

九州大学 工学部地球環境工学科
船舶海洋システム工学コース

海事統計学 第14回 (担当:木村)

検定(4):適合度の検定 K-S検定

授業の資料等は

<http://sysplan.nams.kyushu-u.ac.jp/gen/index.html>

検定によってどんなことが判定できるか？

・比率の検定

- (1) 母比率 P が、ある値 P_0 に等しいといえるか？
- (2) 比率の差の検定： 2つの異なる母集団の間で、母比率に差があるといえるか？

・平均値の検定(正規母集団)

- (1) 母集団の平均値 μ が、ある値 μ_0 に等しいといえるか？
- (2) 平均値の差の検定： 2つの異なる母集団の間で、母平均に差があるといえるか？

・分散の検定

- (1) **正規母集団**の分散 σ^2 が、ある値 σ_0^2 に等しいといえるか？
- (2) 分散の差の検定： 2つの異なる**正規母集団**の間で、分散に差があるといえるか？

・適合度の検定

- (1) 観察されたデータが、特定の分布に一致しているといえるか？

- (2) 2つの母集団の確率分布が異なるものであるかどうか？

分布の種類を問わない

(ノンパラメトリック)

← コルモゴロフ・スミルノフ検定

【復習】 確率変数のとりうる値が離散かつ無限個の場合

ポアソン分布 (Poisson distribution)

発生時間間隔は
指数分布

単位時間あたり平均で λ 回発生する事象が、
単位時間に x 回(ゼロを含む)発生する確率

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

二項分布において、

p が微小で、 n が極めて大きい

$\lambda = n p$ とおく

確率 p を持つ事象が n 回の試行中 x 回おこることを考えた

$$P(x) = {}_n C_x p^x (1-p)^{n-x} = \frac{n(n-1)\cdots(n-x+1)}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{\lambda^x}{x!} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x-1}{n}\right) \left(1 - \frac{\lambda}{n}\right)^{-x}$$

ここで

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n = e^{-\lambda}, \quad \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{x-1}{n}\right) = 1, \quad \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-x} = 1,$$

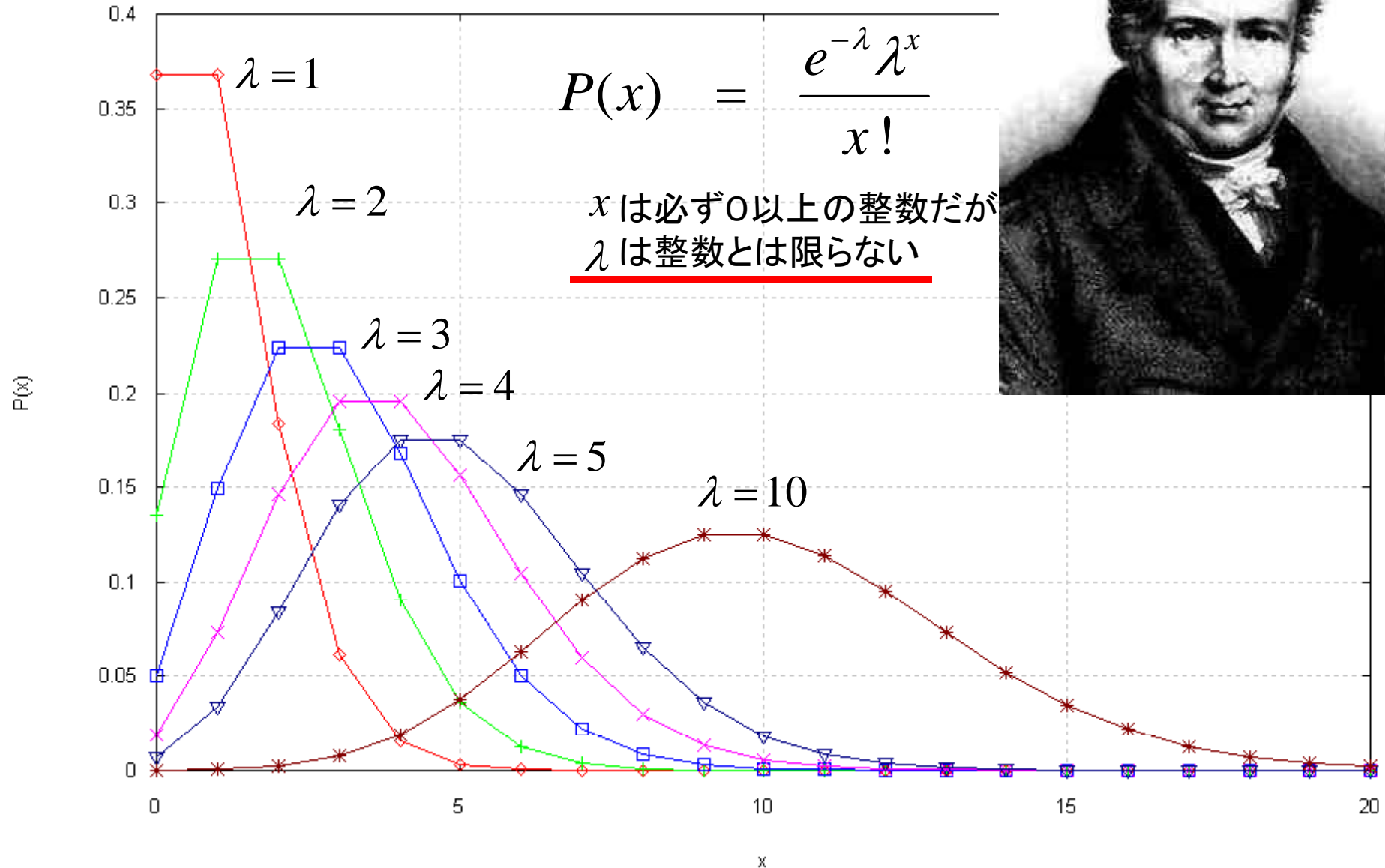
よって

$$P(x) = \lim_{n \rightarrow \infty} {}_n C_x p^x (1-p)^{n-x} = \frac{e^{-\lambda} \lambda^x}{x!}$$

ポアソン分布は
二項分布の極限

Simeon Denis Poisson (1781 – 1840)
(France)

ポアソン分布 (Poisson distribution)



【復習】 確率変数のとりうる値が離散かつ無限個の場合

ポアソン分布 (Poisson distribution)

発生時間間隔は
指数分布

単位時間あたり平均で λ 回発生する事象が、
単位時間に x 回(ゼロを含む)発生する確率

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

ポアソン分布の現れる例:

