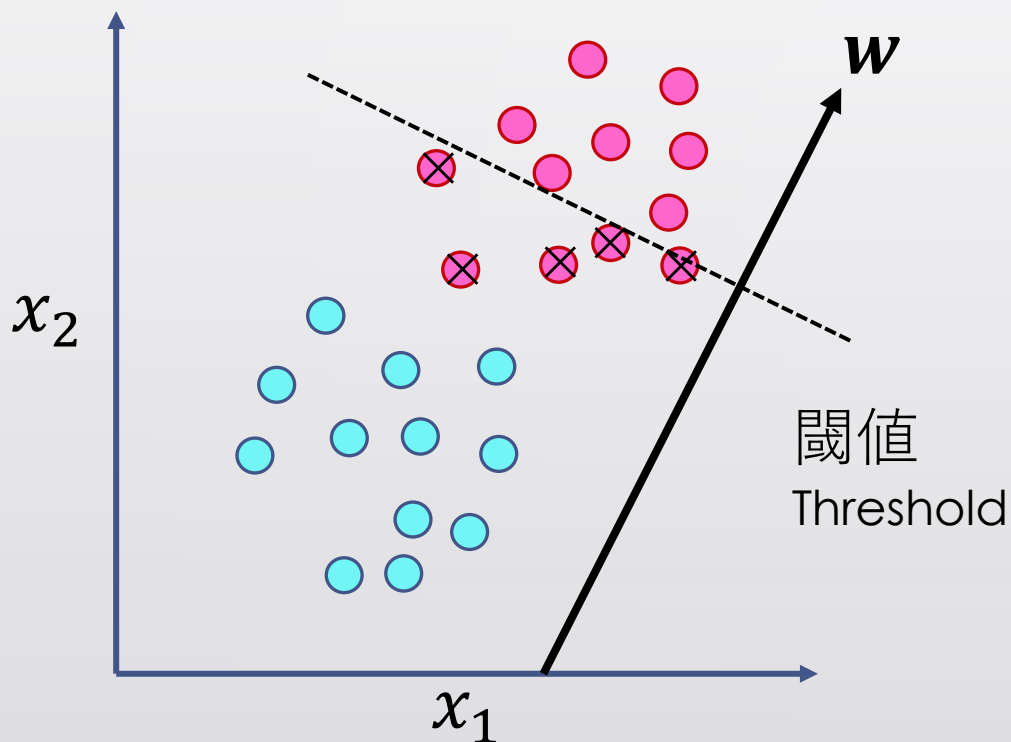


決定境界の向きが決まる

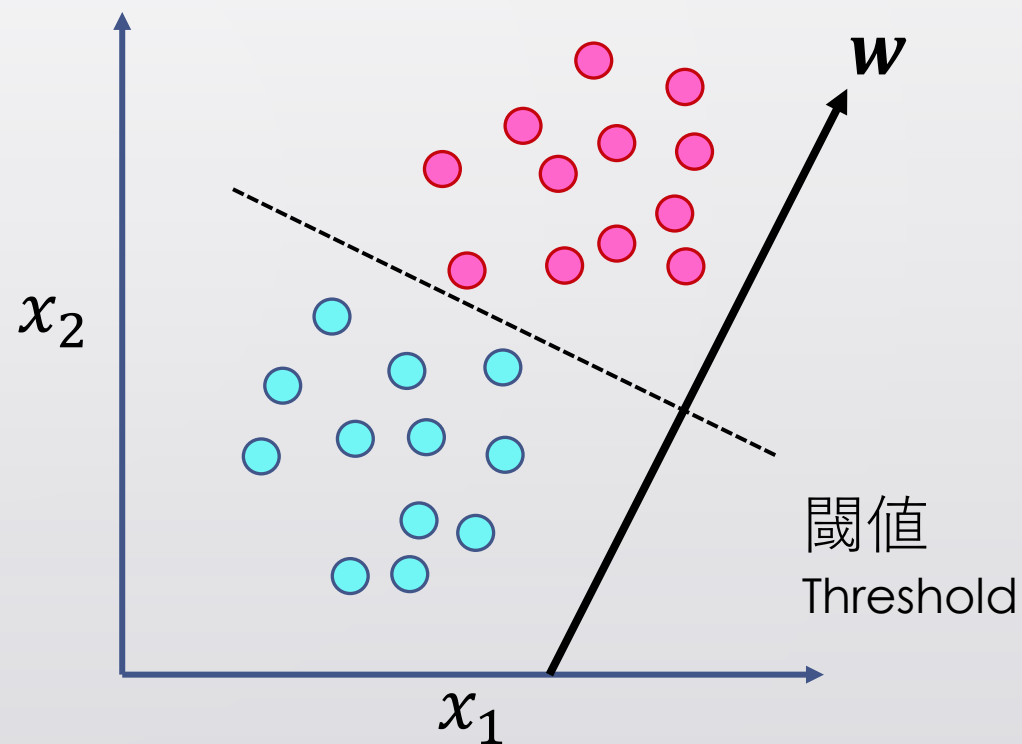
Orientation of decision boundary is determined

閾値をどう決定するか？ How should we determine the threshold?

