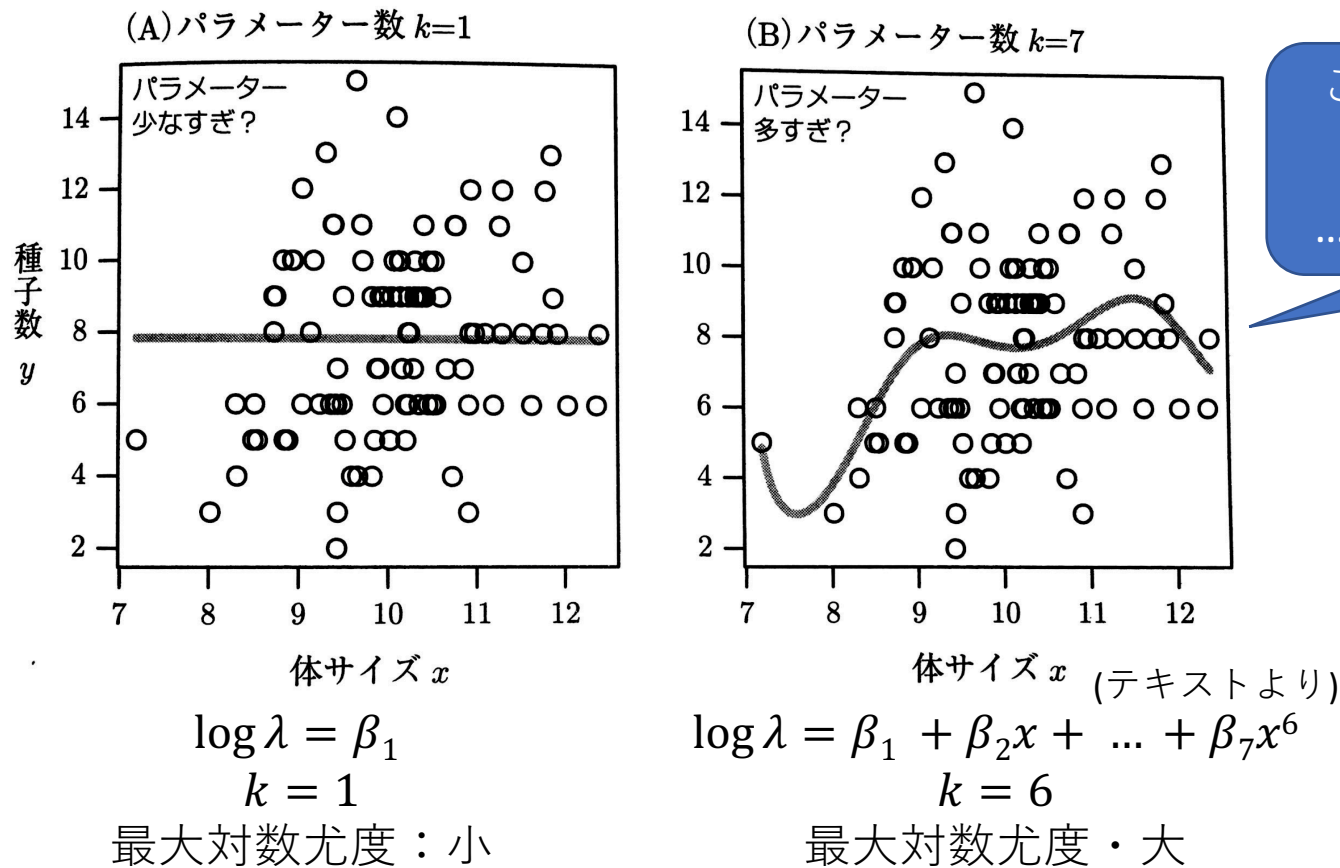


第4章:GLMのモデル選択(前半)

太田研究室 学部4年 和田

章の概要：「良い」モデルとは何か？どのような規準で選択すればいいか？

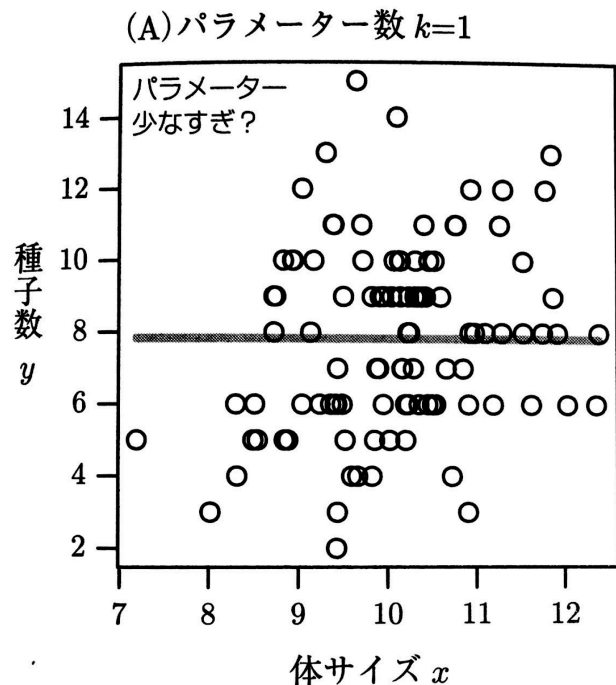
パラメーター数(k)を多くすればする分、最大対数尤度は大きくなる



こちらの方がモデルとして
優れている
..... とは限らない!