備考: $x = 0$ のとき $x! = 1$

- 1分間に放射性物質から放射される粒子が平均2個観測されるとき、1分間に1個も観測されない確率は？

- ある地域において、年間の交通事故件数が平均730件のとき、1日に発生する事故の件数が0件である確率は？

$$\lambda = 2 \quad P(0) = \frac{e^{-2} 2^0}{0!} = e^{-2} = 0.135$$

- ある機械で部品を 10000個作ると、平均 5個の不良品ができるとき、部品を 1000個作ったときに不良品が0個である確率は？ 不良品が1個出る確率は？

指数分布 と ポアソン分布 は同一現象の異なる表現

λ : 平均到着(故障)率(人/時) $\Leftrightarrow \frac{1}{\lambda}$: 平均到着(故障)時間間隔

3. 分布の適合度の検定(カイ2乗検定)

期待度数に関する未知母数をデータから推定する場合(最尤推定)

前回とは
ここが違う!

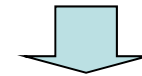
【例】ある地域の交通事故件数を100日間調査:

事故件数(クラス) k	0	1	2	3	4	5	計
日数 (観測度数) f_k	43	31	14	8	3	1	100

問: 1日あたりの事故件数はポアソン分布に従っているといえるかどうか
有意水準5%で検定せよ。

1) まず各クラスにおける期待度数を求めなければならない。

この値が分からない
未知母数



ポアソン分布

単位時間あたり平均で λ 回発生する事象が、
単位時間に x 回(ゼロを含む)発生する確率

$$P(x) =$$

クラス k

期待度数は、これに全度数(100)を
乗じれば求めることができる

3. 分布の適合度の検定(カイ2乗検定)

期待度数に関する未知母数をデータから推定する場合(最尤推定)

前回とは
ここが違う!

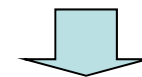
【例】ある地域の交通事故件数を100日間調査:

事故件数(クラス) k	0	1	2	3	4	5	計
日数 (観測度数) f_k	43	31	14	8	3	1	100

問: 1日あたりの事故件数はポアソン分布に従っているといえるかどうか
有意水準5%で検定せよ。

1) まず各クラスにおける期待度数を求めなければならない。

この値が分からない
未知母数



データから推定

ポアソン分布

単位時間あたり平均で λ 回発生する事象が、

単位時間に x 回(ゼロを含む)発生する確率

$$P(x) =$$

クラス k

期待度数は、これに全度数(100)を
乗じれば求めることができる

3. 分布の適合度の検定(カイ2乗検定)

期待度数に関する未知母数をデータから推定する場合(最尤推定)

前回とは
ここが違う!

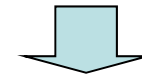
【例】ある地域の交通事故件数を100日間調査:

事故件数(クラス) k	0	1	2	3	4	5	計
日数 (観測度数) f_k	43	31	14	8	3	1	100

問: 1日あたりの事故件数はポアソン分布に従っているといえるかどうか
有意水準5%で検定せよ。

1) まず各クラスにおける期待度数を求めなければならない。

この値が分からない
未知母数



データから推定

ポアソン分布

単位時間あたり平均で λ 回発生する事象が、
単位時間に x 回(ゼロを含む)発生する確率

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

クラス k

期待度数は、これに全度数(100)を
乗じれば求めることができる

期待度数に関する未知母数をデータから推定する場合(つづき1)

標本平均 \bar{x} をポアソン分布の未知母数 λ の推定値 $\hat{\lambda}$ とすると、

$$\hat{\lambda} = \frac{1}{100} (0 \times 43 + 1 \times 31 + 2 \times 14 + 3 \times 8 + 4 \times 3 + 5 \times 1) = 1$$

ポアソン分布の
「最尤推定」を参照

よって $P(x) = \frac{e^{-1}}{x!}$ この式から期待度数を計算すると

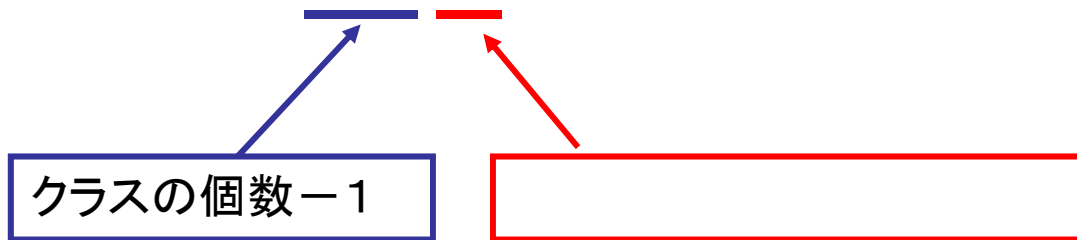
事故件数(クラス) k	0	1	2	3	4	5	6以上	計
日数 (観測度数) f_k	43	31	14	8	3	1	0	100
期待度数 f_k^*	36.8	36.8	18.4	6.13	1.53	0.307	0.0330	100

クラスk=4以上は
期待度数が小さすぎるので
「クラスk=3以上」にまとめる

期待度数に関する未知母数をデータから推定する場合(つづき2)

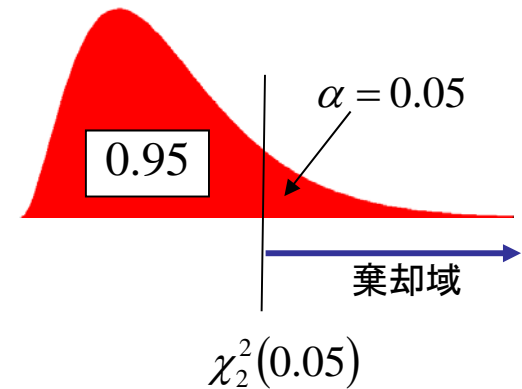
事故件数(クラス) k		0	1	2	3以上	計
日数 (観測度数) f_k		43	31	14	12	100
期待度数 f_k^*		36.8	36.8	18.4	8.0	100

このとき、自由度 $4 - 1 - 1 = 2$ のカイ2乗分布表から有意水準5%の棄却域を読み取る



$$\chi^2 = \frac{(43-36.8)^2}{36.8} + \frac{(31-36.8)^2}{36.8} + \frac{(14-18.4)^2}{18.4} + \frac{(12-8.0)^2}{8.0}$$