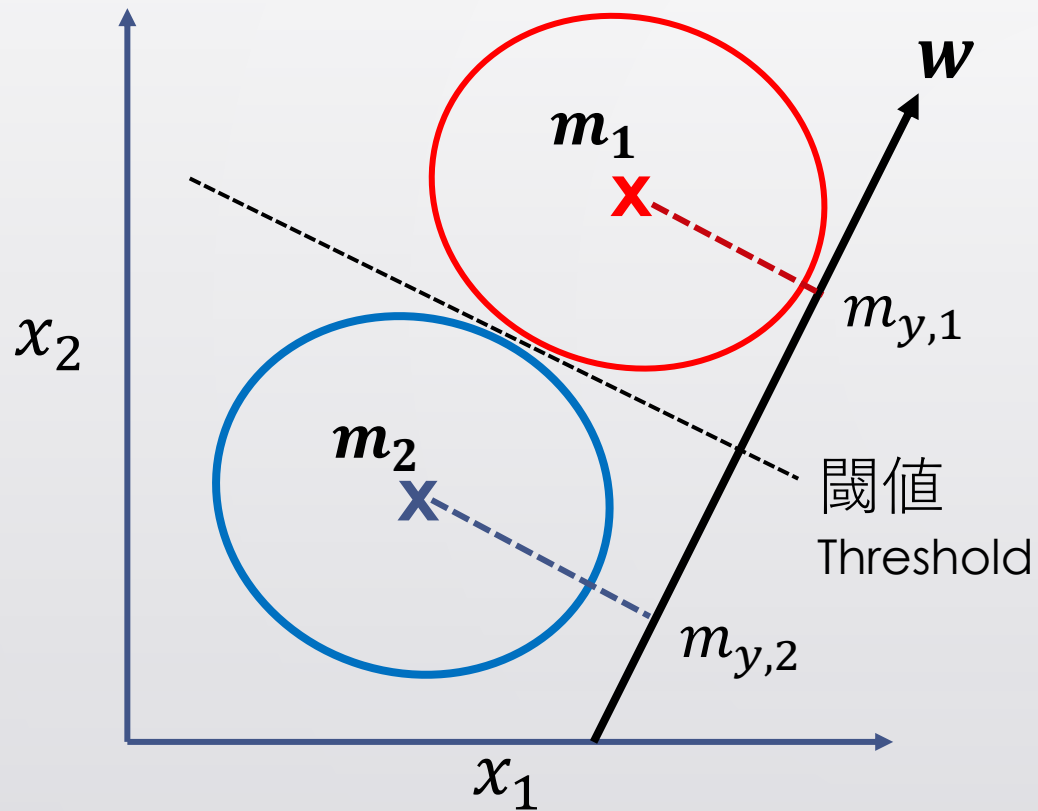
不適切な閾値 Inappropriate threshold



適切な閾値 Appropriate threshold



//////////
閾値をどう決定するか？ How should we determine the threshold?



