

実践法としてのデータ科学



CDS

早稲田大学
データ科学センター

- データ分析の課題に対して、どのように意思決定写像を定めるのかを例で説明

例 構造推定としての回帰問題

目的: 説明変数 x_1, \dots, x_n が与えられたもとでの y_1, \dots, y_n の確率的データ生成観測メカニズムを明らかにしたい

設定: $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$
 ε_i は i.i.d. で正規分布 $\mathcal{N}(0, \sigma_\varepsilon^2)$ に従う

評価基準: 尤度最大化

2変数データ
(量的データと量的データ)
 $(x_1, y_1), \dots, (x_n, y_n)$

意思決定写像

母回帰係数と
誤差項の分散
の推定量
 $\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}_\varepsilon^2$

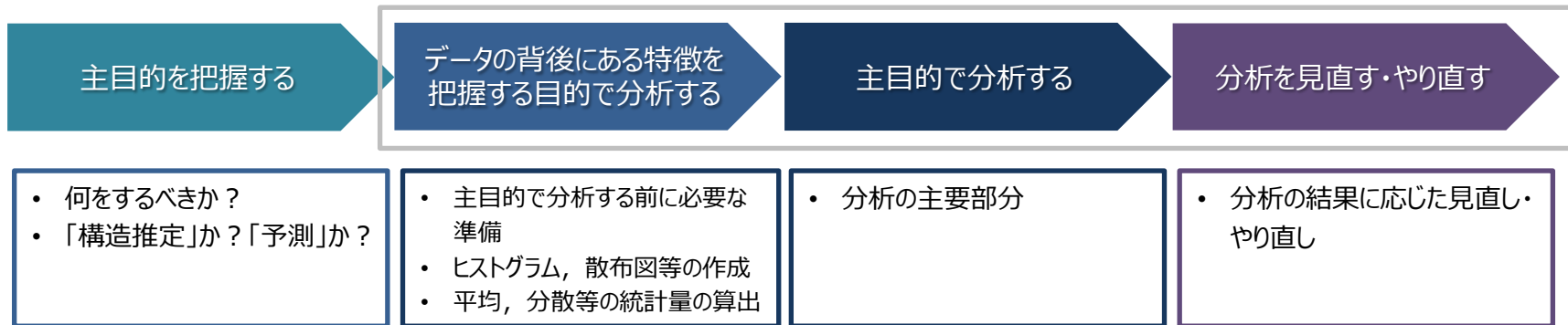
講義科目「データ科学実践」のご紹介

対象者：「データ科学入門 α ・ β ・ γ ・ δ 」を学んだ学生

講義内容：いくつかのデータ分析の課題について、分析の筋道を立て、実際に分析を行って結果を考察するまでの一連の流れを演習する

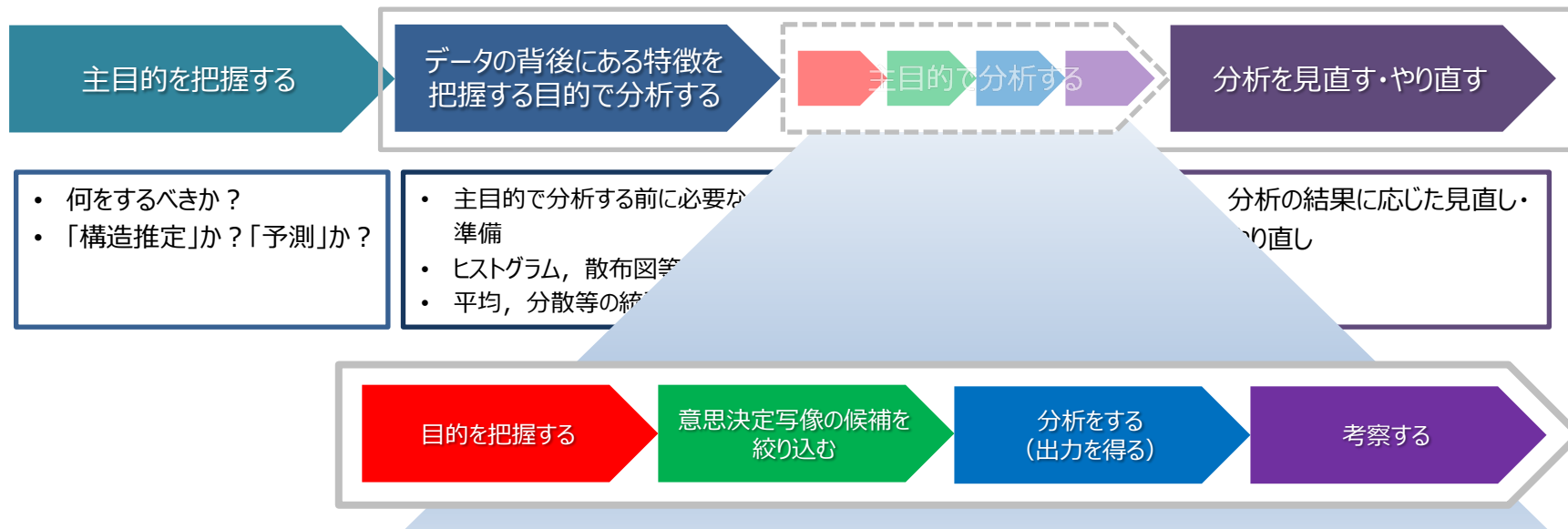
講義形式：フルオンデマンド

主目的での分析：「構造推定」または「予測」



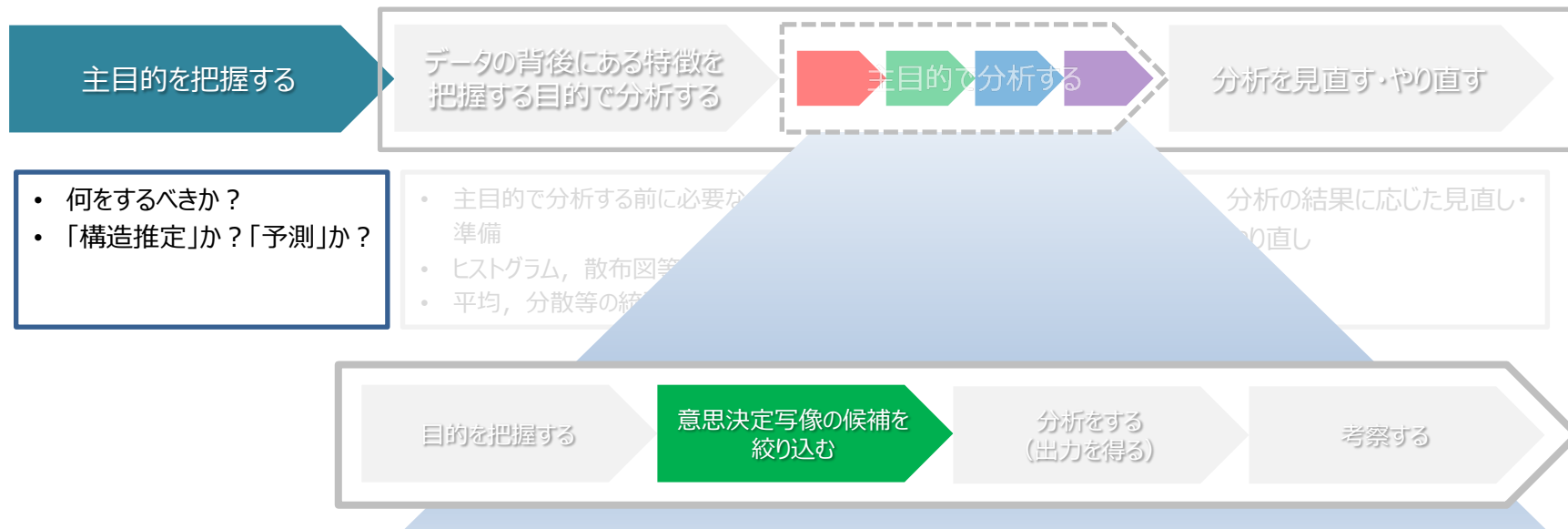
データ分析の全体の流れ

主目的での分析：「構造推定」または「予測」



データ分析の全体の流れ

主目的での分析：「構造推定」または「予測」



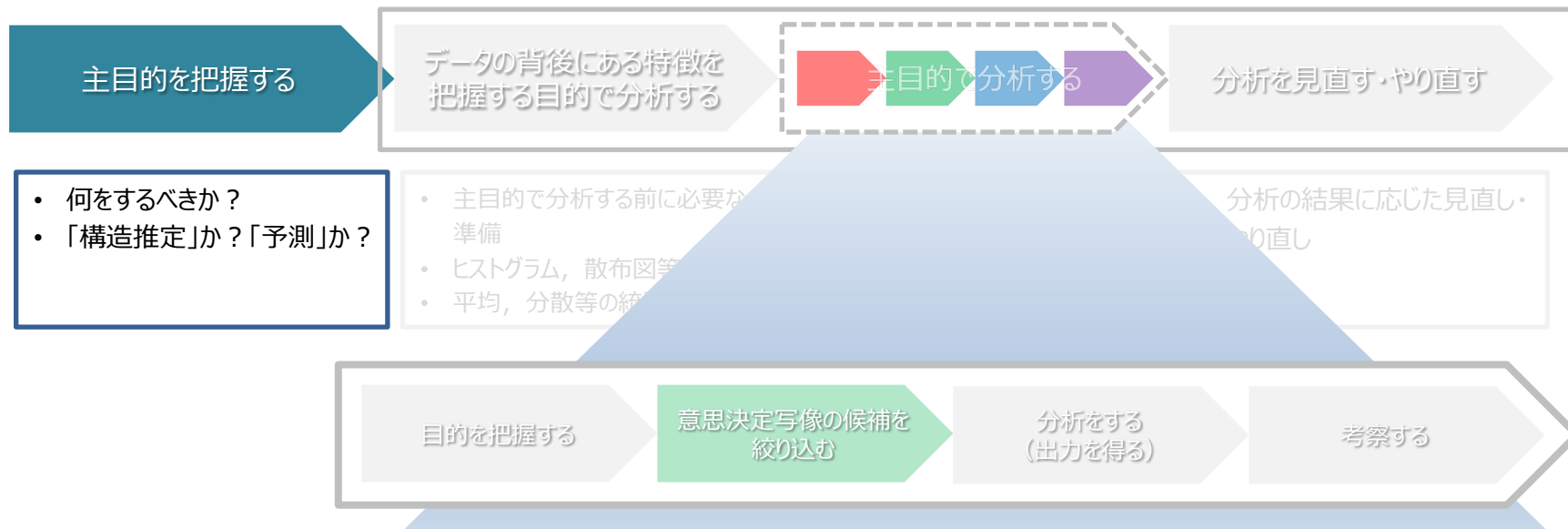
ある不動産会社では、

- 過去に販売した中古住宅に関するデータが得られている
- 物件情報から価格がどのように定まるのかを知りたいと考えている

No.	面積 (m ²)	築年数 (年)	主要駅からの 電車時間 (分)	最寄り駅までの 徒歩時間 (分)	建物種別	駐車施設	価格 (百万円)
1	125	26	25	21	戸建て	有	26.9
2	85	18	19	11	戸建て	有	36.8
3	100	25	30	9	戸建て	有	7.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

データ分析の全体の流れ

主目的での分析：「構造推定」または「予測」



主目的の把握

目的：構造推定 [物件情報から価格が定まる構造を表す式を推定する]