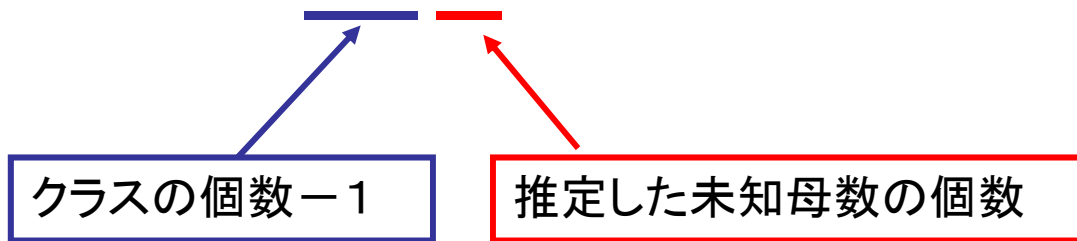
$$= 5.011$$



期待度数に関する未知母数をデータから推定する場合(つづき2)

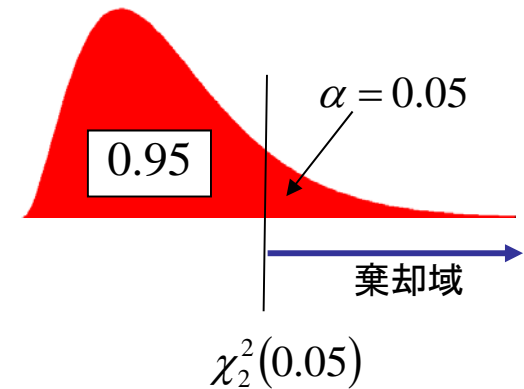
事故件数(クラス) k	0	1	2	3以上	計
日数 (観測度数) f_k	43	31	14	12	100
期待度数 f_k^*	36.8	36.8	18.4	8.0	100

このとき、自由度 $4 - 1 - 1 = 2$ のカイ2乗分布表から有意水準5%の棄却域を読み取る

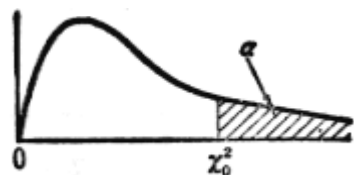


$$\chi^2 = \frac{(43-36.8)^2}{36.8} + \frac{(31-36.8)^2}{36.8} + \frac{(14-18.4)^2}{18.4} + \frac{(12-8.0)^2}{8.0}$$

$$= 5.011$$



χ^2 分布表 $\alpha = P(X > \chi_0^2) \rightarrow \chi_0^2$

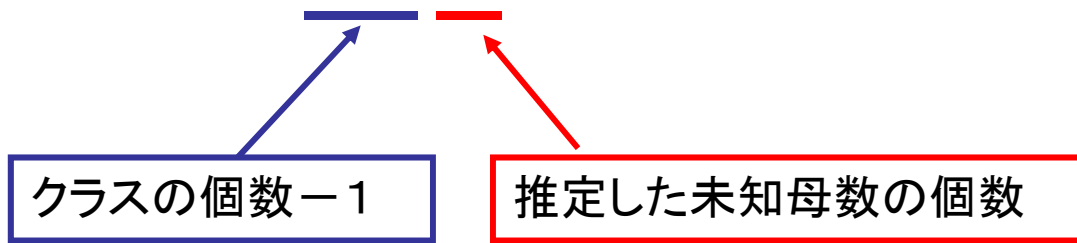


自由度 m	α									
	0.995	0.99	0.975	0.95	0.90	0.10	0.05	0.025	0.01	0.005
1	0.000	0.000	0.001	0.003	0.016	2.71	3.84	5.02	6.63	7.88
2	0.010	0.020	0.051	0.103	0.211	4.61	5.99	7.38	9.21	10.60
3	0.072	0.115	0.216	0.352	0.584	6.25	7.81	9.35	11.34	12.84
4	0.207	0.297	0.484	0.711	1.064	7.78	9.49	11.14	13.28	14.86
5	0.412	0.554	0.831	1.145	1.610	9.24	11.07	12.83	15.09	16.75
6	0.676	0.872	1.237	1.635	2.20	10.64	12.59	14.45	16.81	18.55
7	0.989	1.239	1.690	2.17	2.83	12.02	14.07	16.01	18.48	20.3
8	1.344	1.646	2.18	2.73	3.49	13.36	15.51	17.53	20.1	22.0
9	1.735	2.09	2.70	3.33	4.17	14.68	16.92	19.02	21.7	23.6
10	2.16	2.56	3.25	3.94	4.87	15.99	18.31	20.5	23.2	25.2
11	2.60	3.05	3.82	4.57	5.58	17.28	19.68	21.9	24.7	26.8
12	3.07	3.57	4.40	5.23	6.30	18.55	21.0	23.3	26.2	28.3
13	3.57	4.11	5.01	5.89	7.04	19.81	22.4	24.7	27.7	29.8
14	4.07	4.66	5.63	6.57	7.79	21.1	23.7	26.1	29.1	31.3
15	4.60	5.23	6.26	7.26	8.55	22.3	25.0	27.5	30.6	32.8
16	5.14	5.81	6.91	7.96	9.31	23.5	26.3	28.8	32.0	34.3
17	5.70	6.41	7.56	8.67	10.09	24.8	27.6	30.2	33.4	35.7
18	6.26	7.01	8.23	9.39	10.86	26.0	28.9	31.5	34.8	37.2
19	6.84	7.63	8.91	10.12	11.65	27.2	30.1	32.9	36.2	38.6
20	7.43	8.26	9.59	10.85	12.44	28.4	31.4	34.2	37.6	40.0
21	8.03	8.90	10.28	11.59	13.24	29.6	32.7	35.5	38.9	41.4
22	8.64	9.54	10.98	12.34	14.04	30.8	33.9	36.8	40.3	42.8
23	9.26	10.20	11.69	13.09	14.85	32.0	35.2	38.1	41.6	44.2
24	9.89	10.86	12.40	13.85	15.66	33.2	36.4	39.4	43.0	45.6
25	10.52	11.52	13.12	14.61	16.47	34.4	37.7	40.6	44.3	46.9
26	11.16	12.20	13.84	15.38	17.29	35.6	38.9	41.9	45.6	48.3
27	11.81	12.88	14.57	16.15	18.11	36.7	40.1	43.2	47.0	49.6
28	12.46	13.56	15.31	16.93	18.94	37.9	41.3	44.5	48.3	51.0
29	13.12	14.26	16.05	17.71	19.77	39.1	42.6	45.7	49.6	52.3
30	13.79	14.95	16.79	18.49	20.6	40.3	43.8	47.0	50.9	53.7