各クラスの中心の射影の加重平均を
閾値にする

Adopt as threshold value the weighted mean of
projection of center of each class

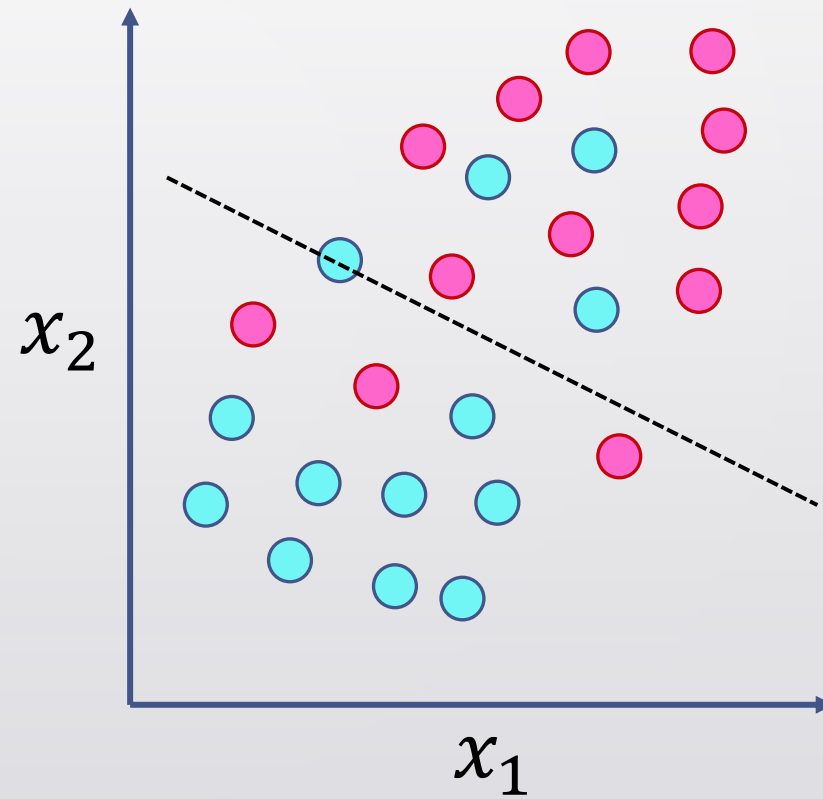
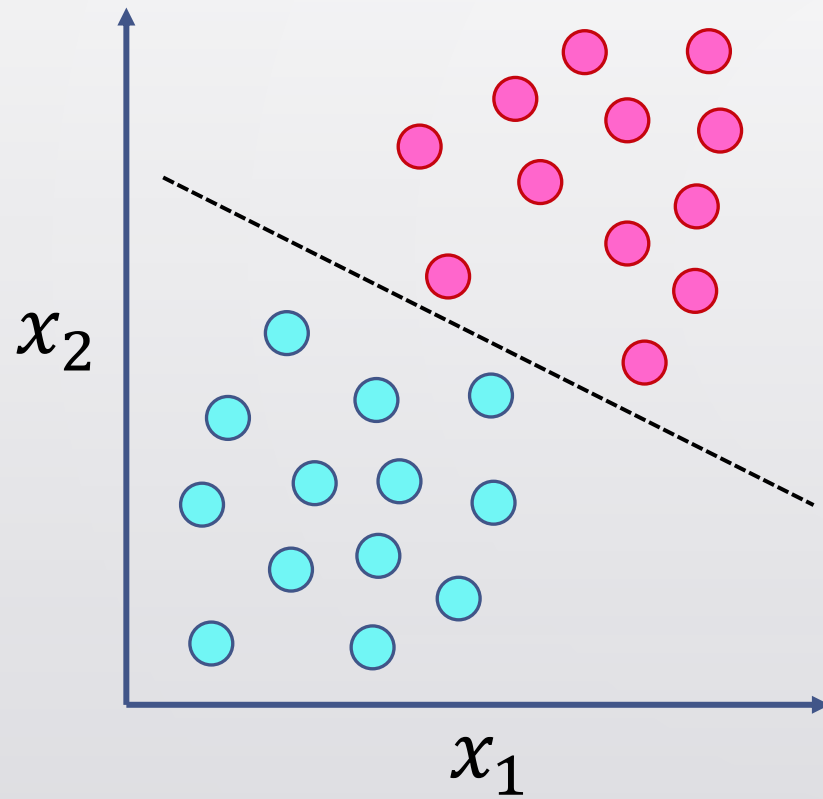
$$Threshold = \frac{N_1 s_{y,1}^2 m_{y,1} + N_2 s_{y,2}^2 m_{y,2}}{N_1 s_{y,1}^2 + N_2 s_{y,2}^2}$$

$s_{y,j}^2$: クラスjのデータの射影の分散

$m_{y,j}$: クラスjの重心の射影

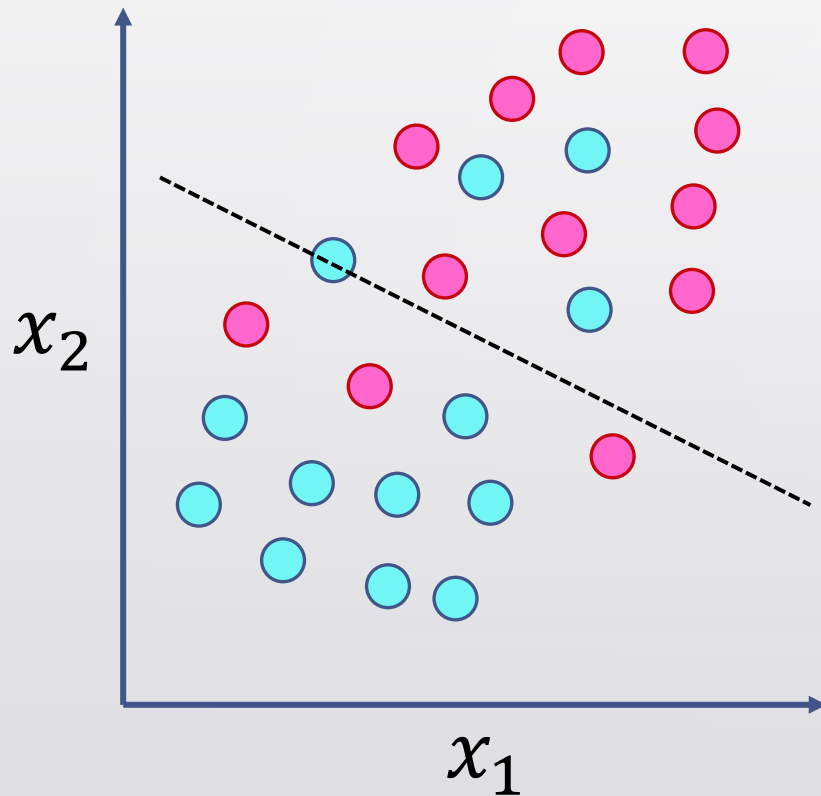
//////
“線型分離可能”とは？

What does “Linearly Separable” Mean?