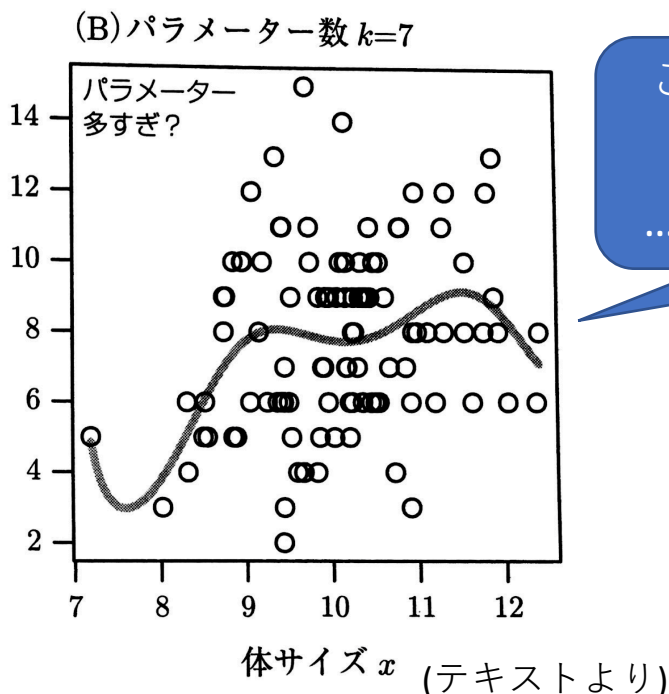
章の概要：「良い」モデルとは何か？どのような基準で選択すればいいか？

パラメーター数(k)を多くすればする分、**最大対数尤度**は大きくなる
観測データへのあてはまりの良さ



$$\log \lambda = \beta_1$$
$$k = 1$$

最大対数尤度：小



$$\log \lambda = \beta_1 + \beta_2 x + \dots + \beta_7 x^6$$
$$k = 6$$

最大対数尤度：大

こちらの方がモデルとして
優れている
..... とは限らない!

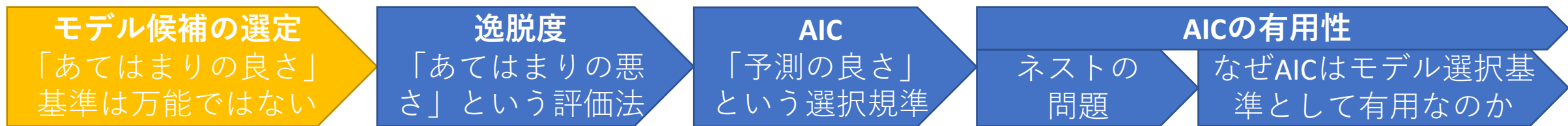
【理由】**要確認**

- ・ 計算処理
- ・ 実際の現象との乖離

最大対数尤度以外の、新しいモデルの評価法・選択規準：

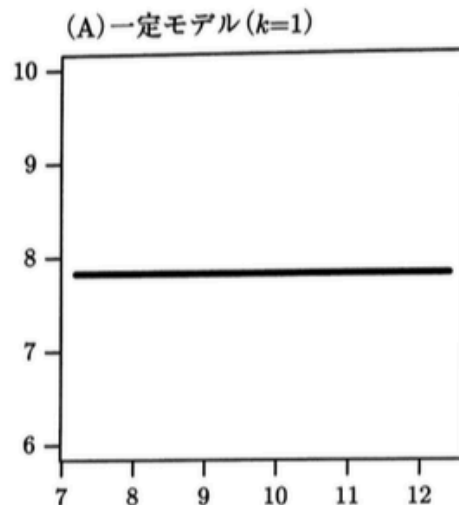
「当てはまりの**悪さ**」→**逸脱度**
「そのモデルは良い**予測**をするのか？」→ **AIC**