説明変数：面積，築年数，主要駅からの電車時間，徒歩時間，
建物種別，駐車施設

目的変数：価格（量的変数，観測可能）

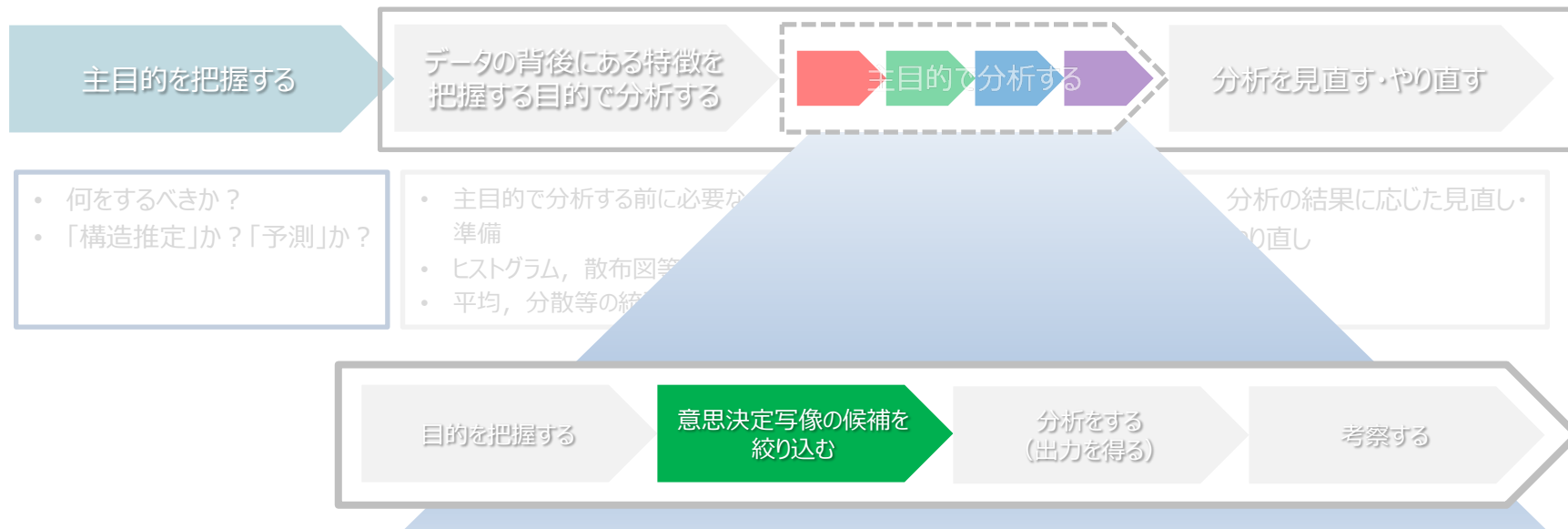
問題の種類：回帰



目的変数	量的	質的
観測可能	回帰	分類
観測不可能	縮約	クラスタリング

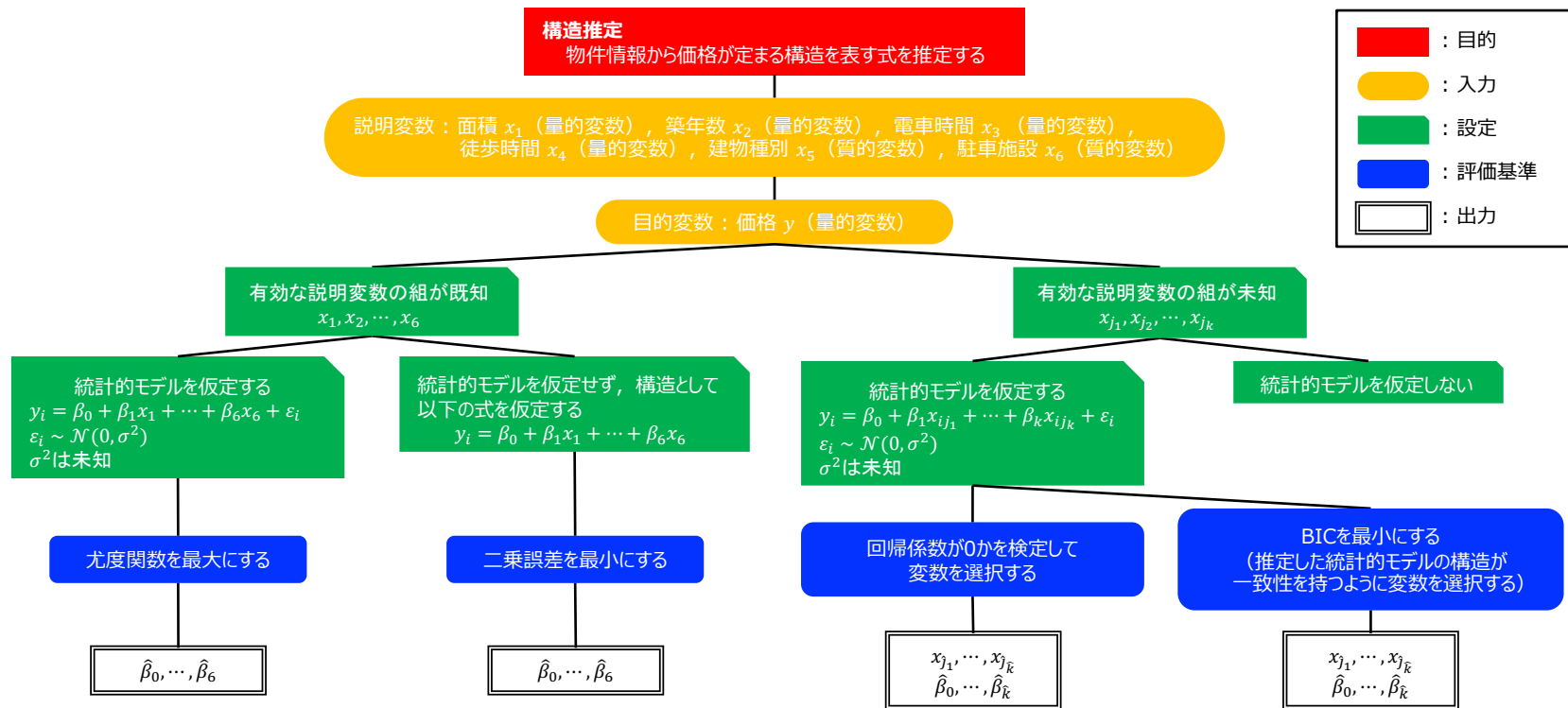
データ分析の全体の流れ

主目的での分析：「構造推定」または「予測」