期待度数に関する未知母数をデータから推定する場合(つづき2)

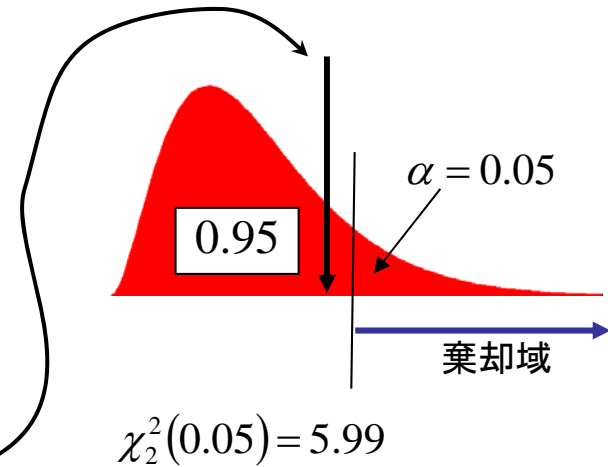
事故件数(クラス) k	0	1	2	3以上	計
日数 (観測度数) f_k	43	31	14	12	100
期待度数 f_k^*	36.8	36.8	18.4	8.0	100

このとき、自由度 $4 - 1 - 1 = 2$ のカイ2乗分布表から有意水準5%の棄却域を読み取る



$$\chi^2 = \frac{(43-36.8)^2}{36.8} + \frac{(31-36.8)^2}{36.8} + \frac{(14-18.4)^2}{18.4} + \frac{(12-8.0)^2}{8.0}$$

$$= 5.011$$



帰無仮説は棄却されない = ポアソン分布に従っていないとは言えない

適合度の検定(カイ2乗検定):まとめ

標本分布および仮説による分布が度数分布で表されているとする。
それぞれの第 k クラスの度数を f_k および f_k^* ただし $k = 1, 2, \dots, m$ とする。
このとき、仮説が正しいならば

統計量 $\chi^2 = \sum_{k=1}^m \left(\frac{(f_k - f_k^*)^2}{f_k^*} \right)$ は近似的に自由度 $m-1$ のカイ2乗分布に従う
ただし、期待度数 $f_k^* \geq 5$

観測度数 ← f_k ← f_k^* 期待度数

- カイ2乗分布表を使って片側検定
- (期待度数 < 5) の場合は、クラスを統合して度数を5以上にする
- 期待度数分布に関して未知母数が存在する場合、
 - (1) データから未知母数を推定
 - (2) カイ2乗分布の自由度は、 とする

適合度の検定(カイ2乗検定):まとめ

標本分布および仮説による分布が度数分布で表されているとする。
それぞれの第 k クラスの度数を f_k および f_k^* ただし $k = 1, 2, \dots, m$ とする。
このとき、仮説が正しいならば

統計量 $\chi^2 = \sum_{k=1}^m \left(\frac{(f_k - f_k^*)^2}{f_k^*} \right)$ は近似的に自由度 $m-1$ のカイ2乗分布に従う
ただし、期待度数 $f_k^* \geq 5$

観測度数 ← f_k ← f_k^* 期待度数

- カイ2乗分布表を使って片側検定
- (期待度数 < 5) の場合は、クラスを統合して度数を5以上にする
- 期待度数分布に関して未知母数が存在する場合、
 - (1) データから未知母数を推定
 - (2) カイ2乗分布の自由度は、(クラス数) - 1 - (未知母数の个数) とする

検定によってどんなことが判定できるか？

・比率の検定

- (1) 母比率 P が、ある値 P_0 に等しいといえるか？
- (2) 比率の差の検定： 2つの異なる母集団の間で、母比率に差があるといえるか？

・平均値の検定(正規母集団)

- (1) 母集団の平均値 μ が、ある値 μ_0 に等しいといえるか？
- (2) 平均値の差の検定： 2つの異なる母集団の間で、母平均に差があるといえるか？

・分散の検定

- (1) **正規母集団**の分散 σ^2 が、ある値 σ_0^2 に等しいといえるか？
- (2) 分散の差の検定： 2つの異なる**正規母集団**の間で、分散に差があるといえるか？

・適合度の検定

- (1) 観察されたデータが、特定の分布に一致しているといえるか？

- (2) 2つの母集団の確率分布が異なるものであるかどうか？

分布の種類を問わない
(ノンパラメトリック)

← コルモゴロフ・スミルノフ検定

Kolmogorov-Smirnov検定 (K-S検定)とは？

(コルモゴロフスミルノフ検定)