4.1 データはひとつ、モデルはたくさん

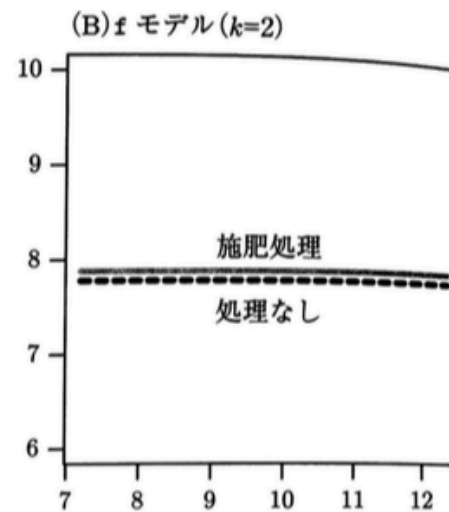


一つのデータに対し、考慮する説明変数のパターン(=候補となるモデル)はたくさんある

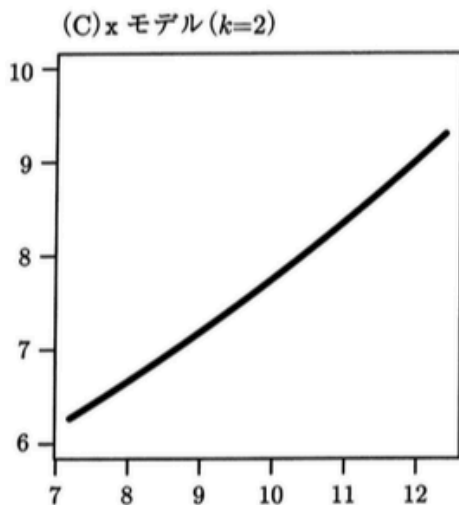
体のサイズ x_i も施肥の有無も、種子の量 y_i に影響しない



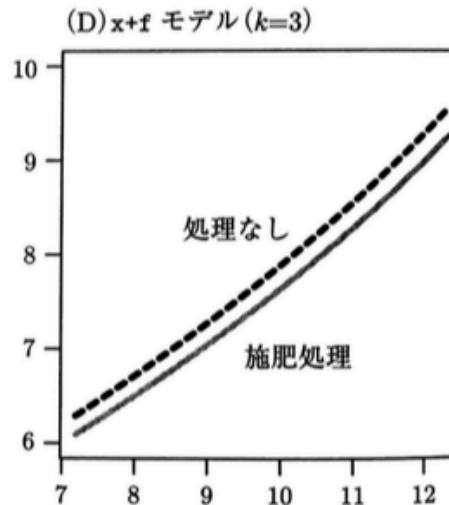
施肥の有無 f_i のみが種子の量 y_i に影響する



体のサイズ x_i のみが種子の量 y_i に影響する



施肥の有無 f_i と、体のサイズ x_i の両方が、種子の量 y_i に影響する



(テキストより)

一つのデータに対し、考慮する説明変数のパターン(=候補となるモデル)はたくさんある

このうち、どれを採用すべきか？
を考えたときに.....