//////
“線型分離可能”とは？

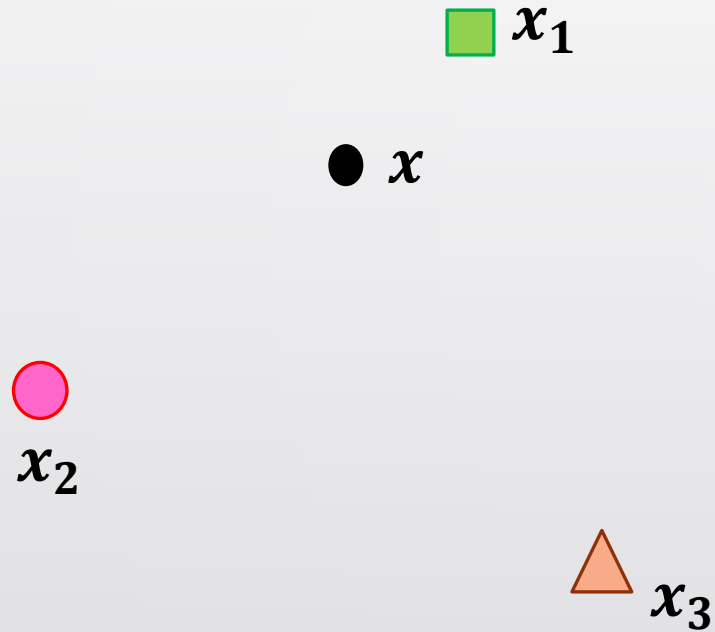
What does “Linearly Separable” Mean?



線型分離可能でない問題には、LDA
がうまく機能しない

LDA does not work well for linearly
inseparable problems

最近傍法 Nearest Neighbor Method



データ x は最も近くにある鋳型データ x_j と同じクラスに属するとみなす

Data x is judged to belong to the same class as template data x_j

最近傍法 Nearest Neighbor Method

クラス C_i の鋳型と \mathbf{x} の最小距離

Shortest distance between \mathbf{x} and templates of class C_i

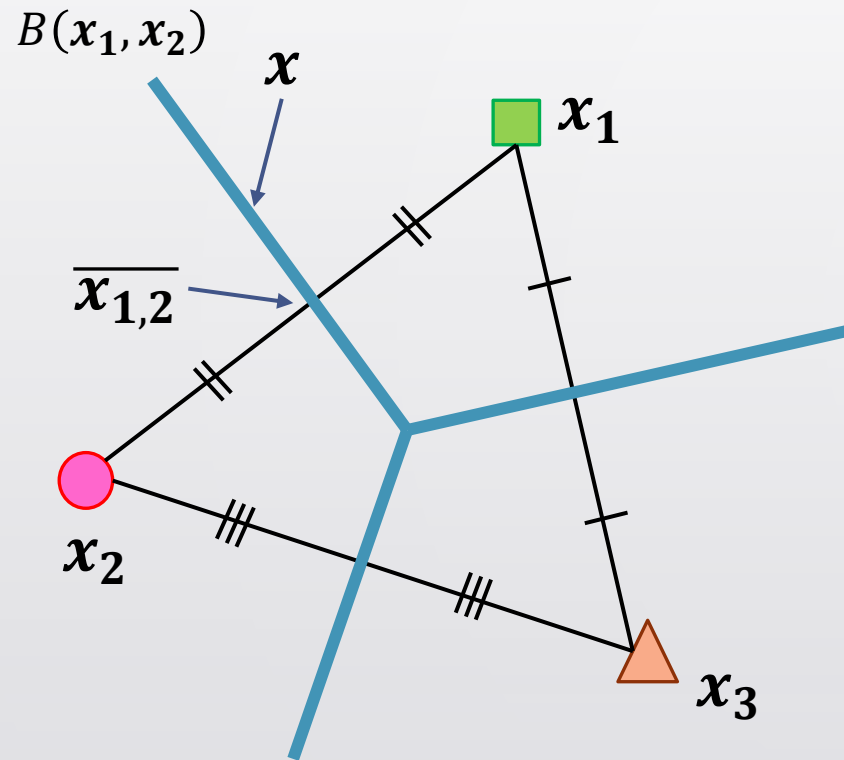
$$\underline{\operatorname{argmin}_i \{ \min_j d(\mathbf{x}, \mathbf{x}_j^i) \}} \quad \text{if } \min_{i,j} d(\mathbf{x}, \mathbf{x}_j^i) < t$$