【1標本 K-S検定】

標本X

x_1, x_2, \dots, x_n

標本Xが1独立変数の確率密度関数 $f(x)$ と一致していると言えるか？

ここが凄い！

特に条件が無い
正規分布でも指数分布でも何でも良い！

【2標本 K-S検定】

標本X

x_1, x_2, \dots, x_n

標本Y

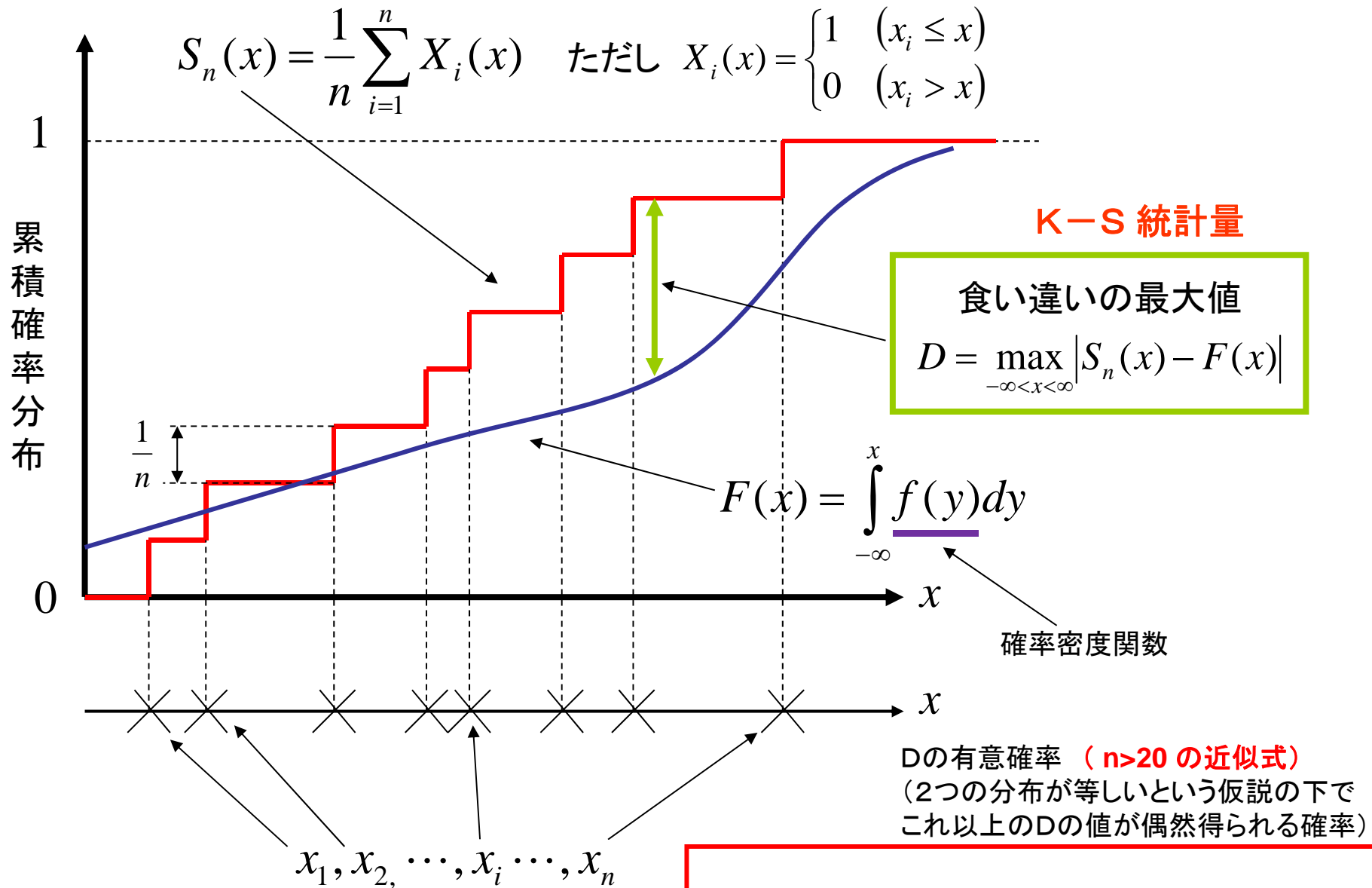
y_1, y_2, \dots, y_m

標本Xと標本Yが同一の母集団の確率密度関数より生じていると言えるか？

※ 非常に使い勝手の良い検定だが、サンプル数に制限あり。
計算には一般にコンピュータの支援が必要

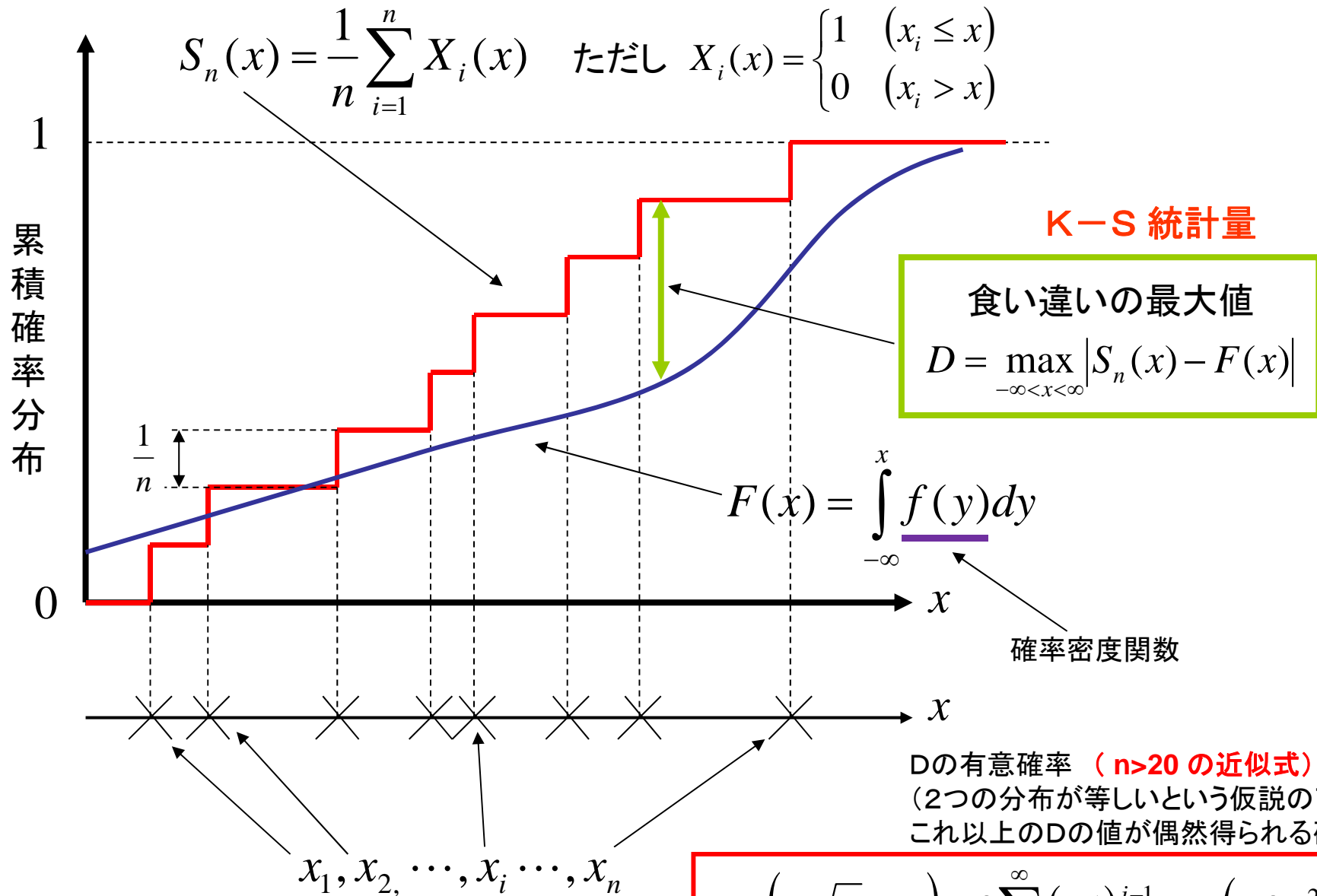
K-S検定の計算手順(1標本)