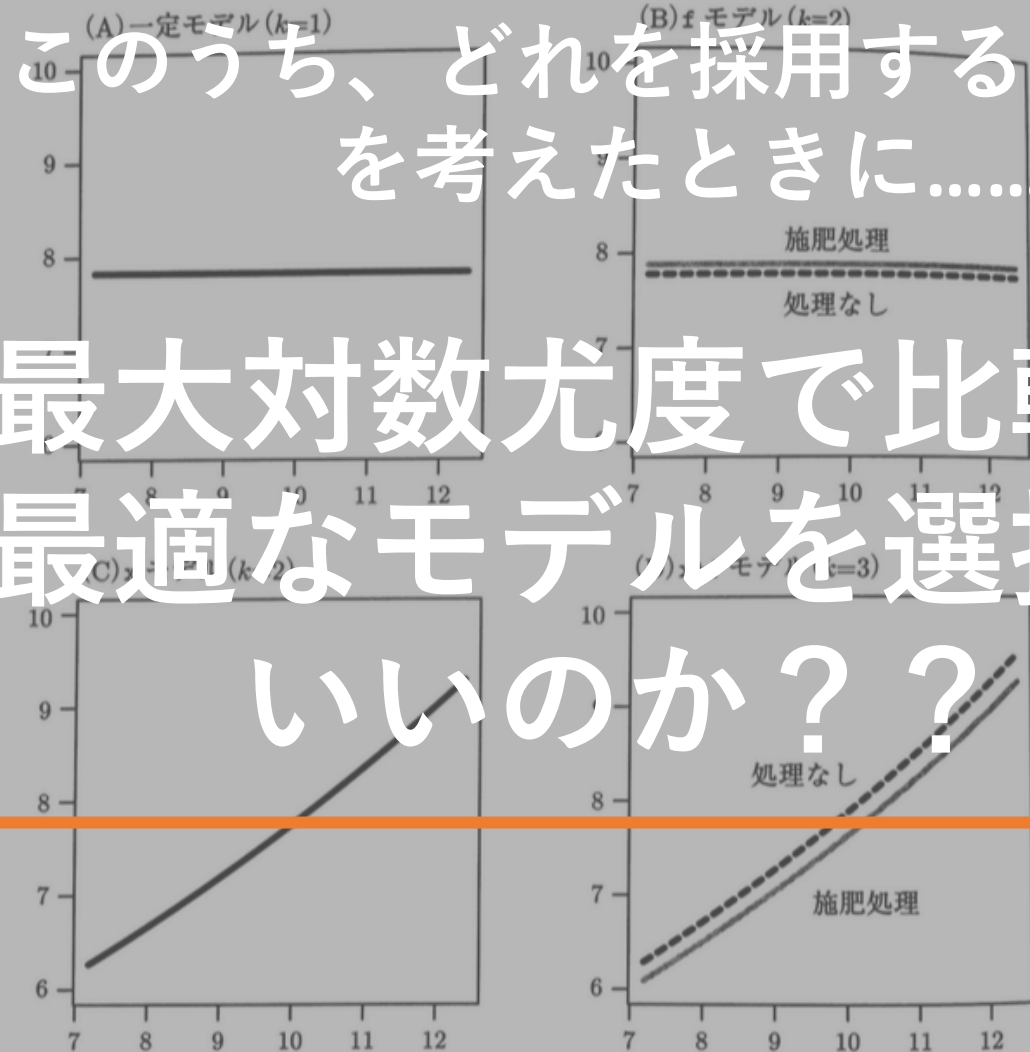
体のサイズ x_i も施肥の有無も、種子の量 y_i に影響しない

施肥の有無 f_i のみが種子の量 y_i に影響する

最大対数尤度で比較し、
最適なモデルを選択して
いいのか？？？

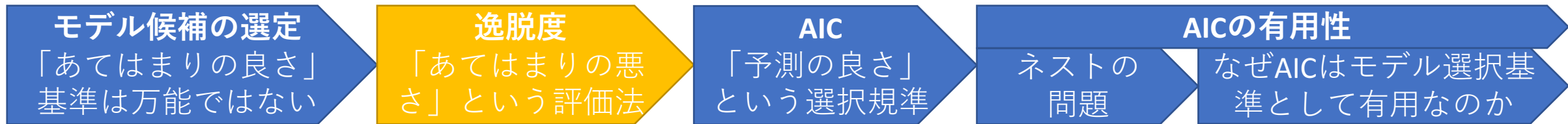
体のサイズ x_i のみが種子の量 y_i に影響する

施肥の有無 f_i と、体のサイズ x_i の両方が、種子の量 y_i に影響する



(テキストより)

4.2 統計モデルのあてはまりの悪さ：逸脱度



モデルのデータへのあてはまりの悪さ 「逸脱度」は、最大対数尤度の変形

- 「逸脱度」 (Deviance) D
- 統計モデルの、データへの「あてはまりの悪さ」の指標

$$D = -2 \log L^*$$

$\log L(\{\beta_j\})$ を $\log L$ 、
その最大対数尤度を $\log L^*$ と表記

glm() コマンドの出力結果に表示

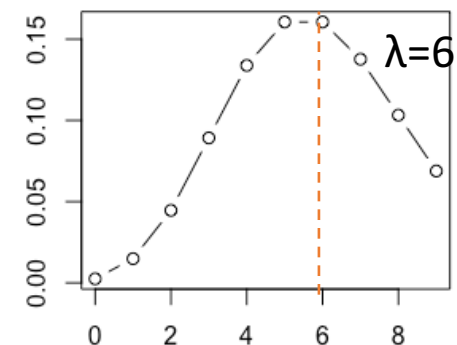
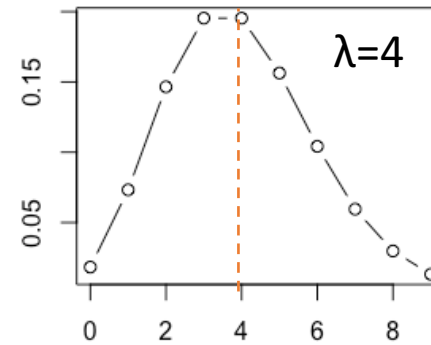
名前	In English	定義
逸脱度 (D)	Deviance	$-2 \log L^*$
最小の逸脱度	Minimum deviance	フルモデル(後述)をあてはめたときのD
残差逸脱度	Residual deviance	D – 最小のD
最大の逸脱度	Maximum deviance	Nullモデル(後述)をあてはめたときのD
Null 逸脱度	Null deviance	最大のD – 最小のD

フルモデル、Nullモデルはそれぞれ、パラメータ数を最大、最小(1)にした場合のモデルである①

「フルモデル」 (full model) ... 最もあてはまりが**いい**モデル

- 個々のデータに、一対一対応でパラメータ λ が定まっている
 - 100個のデータがあれば100個の λ を定めている

例: $y_i = \{4, 4, 4, 6, \dots\}$ のとき、
 $i \in \{1, 2, 3\}$ の y_i は4なので、 $\{\lambda_1, \lambda_2, \lambda_3\} = \{4, 4, 4\}$
 $i = 4$ の y_4 は6なので、 $\lambda_4 = 6$
...(以後同上)