\mathbf{x} と最も近い鋳型が属するクラスの識別番号を返す

Returns the identifier of the class to which the template closest to \mathbf{x}

Reject if $\min_{i,j} d(\mathbf{x}, \mathbf{x}_j^i) \geq t$

ボロノイ境界 Voronoi Boundary



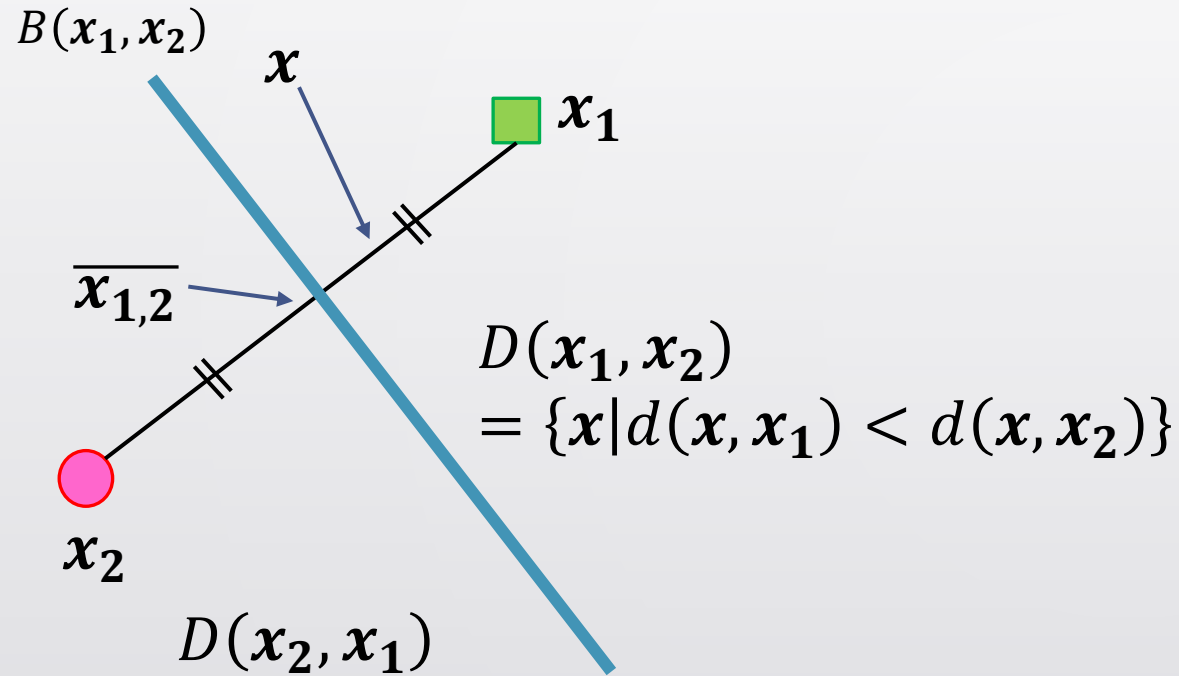
一対の鋳型から等距離にある点の集合

Set of points equidistant from a pair of templates