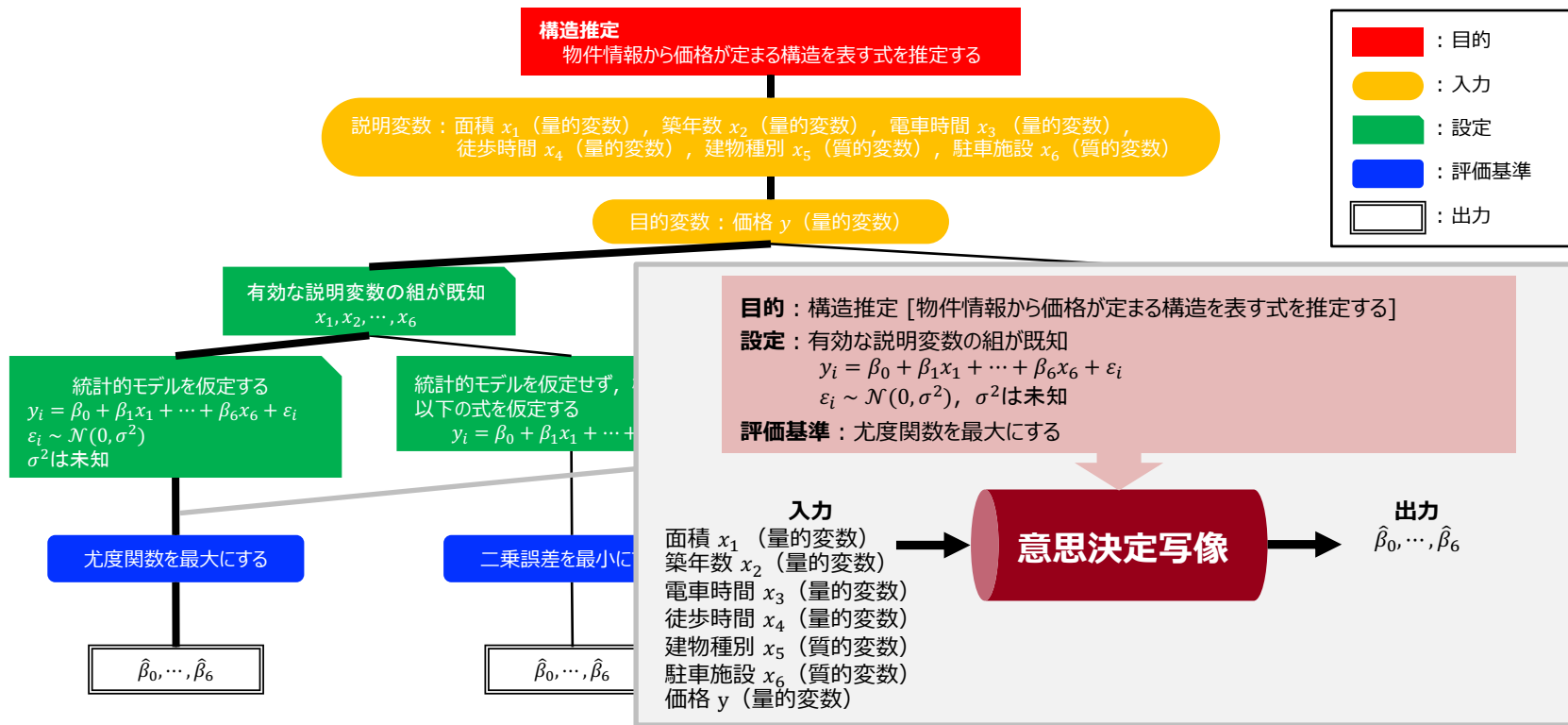
意思決定写像の候補の絞り込み

意思決定写像構成木



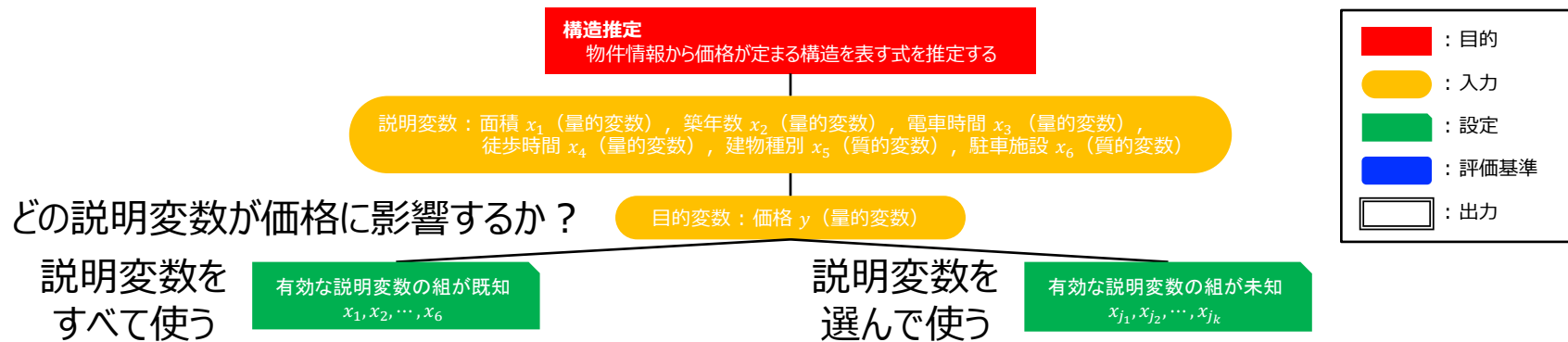
意思決定写像の候補の絞り込み

意思決定写像構成木



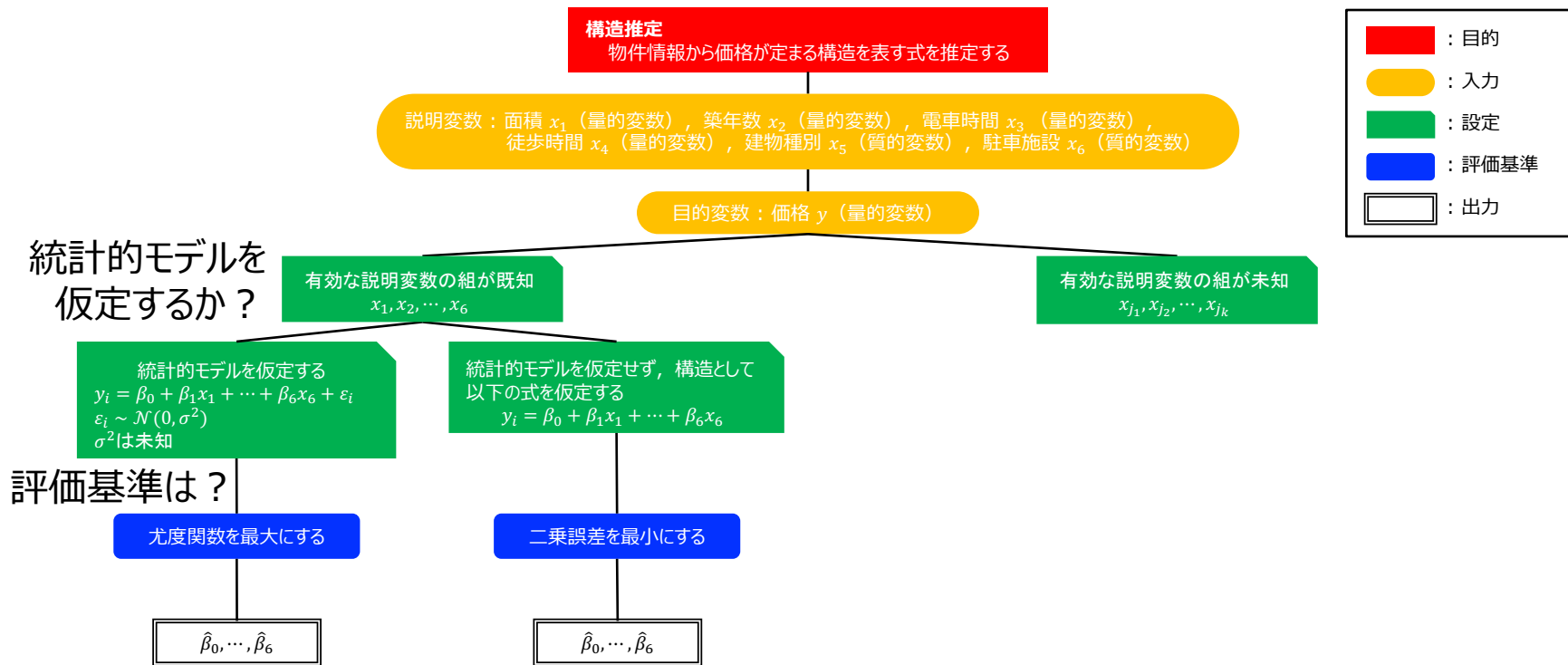
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