- (同じ回帰で)他のどのモデルを使った時よりも、必然的に最大対数尤度は最大、逸脱度は最小になる
 - フルモデルを当てはめた時の逸脱度 = 最小のD(minimum deviance)
- 「現象を説明しうる理想のモデルを考えている」のではなく、「現在のデータ(のみ)にモデルを近づけている」ので、モデルとしての価値はない

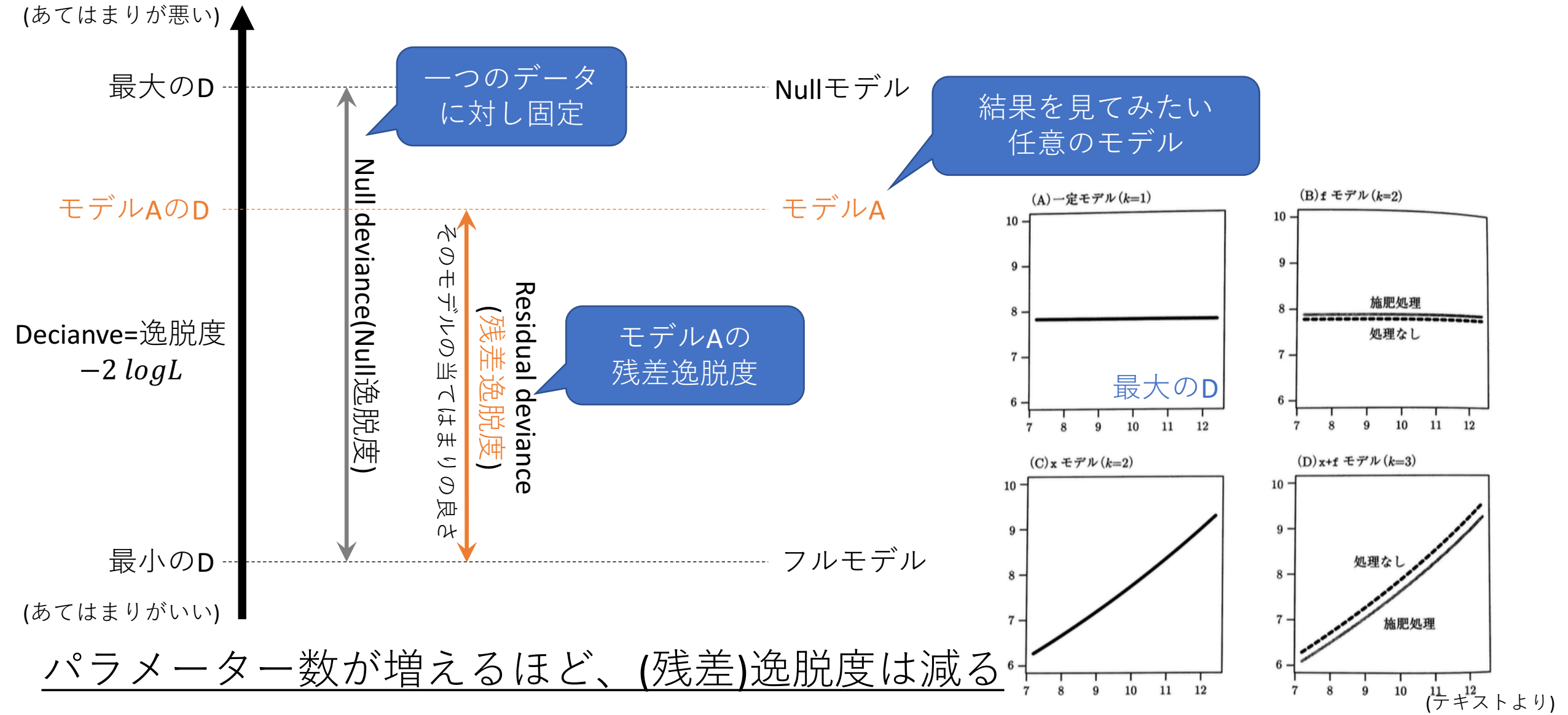
フルモデル、Nullモデルはそれぞれ、パラメーター数を最大、最小(1)にした場合のモデルである②

「Null モデル」 (Null model) ... 最もあてはまりが**悪い**モデル

- パラメーター数が1
 - つまり、この文脈においては $\lambda = \exp(\beta_1)$
- パラメーターは、全ての説明変数から完全に独立である
- (同じ回帰で)他のどのモデルを使った時よりも、必然的に対数尤度は最小、逸脱度は最大になる