$$B(x_i, x_j) = \{x \mid d(x, x_i) = d(x, x_j)\}$$

$$(x - \overline{x_{i,j}}) \cdot (x_i - x_j) = 0$$

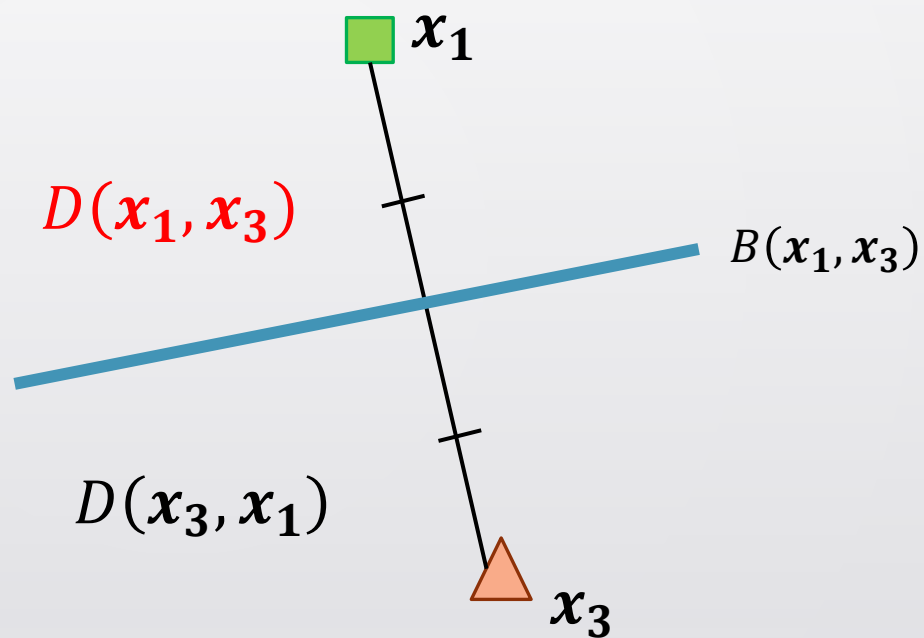
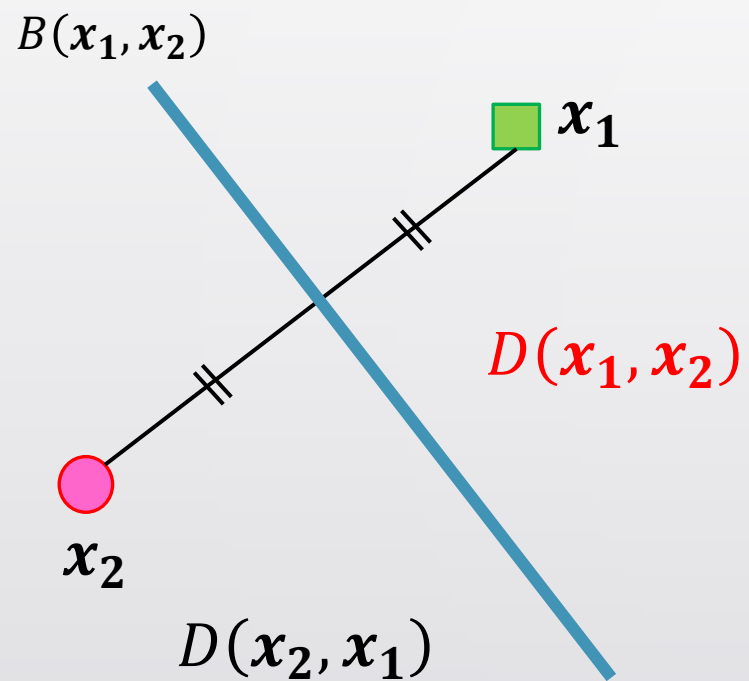
ボロノイ領域 Voronoi Region



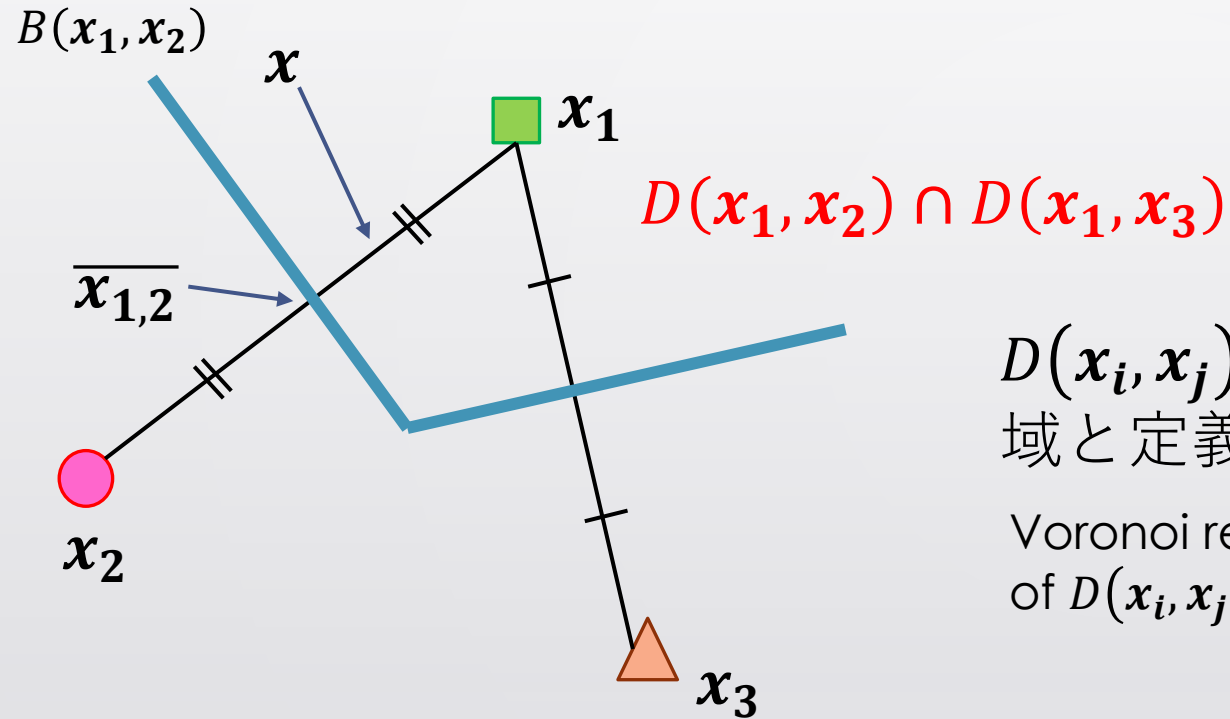
$D(x_i, x_j)$:

点 x_j より点 x_i に距離が近い点の集合
Set of points closer to x_i than x_j

ボロノイ領域 Voronoi Region



ボロノイ領域 Voronoi Region



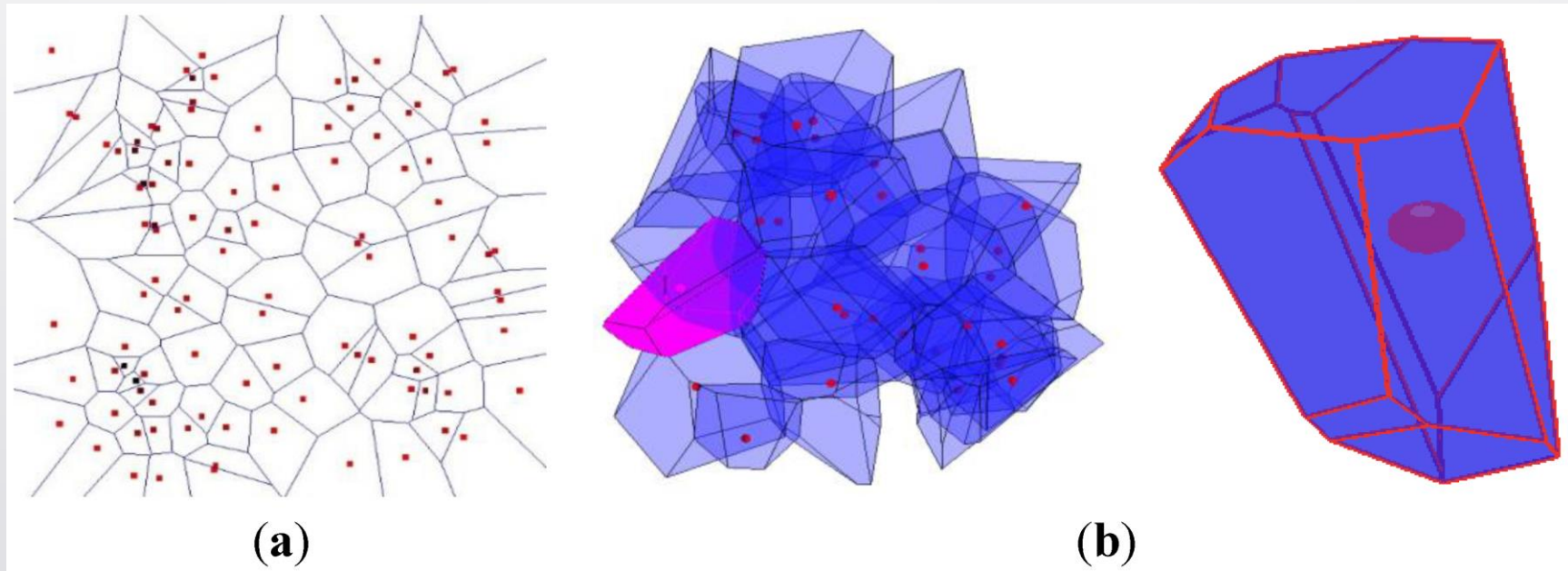
$D(x_i, x_j) (j \neq i)$ の積集合を x_i のボロノイ領域と定義する

Voronoi region of x_i is defined as set intersection of $D(x_i, x_j) (j \neq i)$

ボロノイ図 Voronoi Diagram

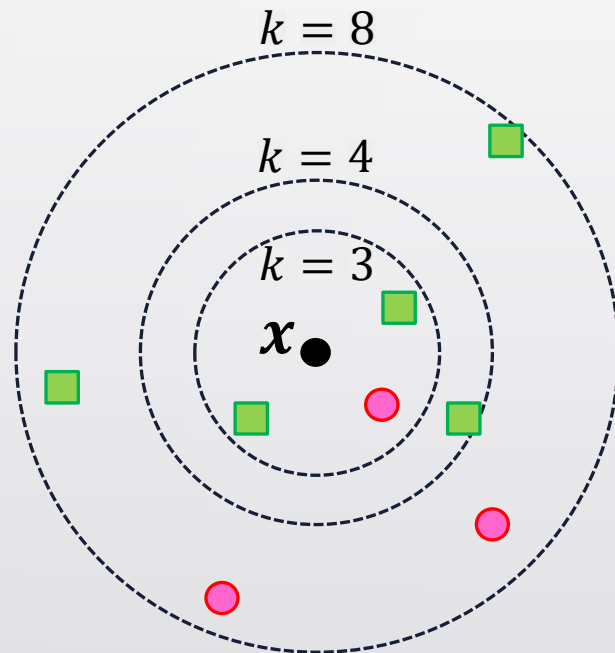
最近傍法の決定境界はボロノイ図を描く

Voronoi diagram shows configuration of decision boundaries in nearest neighbor method