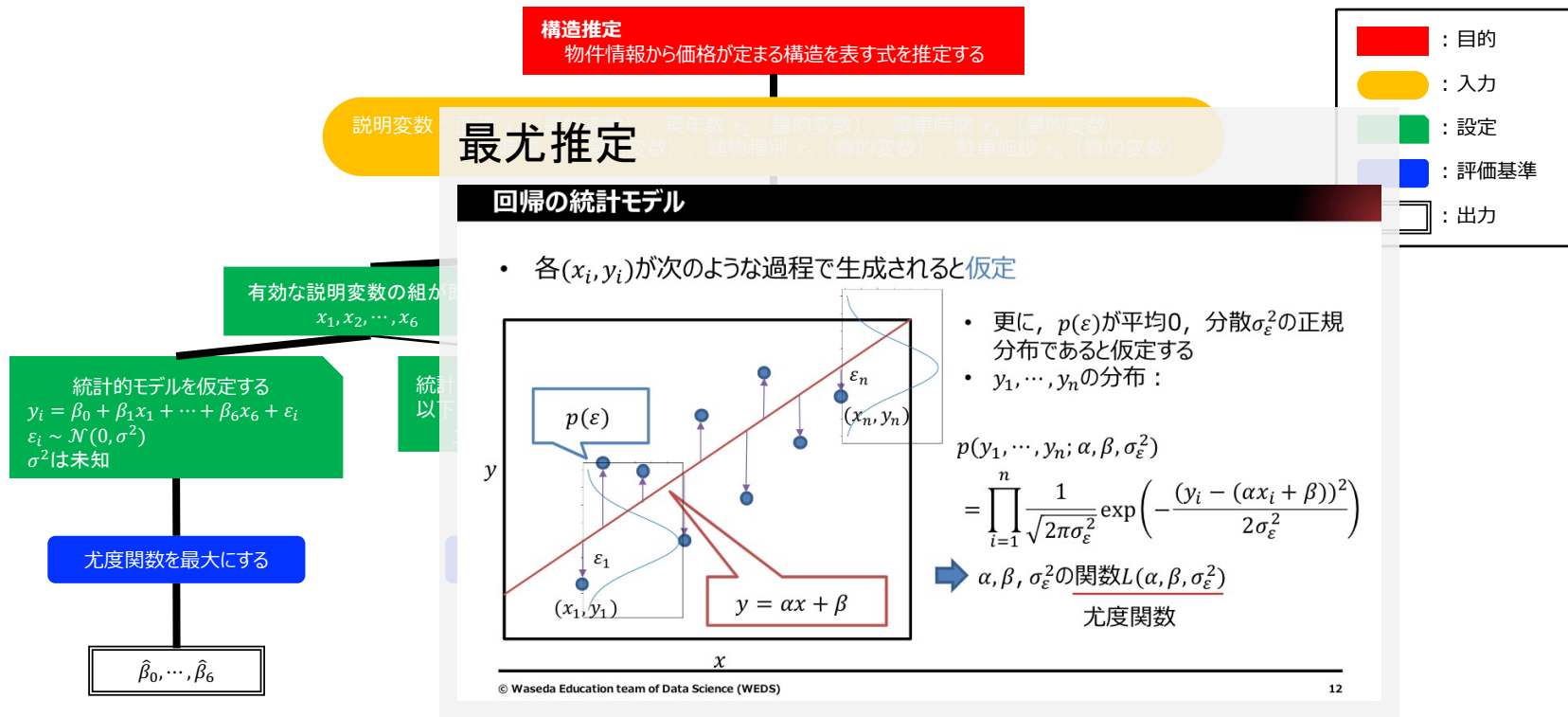
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



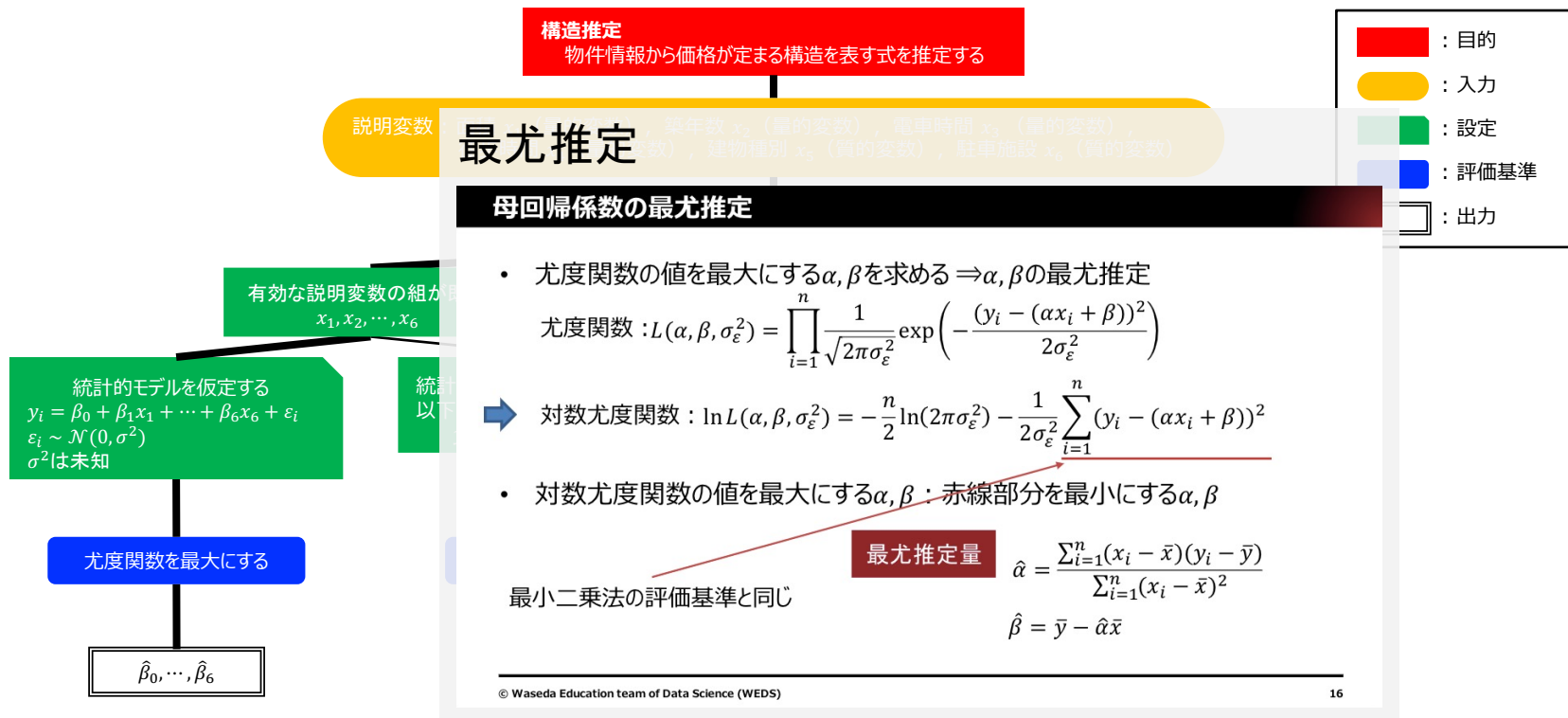
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