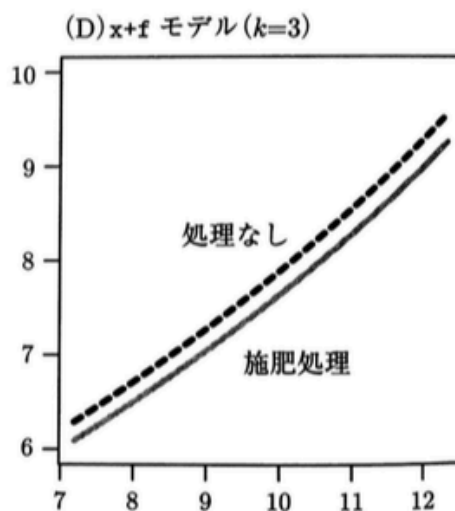
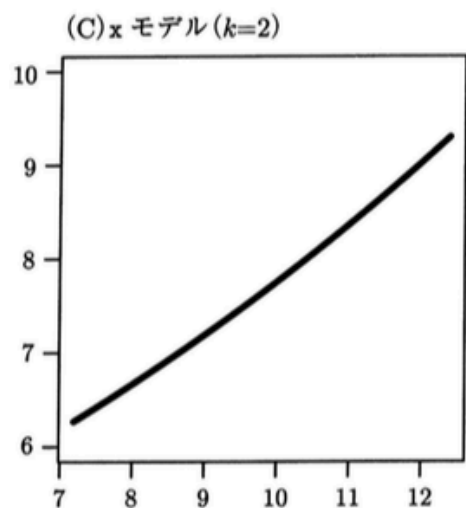
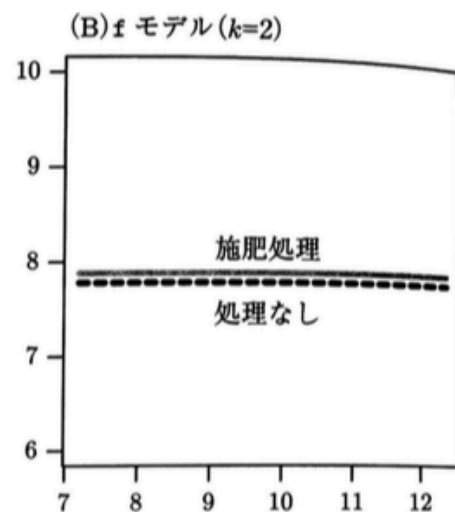
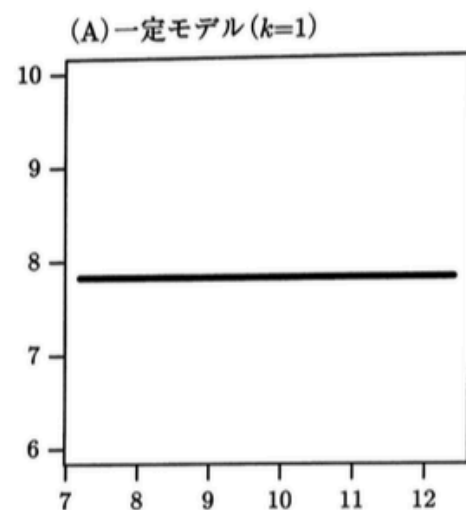
Null モデルを当てはめた場合の逸脱度
= 最大のD (Maximum deviance)

種々の逸脱度の関係性は以下



R実践：残差逸脱度の算出・比較

R実践：残差逸脱度の算出・比較



(テキストより)

```
> fit.xf

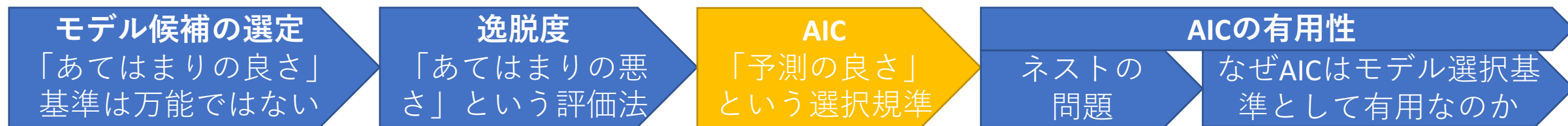
Call:  glm(formula = y ~ x + f, family = poisson, data = d)

Coefficients:
(Intercept)          x          fT
    1.26311      0.08007     -0.03200

Degrees of Freedom: 99 Total (i.e. Null);  97 Residual
Null Deviance:      89.51
Residual Deviance:  84.81      AIC: 476.6
```

モデル	k	logL*	Deviance	Residual D
A:一定	1	-237.6	475.3	89.51
B: f	2	-237.6	475.3	89.48
C: x	2	-235.4	470.8	84.99
D: x+f	3	-235.3	470.6	84.81
フル	100	-192.9	385.8	0.0

4.3 モデル選択規準AIC



AICの比較により「予測の良さ」を重視した モデル選択を行うことができる

AIC (Akaike's information criterion)

- 「モデル選択規準」(model selection criterion)の一つ
- 予測の良さを重視する (~~当てはまりの良さ~~)
- 小さい方が「良い」モデル

$$\begin{aligned} AIC &= -2\{(\text{最大対数尤度}) - (\text{パラメーター数})\} \\ &= -2(\log L^* - k) \\ &= \boxed{-2\log L^*} + 2k \end{aligned}$$