累積確率分布の差を計算

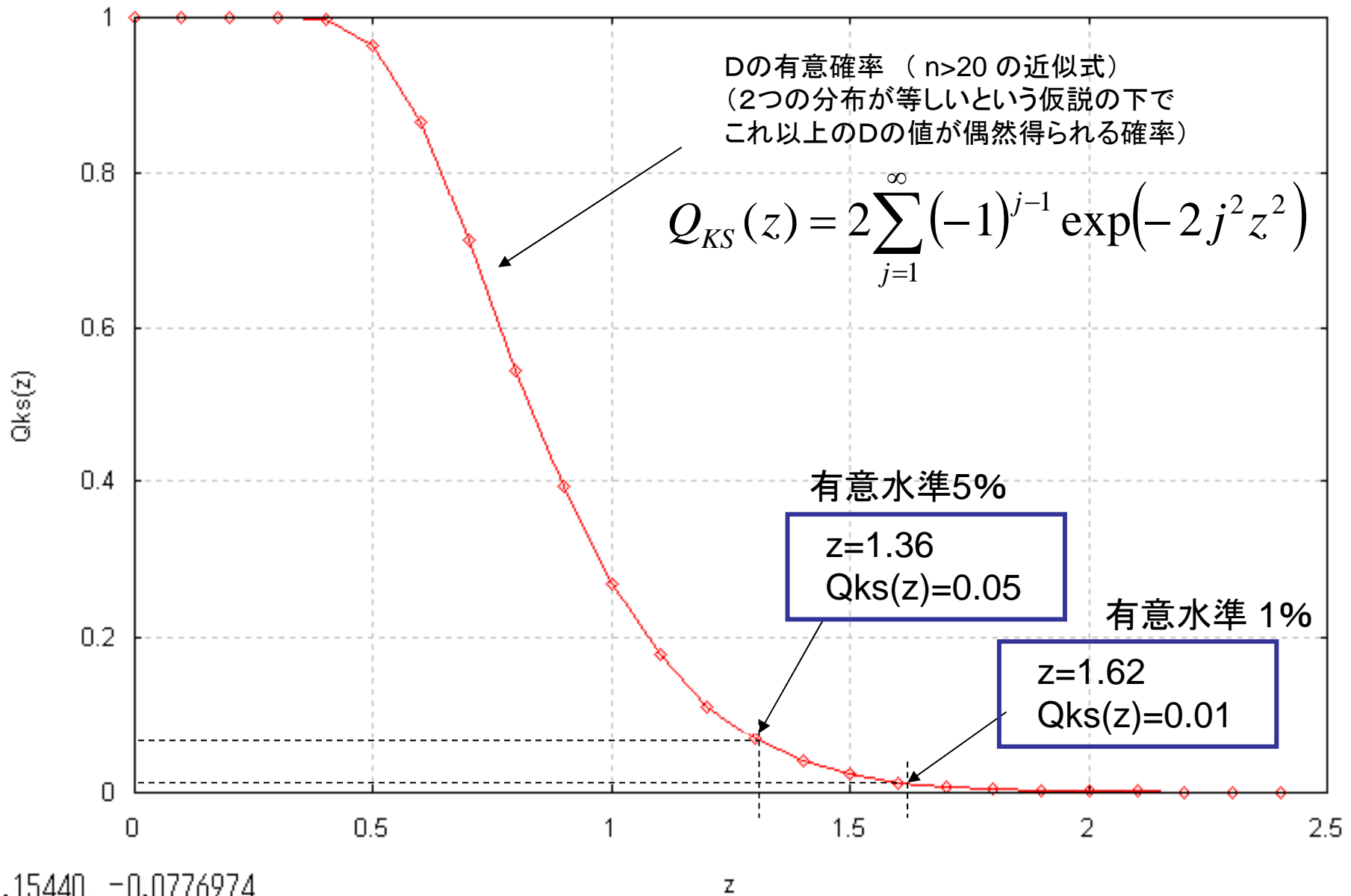


K-S検定の計算手順(1標本)

累積確率分布の差を計算



$$\Pr(D\sqrt{n} > z) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2j^2 z^2)$$



K-S検定の計算手順(1標本):まとめ

【1標本 K-S検定】

標本X

x_1, x_2, \dots, x_n

標本Xが1独立変数の確率密度関数 $f(x)$ と一致していると言えるか？

(ただし $n > 20$)

- (1) 帰無仮説を「標本Xが確率密度関数 $f(x)$ から発生」とする。
- (2) 標本Xの累積確率分布 $S_n(x)$ と確率密度関数 $f(x)$ の累積確率分布 $F(x) = \int_{-\infty}^x f(y)dy$ を求める。
- (3) 上記の2つの累積確率分布の差の絶対値の最大値であるKS統計量 $D = \max_{-\infty < x < \infty} |S_n(x) - F(x)|$ を求める。
- (4) 標本個数 n と上記のKS統計量 D を用いて の値を計算
- (5) 有意水準5%の場合 の値が1.36 以上だったら帰無仮説を棄却し、「一致しない」と結論
有意水準1%の場合 の値が 1.62 以上だったら帰無仮説を棄却し、「一致しない」と結論