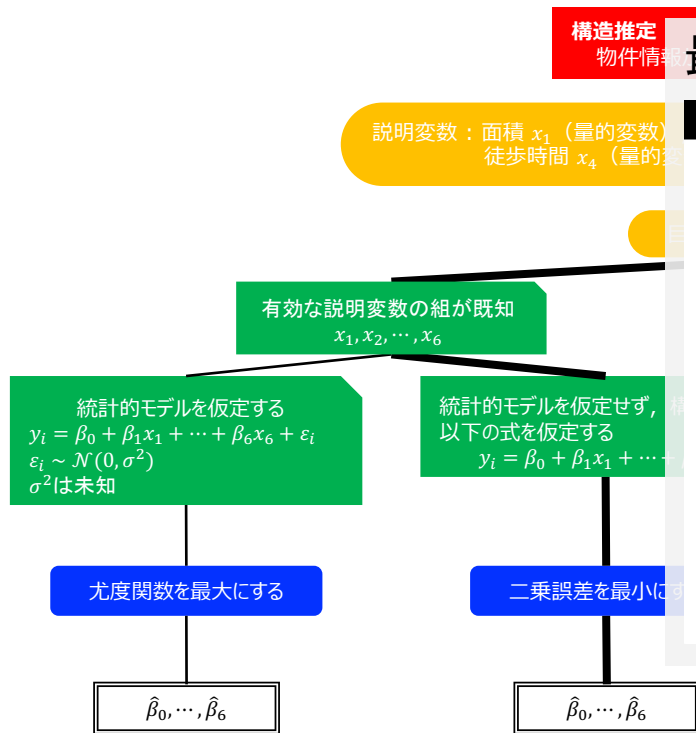
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる

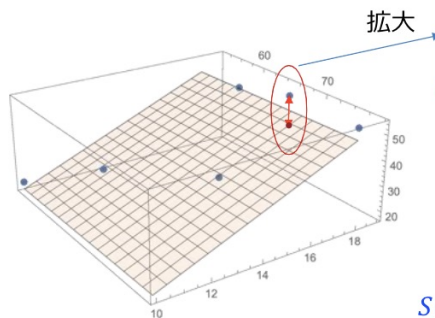


最小二乗法

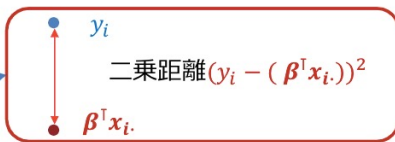
重回帰分析

- 最小二乗法による $\hat{\beta}$ の決定

回帰式上の点 $\beta^T x_i$ とデータ上で x_i に対応した点 y_i の間の「距離」として
二乗距離: $(y_i - (\beta^T x_i))^2$



拡大



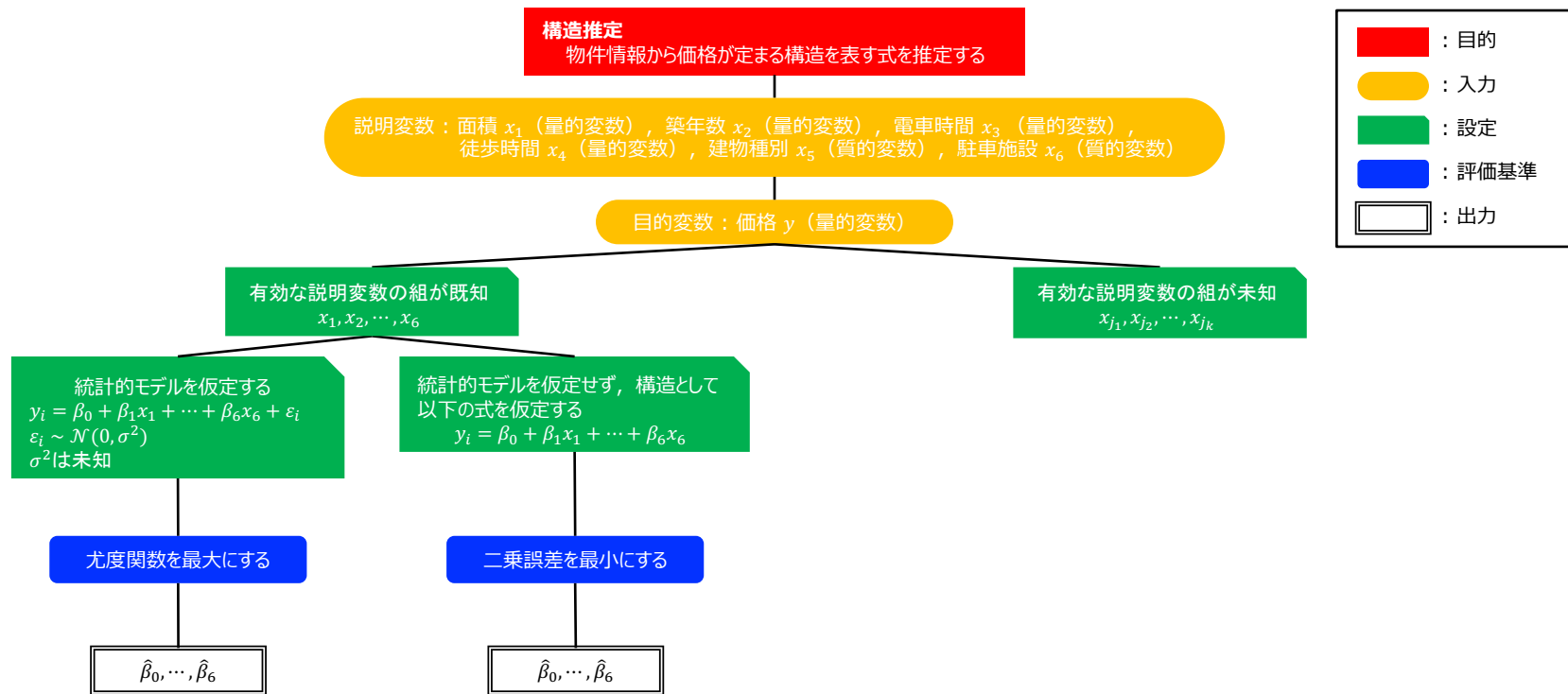
全ての点に対する二乗距離の合計値
(二乗誤差損失):