要確認

「基準」と「規準」 の違いとは？

「**規準**」...何を測定するか。
何の指標を用いるか。
辞書「何かを行う際に手本・標準
とすべきもの」
Eg「**道徳の規準**」「**社会生活の規
準**」

「**基準**」...どこまで達成でき
たか。測定された値に基づく
評価。
辞書「物事を判断するためのより
どころ。標準と見なす数値など」
Eg「**選考基準**」「**前年を基準に予
算を決める**」

AICの比較により「予測の良さ」を重視した モデル選択を行うことができる

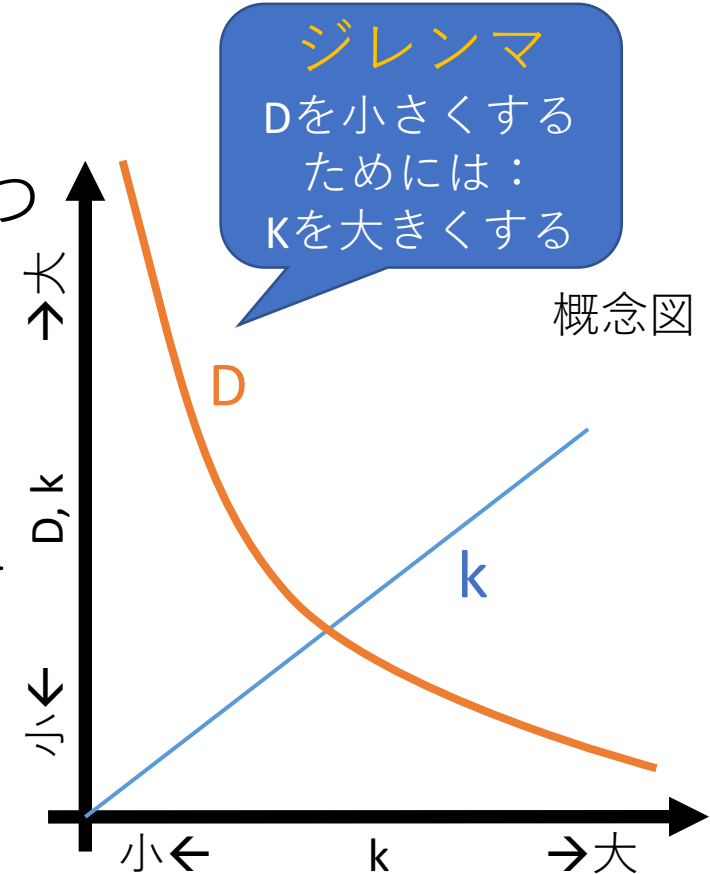
AIC (Akaike's information criterion)

- 「モデル選択規準」(model selection criterion)の一つ
- 予測の良さを重視する (~~当てはまりの良さ~~)
- 小さい方が「良い」モデル

$$\begin{aligned} AIC &= -2\{(\text{最大対数尤度}) - (\text{パラメーター数})\} \\ &= -2(\log L^* - k) \\ &= \mathbf{D + 2k} \end{aligned}$$

これを小さくしたい

右辺の二項、両方を小さくしたい



逸脱度とパラメーター数の両方が小さい→「良い」モデル

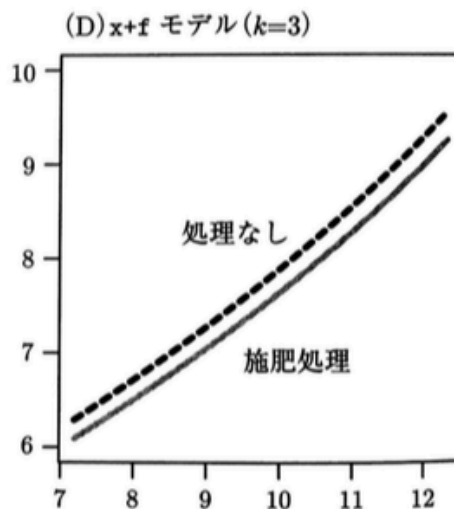
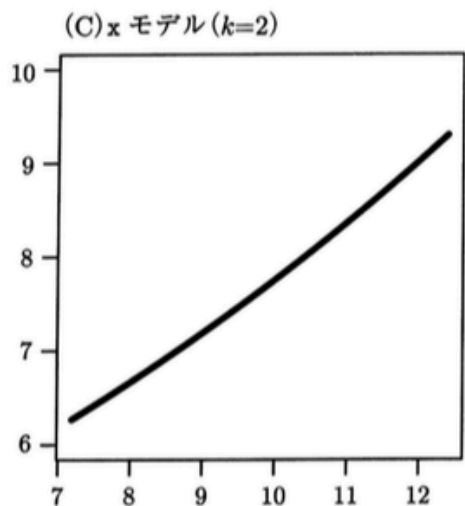
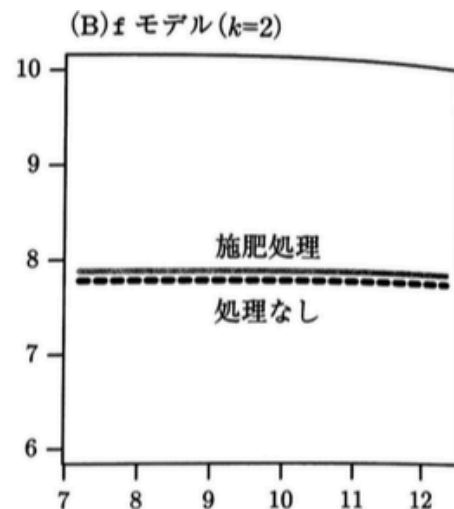
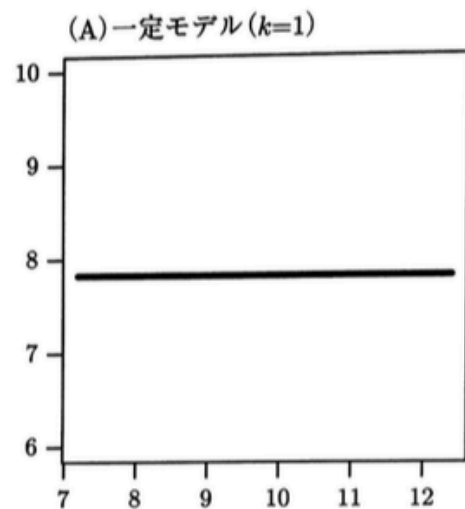
疑問