<https://www.mdpi.com/2220-9964/4/3/1480>

k 最近傍法 k Nearest Neighbor Method

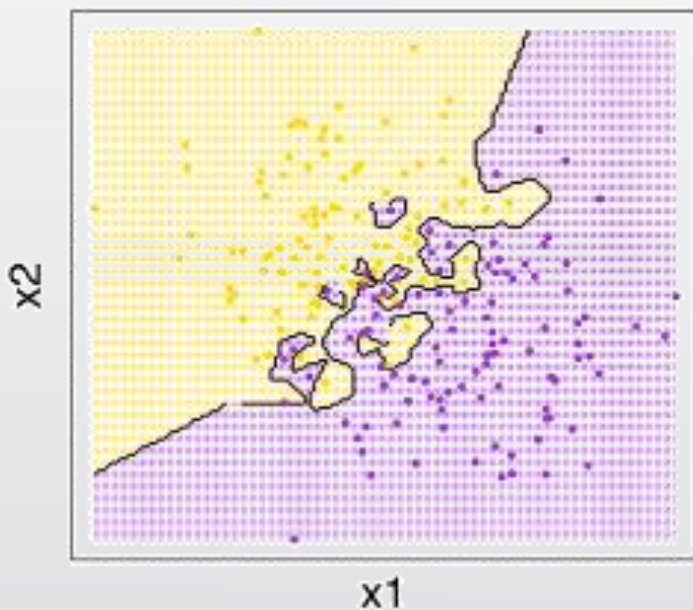


データ x のクラスを最近傍にある k 個のデータの多数決投票により決定する

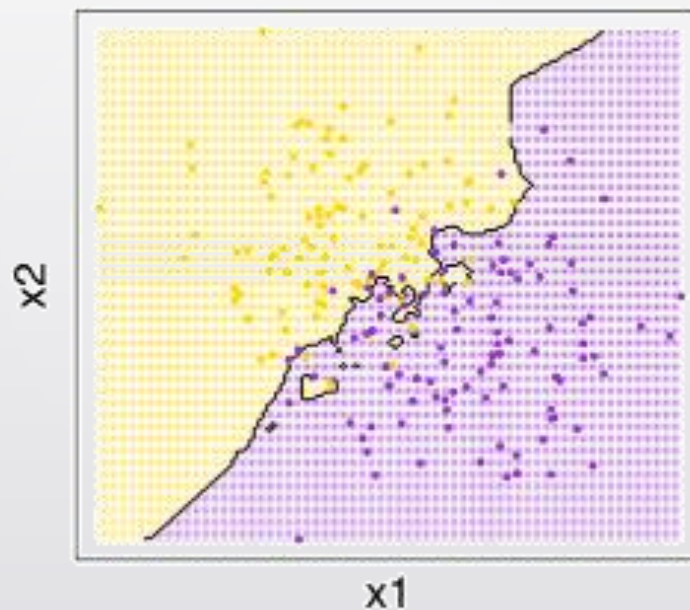
Class of data x is determined by majority voting of k data points closest to x

k 最近傍法 k Nearest Neighbor Method

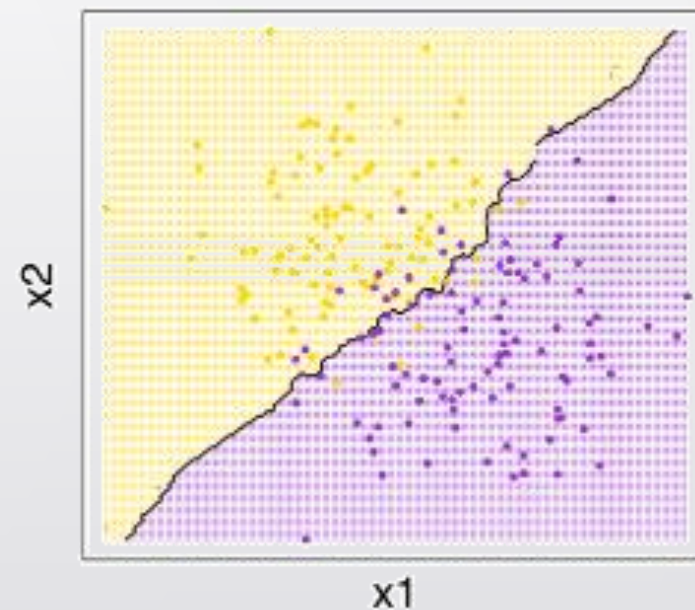
Binary kNN Classification (k=1)



Binary kNN Classification (k=5)



Binary kNN Classification (k=25)



<https://elvyna.github.io/2019/knn-explained/>