$$S(\beta) = \sum_{i=1}^n (y_i - \beta^T x_i)^2$$

$S(\beta)$ を最小にする β を決定 = 最小二乗法

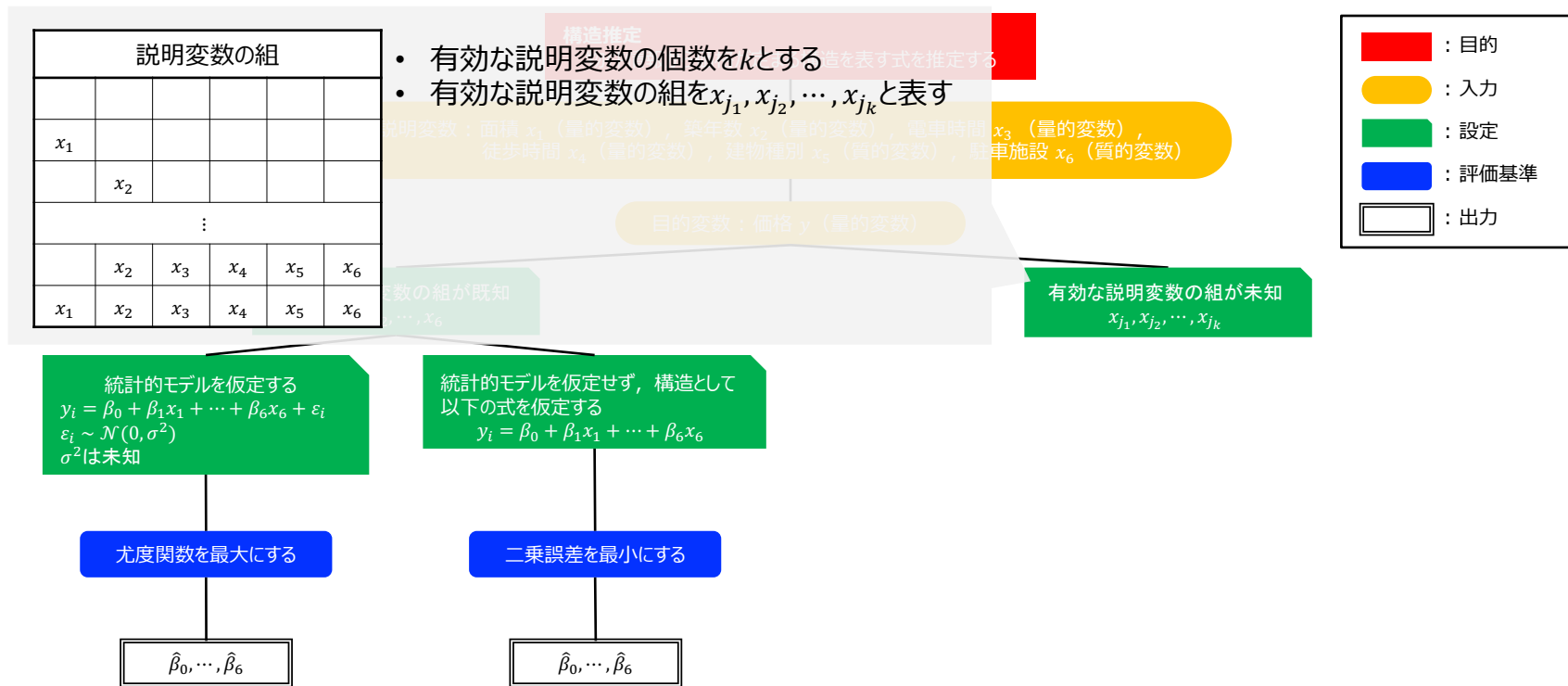
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



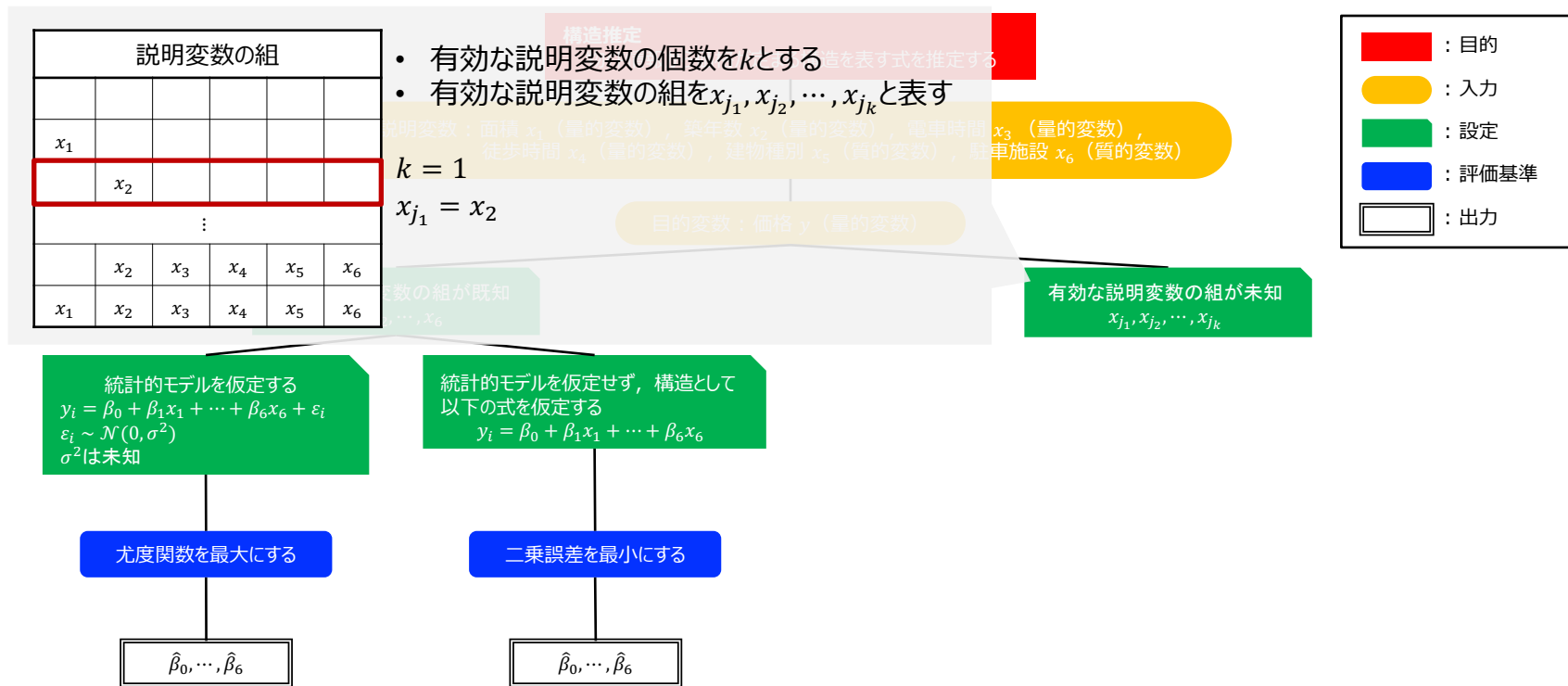
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる



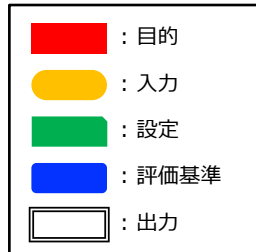
意思決定写像の候補の絞り込み

- 目的, 設定, 評価基準を具体化し, それらの候補を挙げる

説明変数の組						統計的モデル	
						$y_i = \beta_0$	$+ \varepsilon_i$
x_1						$y_i = \beta_0 + \beta_1 x_{i1}$	$+ \varepsilon_i$
	x_2					$y_i = \beta_0 + \beta_2 x_{i2}$	$+ \varepsilon_i$
					\vdots	\vdots	
	x_2	x_3	x_4	x_5	x_6	$y_i = \beta_0 + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \varepsilon_i$	
x_1	x_2	x_3	x_4	x_5	x_6	$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \varepsilon_i$	

を推定する

電車時間 x_3 (量的変数),
駅数, 駐車施設 x_6 (質的変数)



有効な説明変数の組が未知

$$x_{j_1}, x_{j_2}, \dots, x_{j_k}$$

統計的モデルを仮定する
 $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_6 x_{i6} + \varepsilon_i$
 $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$
 σ^2 は未知

尤度関数を最大にする

$$\hat{\beta}_0, \dots, \hat{\beta}_6$$

統計的モデルを仮定せず, 構造として
以下の式を仮定する
 $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_6 x_{i6}$