K-S検定の計算手順(1標本):まとめ

【1標本 K-S検定】

標本X

$$x_1, x_2, \dots, x_n$$

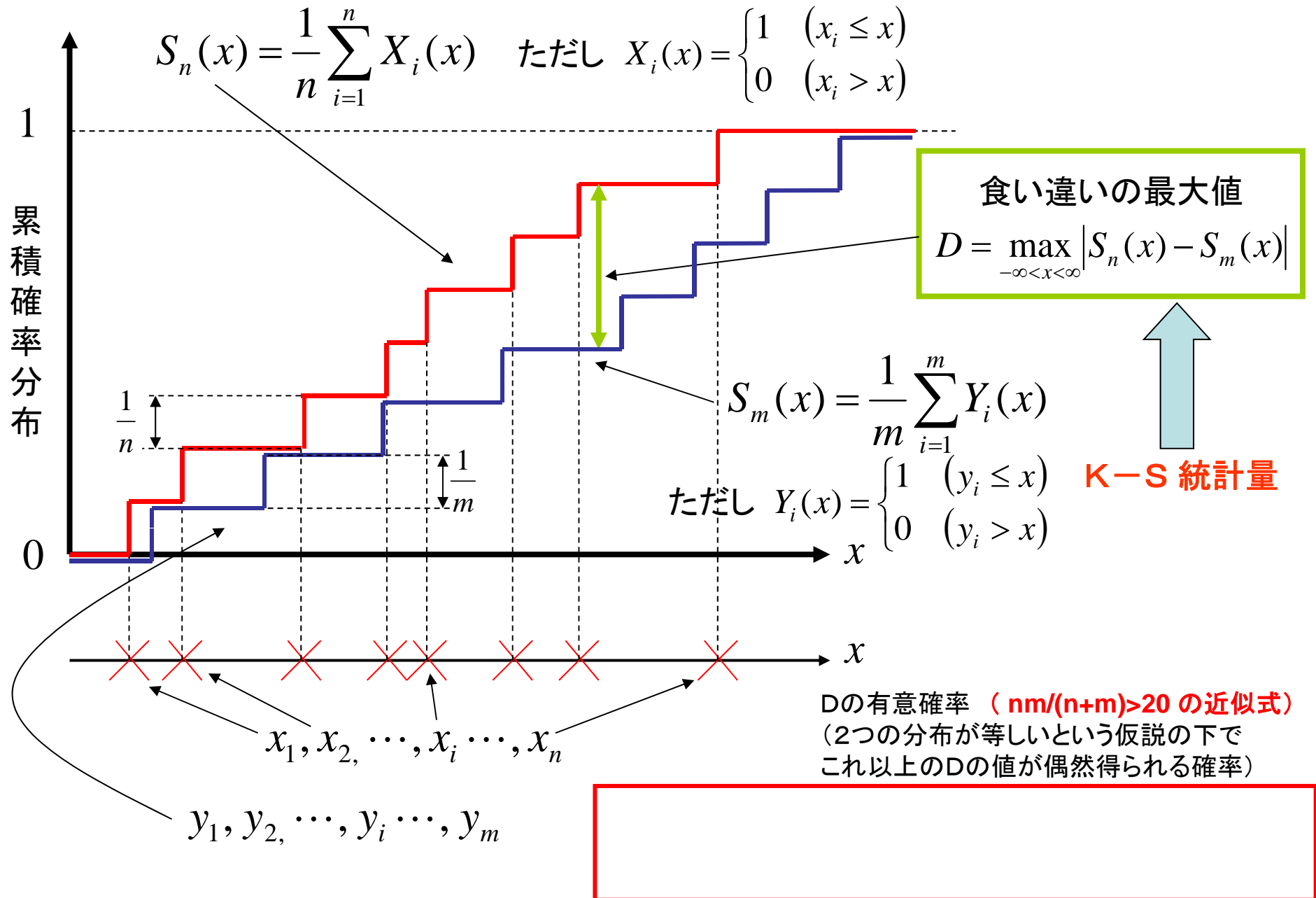
標本Xが1独立変数の確率密度関数 $f(x)$ と一致していると言えるか？

(ただし $n > 20$)

- (1) 帰無仮説を「標本Xが確率密度関数 $f(x)$ から発生」とする。
- (2) 標本Xの累積確率分布 $S_n(x)$ と確率密度関数 $f(x)$ の累積確率分布 $F(x) = \int_{-\infty}^x f(y)dy$ を求める。
- (3) 上記の2つの累積確率分布の差の絶対値の最大値であるKS統計量 $D = \max_{-\infty < x < \infty} |S_n(x) - F(x)|$ を求める。
- (4) 標本個数 n と上記のKS統計量 D を用いて $D\sqrt{n}$ の値を計算
- (5) 有意水準5%の場合、 $D\sqrt{n}$ の値が1.36 以上だったら帰無仮説を棄却し、「一致しない」と結論
有意水準1%の場合、 $D\sqrt{n}$ の値が 1.62 以上だったら帰無仮説を棄却し、「一致しない」と結論

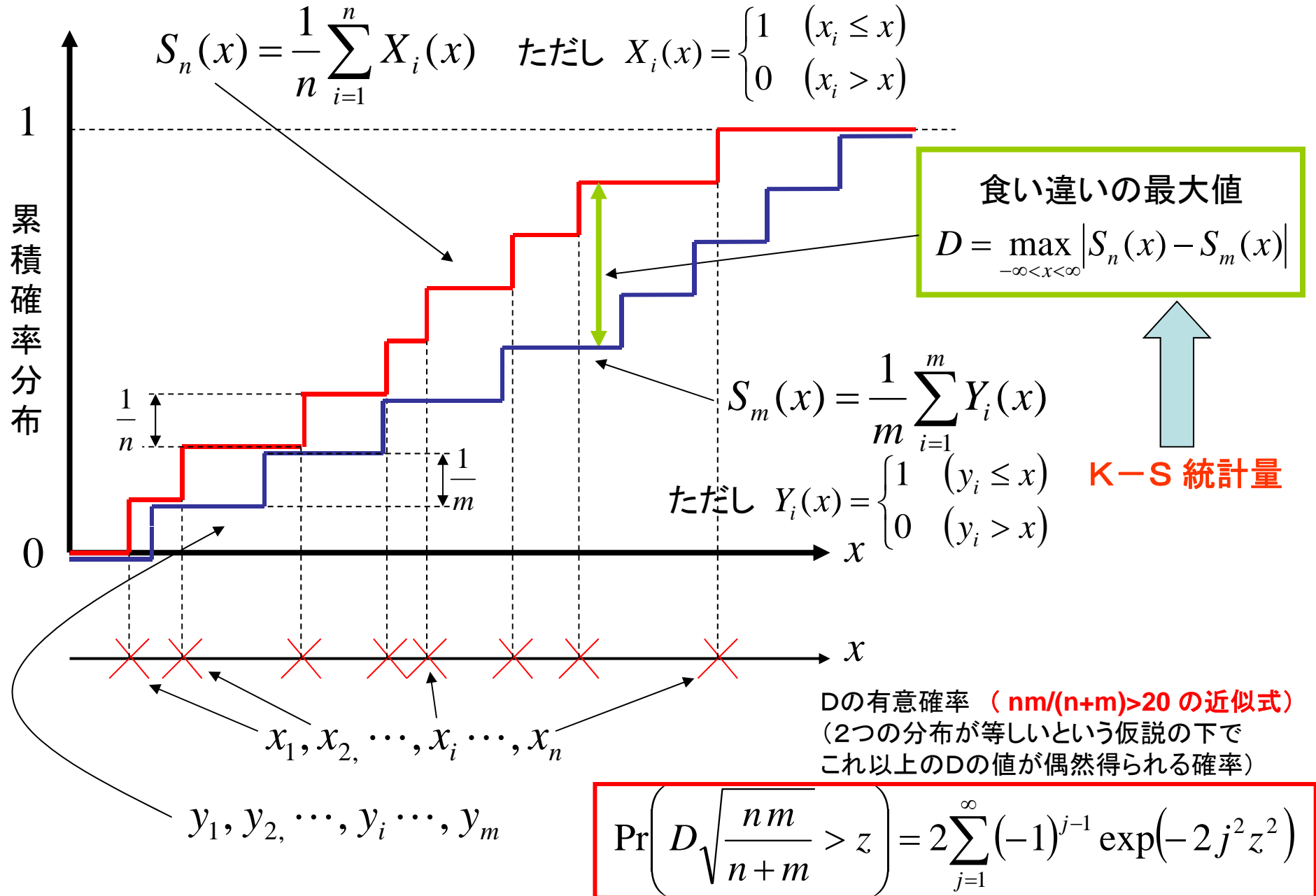
K-S検定の計算手順(2標本)