逸脱度の概念は必要か？ 最大対数尤度じゃダメなのか？
残差逸脱度の概念が必要なのは分かるのだが...

(AICが選択規準として有用である理由は以後)

R実践：AICの算出・比較

R実践：AICの算出・比較



(テキストより)

```
> fit.xf
```

```
Call: glm(formula = y ~ x + f, family = poisson, data = d)
```

```
Coefficients:
```

```
(Intercept)          x          fT  
    1.26311    0.08007   -0.03200
```

```
Degrees of Freedom: 99 Total (i.e. Null); 97 Residual
```

```
Null Deviance:      89.51
```

```
Residual Deviance: 84.81
```

```
AIC: 476.6
```

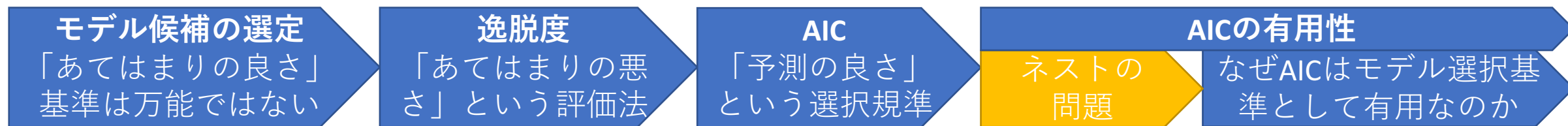
モデル	k	logL*	Deviance	Residual D	AIC
A:一定	1	-237.6	475.3	89.5	477.3
B: f	2	-237.6	475.3	89.5	479.3
C: x	2	-235.4	470.8	85.0	474.8
D: x+f	3	-235.3	470.6	84.8	476.6
フル	100	-192.9	385.8	0.0	585.8

残差逸脱度がDより大きいのに、AICはCが最小

ここまでのまとめ

- 最大対数尤度は「あてはまりの良さ」を表現している
- 最大対数尤度は、モデルの良さの評価規準としてふさわしいのか？？
 - 「一つのデータセット」の説明に終始→実際の現象と乖離する可能性が高い
- 逸脱度 = 「あてはまりの悪さ」という値でモデルを評価可能
- 残差逸脱度は、他のモデルの逸脱度と比較した相対的な規準である
 - フルモデル採用時の逸脱度=最小のDとの比較
- モデル選択規準 AICの値が小さいほど「良い」モデルと言える
- AICは「モデルの予測の良さ」を判断規準とする
- パラメーター数 k と逸脱度 D の双方が程よく小さい時：AICが最小
- 「良い予測をする」モデルは、 k と D が小さい

4.4 AICを説明するためのまた別の問題



モデルAがモデルBの部分集合であるとき、
モデルA,Bは「**ネストしている**」

「ネストしている/した」 (nested)

- 一方のモデルが他方に含まれている状態

$$\begin{array}{ll} \text{モデルA:} & \log \lambda = \beta_1 + \beta_2 x_1 \\ \text{モデルB:} & \log \lambda = \beta_1 + \beta_2 x_1 + \beta_3 x_2 \end{array}$$

- 上記の場合、モデルAはモデルBの部分集合
= 「モデルAとBはネストしている」
- どんなモデルも、フルモデルの部分集合
= フルモデルは、必然的にどんなモデルともネストしている
- ネストしたモデル間の検定：Wald検定、尤度比検定など

4.5 なぜAICでモデル選択して良いのか？