二乗誤差を最小にする

$$\hat{\beta}_0, \dots, \hat{\beta}_6$$

統計的モデルを仮定する
 $y_i = \beta_0 + \beta_1 x_{ij_1} + \dots + \beta_k x_{ij_k} + \varepsilon_i$
 $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$
 σ^2 は未知

回帰係数が0かを検定して
変数を選択する

$$x_{j_1}, \dots, x_{j_{\hat{k}}}$$

$$\hat{\beta}_0, \dots, \hat{\beta}_{\hat{k}}$$

統計的モデルを仮定しない

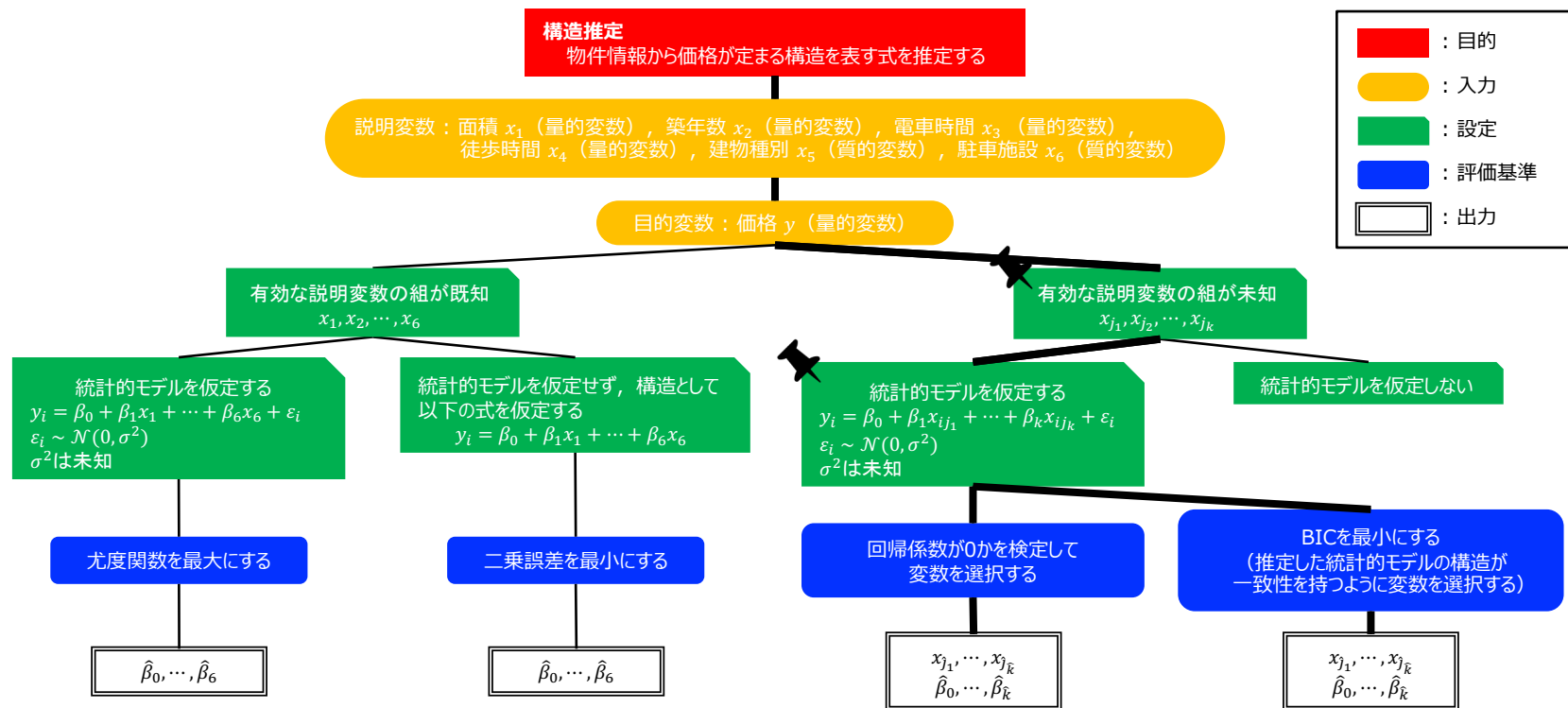
BICを最小にする
(推定した統計的モデルの構造が
一致性を持つように変数を選択する)