累積確率分布の差を計算



K-S検定の計算手順(2標本)

累積確率分布の差を計算



K-S検定の計算手順(2標本):まとめ

標本X

x_1, x_2, \dots, x_n

標本Y

y_1, y_2, \dots, y_m

標本Xと標本Yが同一の母集団の確率密度関数より生じていると言えるか？

(ただし $nm/(n+m) > 20$)

(1) 帰無仮説を「標本XとYが同一の母集団から発生」とする。

(2) 標本Xの累積確率分布 $S_n(x)$ と、標本Yの累積確率分布 $S_m(x)$ を求める。

(3) 上記の2つの累積確率分布の差の絶対値の最大値であるKS統計量

$$D = \max_{-\infty < x < \infty} |S_n(x) - S_m(x)| \text{ を求める。}$$

(4) 標本個数 n および m と上記のKS統計量 D を用いて の値を計算

(5) 有意水準5%の場合、 の値が1.36 以上だったら帰無仮説を棄却し、「一致しない」と結論

有意水準1%の場合、 の値が 1.62 以上だったら帰無仮説を棄却し、「一致しない」と結論

K-S検定の計算手順(2標本):まとめ

標本X

$$x_1, x_2, \dots, x_n$$

標本Y

$$y_1, y_2, \dots, y_m$$

標本Xと標本Yが同一の母集団の確率密度関数より生じていると言えるか？

(ただし $nm/(n+m) > 20$)

(1) 帰無仮説を「標本XとYが同一の母集団から発生」とする。

(2) 標本Xの累積確率分布 $S_n(x)$ と、標本Yの累積確率分布 $S_m(x)$ を求める。

(3) 上記の2つの累積確率分布の差の絶対値の最大値であるKS統計量

$$D = \max_{-\infty < x < \infty} |S_n(x) - S_m(x)| \text{ を求める。}$$

(4) 標本個数 n および m と上記のKS統計量 D を用いて $D \sqrt{\frac{nm}{n+m}}$ の値を計算

(5) 有意水準5%の場合、 $D \sqrt{\frac{nm}{n+m}}$ の値が1.36 以上だったら帰無仮説を棄却し、「一致しない」と結論

有意水準1%の場合、 $D \sqrt{\frac{nm}{n+m}}$ の値が 1.62 以上だったら帰無仮説を棄却し、「一致しない」と結論

【演習問題】 2016.07.26

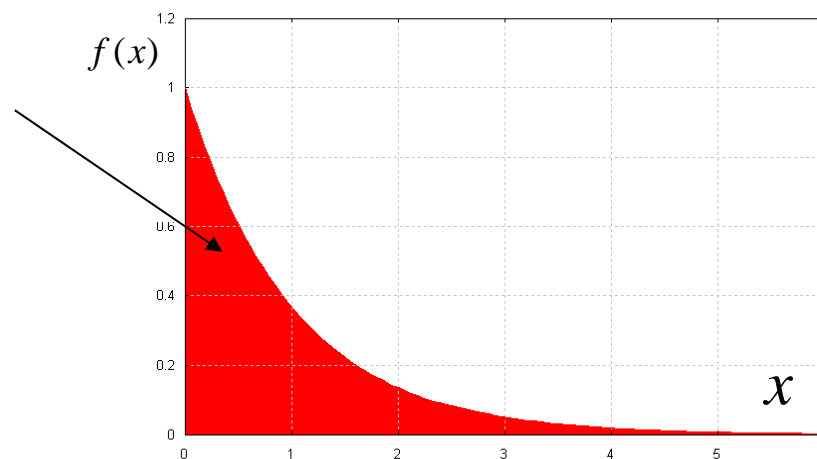
学籍番号

氏名

指数分布関数

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & (x \geq 0) \\ 0 & (x < 0) \end{cases}$$

の累積確率分布関数 $F(x)$ を求めよ。



【演習問題】

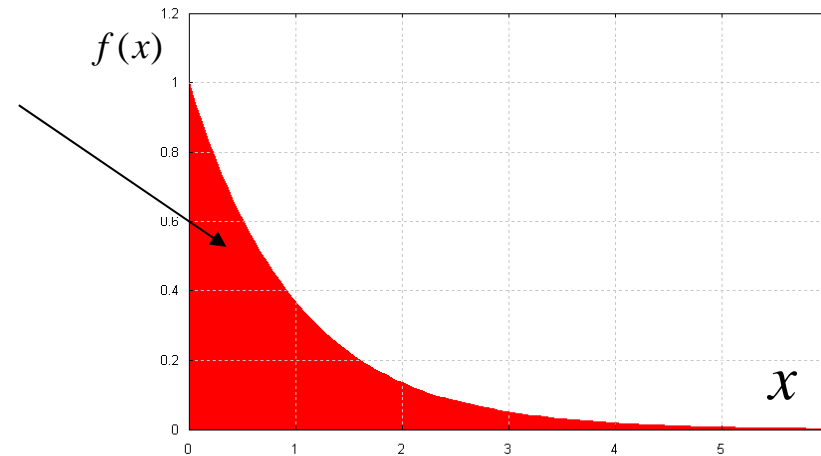
学籍番号

氏名

指数分布関数

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & (x \geq 0) \\ 0 & (x < 0) \end{cases}$$

の累積確率分布関数 $F(x)$ を求めよ。



$x \geq 0$ の場合

$$\begin{aligned} F(x) &= \int_{-\infty}^x f(y) dy = \int_0^x \lambda e^{-\lambda y} dy \\ &= \left[-e^{-\lambda y} \right]_0^x = 1 - e^{-\lambda x} \end{aligned}$$

$x < 0$ の場合

$$F(x) = 0$$