$$x_{j_1}, \dots, x_{j_{\hat{k}}}$$

$$\hat{\beta}_0, \dots, \hat{\beta}_{\hat{k}}$$

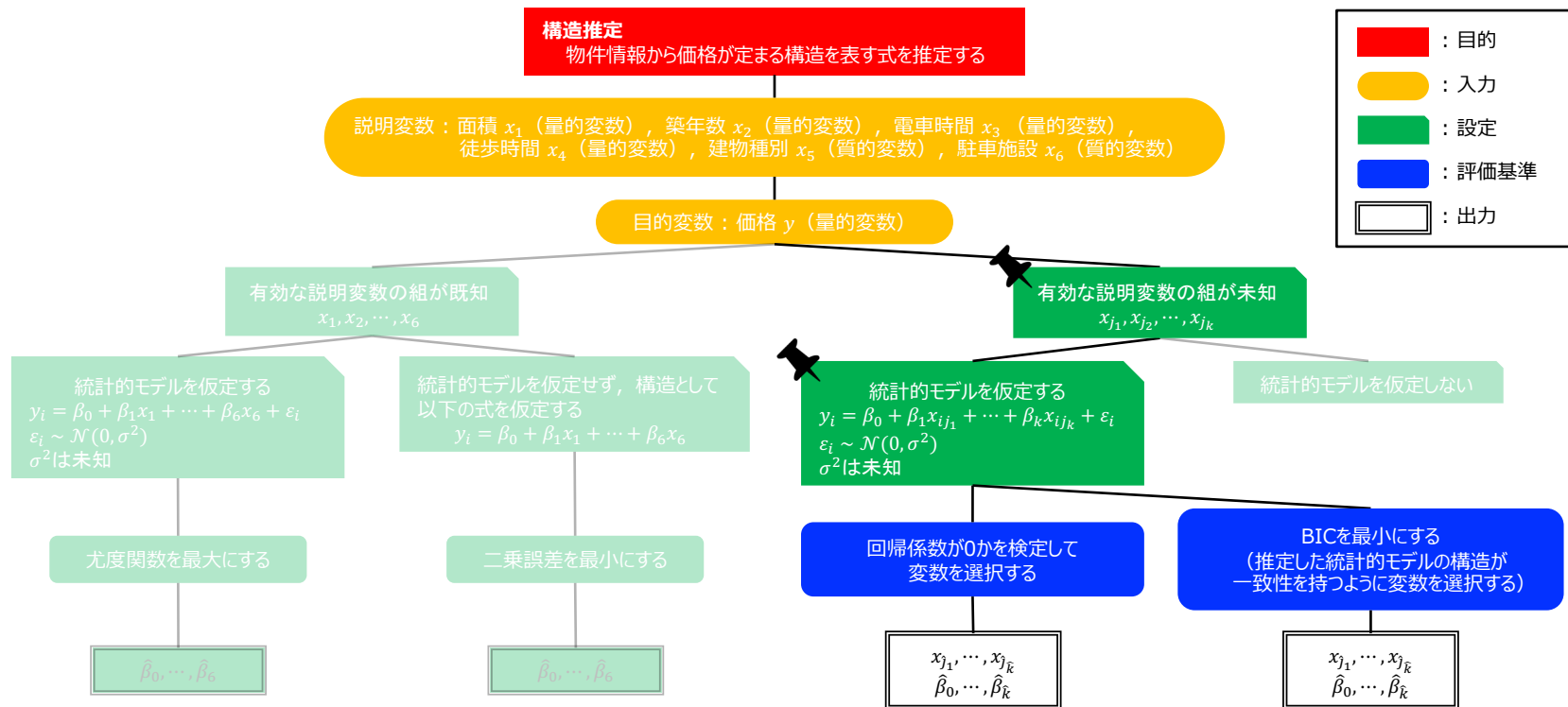
意思決定写像の候補の絞り込み

- 候補の中から、目的、設定、評価基準を定める



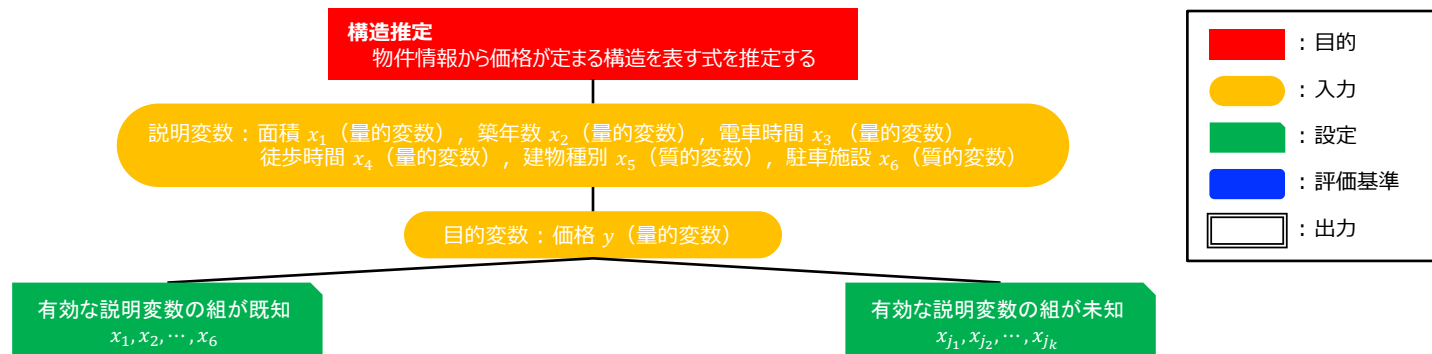
意思決定写像の候補の絞り込み

- 候補の中から、目的、設定、評価基準を定める



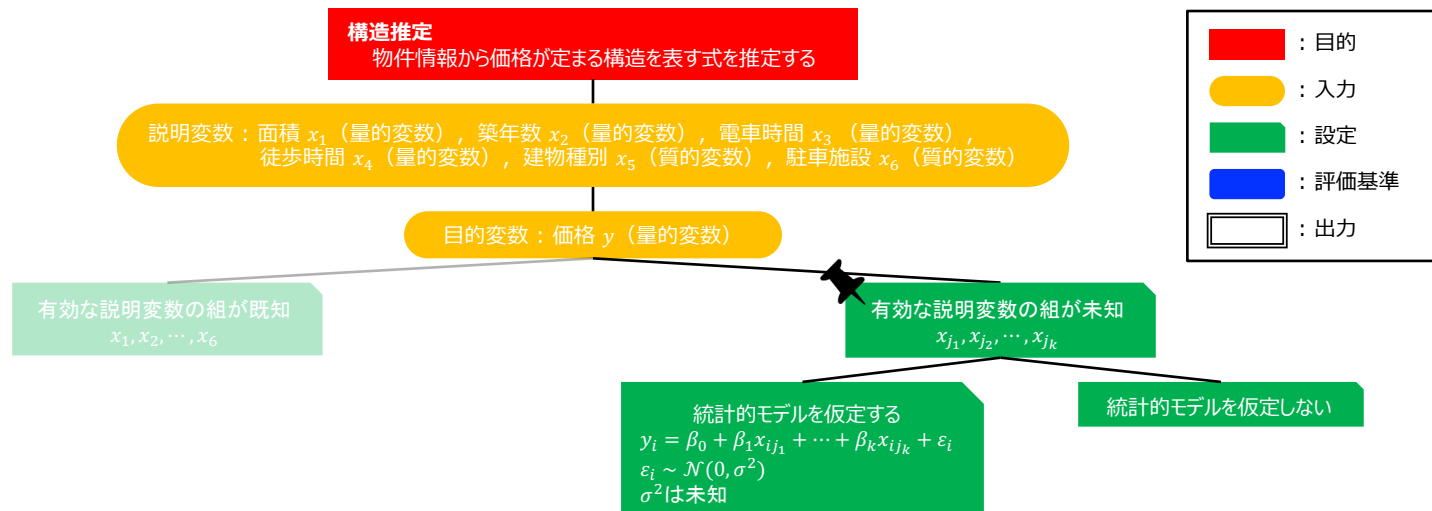
意思決定写像の候補の絞り込み

すべての分岐を考えなくてもよい



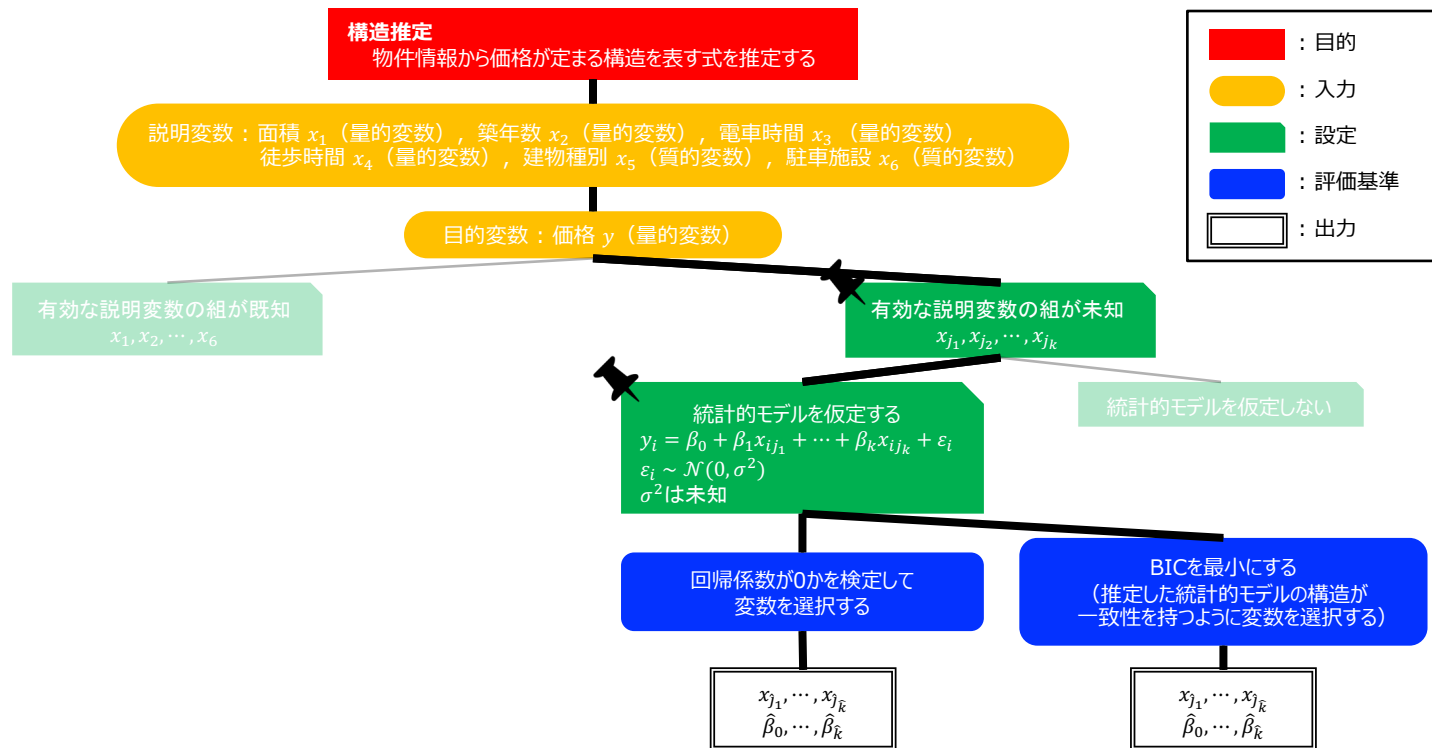
意思決定写像の候補の絞り込み

すべての分岐を考えなくてもよい



意思決定写像の候補の絞り込み

すべての分岐を考えなくてもよい



- 「主目的の把握」と「意思決定写像の候補の絞り込み」のプロセスを通して、どのように意思決定写像を定めるかを説明
- 「意思決定写像構成木」のような考え方によって、意思決定写像を定めることができる