

2020/4/13 北大SPH・統計解析の基礎③④

推定と検定



北海道大学 医学統計学
横田 勲

1

今回の内容

2

- ▶ 統計的仮説検定
 - ▶ カイ二乗検定
- ▶ 区間推定
- ▶ 連続量データの比較
 - ▶ 対応のない場合；t検定とWilcoxon順位和検定
 - ▶ Student型と Welch型
 - ▶ 対応のある場合；対応のある t 検定と Wilcoxon符号付き順位検定、McNamer検定

2

29連勝

3

- ▶ 藤井聡太四段が29連勝した(2017/6/26)

3

炎上したが・・・

5

- ▶ 統計的に話を整理する
- 1. 対戦相手と実力が互角という仮定
 - ▶ 勝利する確率が0.5
- 2. 対戦結果は29勝0敗
- 3. 1.の仮定の下で、29勝0敗が起こることはどのくらいまれなことかを計算
 - ▶ 確率ではなくp値と言えば炎上しなかった！？

5

統計的仮説検定

6

- ▶ 科学的方法である背理法を導入
- 1. 「対戦相手との実力に差はない」と仮定
 - ▶ 帰無仮説という
 - ▶ 群間比較の場面では「比較群間に差がない」
- 2. 仮定が偽であることをいう
 - ▶ P値が事前に定めた有意水準より低い
- 3. 「実力に差がある」と結論づけ
 - ▶ 対立仮説を受容

6

帰無仮説と対立仮説

7

- ▶ 帰無仮説 H_0 (否定したい仮説)
 - ▶ 勝率は〇〇である
 - ▶ 比較群間で効果の大きさに違いがない
 - ▶ 曝露とアウトカムは無関係だ ...etc.
- ▶ 対立仮説 H_1
 - ▶ 帰無仮説と逆の内容
- ▶ 観察データが帰無仮説に反するかに注目
 - ▶ 対立仮説に反するかには注目しない

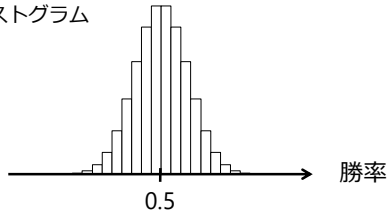
7

帰無仮説が正しい場合

8

- ▶ 同様の研究を繰り返したならば、観察される勝率は、0.5を中心として、左右対称に分布

「29試合行う」を繰り返したところ
得られた勝率のヒストグラム

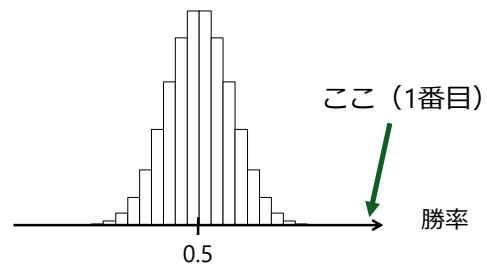


8

観察結果はどこにある？

9

- ▶ 勝率100% (29勝0敗) は、ヒストグラムのどこに位置するのか？

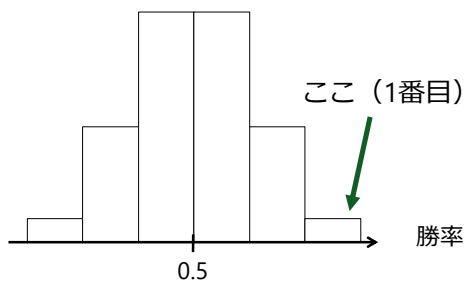


9

たとえば5連勝だったら？

10

- ▶ まだありえそう



10

同じ「1番目」にまれな出来事

11

- ▶ 試合数に応じて、「まれ」っぷりが違う
- ▶ 帰無仮説の下で、ありえるパターン数を分母において考える
 - ▶ 勝率が0.5だとして、5試合行うと、対戦結果は $2^5=32$ 通り考えられ、そのうち1番まれな結果が得られた
 - ▶ 勝率が0.5だとして、29試合行うと、対戦結果は $2^{29}=536870912$ 通り考えられ、そのうち1番まれな結果が得られた

11

どちらがよまれ？

12

- ▶ 勝率が0.5である下で、
 - ▶ 29勝1敗
 - ▶ 勝率 97%
 - ▶ 全パターンのうち、2番目にまれ
 - ▶ 5勝0敗
 - ▶ 勝率 100%
 - ▶ 全パターンのうち、1番目にまれ

12

p値

13

- ▶ 帰無仮説の下で、考えられうる全ての結果パターンに対し、観察結果が得られた順位
 - ▶ 順位を0から1の間で基準化
- ▶ 事前に定めた有意水準より低ければ、帰無仮説は正しくないと思決定する
 - ▶ 有意水準は、通常片側2.5%

13

正規近似を用いた検定

14

- ▶ 以下の検定統計量 Z が標準正規分布に従うことを利用
 - ▶ 平均0、分散 1^2 の正規分布

$$Z = \frac{\hat{p} - p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim N(0,1^2)$$

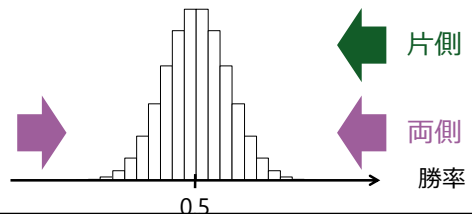
ただし、 \hat{p} は観察された勝率
 p_0 は帰無仮説の下での勝率
 n は試合数（サンプルサイズ）

14

片側p値？両側p値？

15

- ▶ 勝率が高い方からみるか、勝率が低い方からみるものも加えるか
 - ▶ 通常、両側p値は片側p値の2倍
 - ▶ 両側p値で報告することが多い



15

片側検定と両側検定

16

- ▶ 言いたい対立仮説に応じて選ぶべき
- ▶ 片側検定
 - ▶ 帰無仮説：勝率は0.5である
 - ▶ 対立仮説：勝率は0.5より高い
- ▶ 両側検定
 - ▶ 帰無仮説：勝率は0.5である
 - ▶ 対立仮説：勝率は0.5より高い、もしくは低い

16

脚気論争

17

- ▶ 長期航海において脚気患者が続出
 - ▶ 食事内容を変更し、同一航路で訓練航海

食事	脚気の発生		合計
	あり	なし	
洋食	14 (4.2%)	319	333
米食	169 (44.9%)	207	376

- ▶ 洋食にすれば脚気は減る？

17

割合の95%信頼区間を計算

18

- ▶ 洋食
 - ▶ $\frac{14}{333} \pm 1.96 \sqrt{\frac{\frac{14}{333} \times \frac{333-14}{333}}{333}} \approx (0.020, 0.064)$
- ▶ 米食
 - ▶ $\frac{169}{376} \pm 1.96 \sqrt{\frac{\frac{169}{376} \times \frac{376-169}{376}}{376}} \approx (0.399, 0.500)$
- ▶ どうやら差はありそう

18

統計的仮説検定を導入

19

- ▶ 帰無仮説 H_0 ：
食事によって脚気発生割合は変わらない
- ▶ 対立仮説 H_1 ：
洋食は米食より脚気発生割合が低い
 - ▶ 片側検定を用いる
- ▶ 脚気発生割合の群間差に注目

19

仮説をパラメータ化

20

- ▶ 洋食での脚気発生割合 p_1
- ▶ 米食での脚気発生割合 p_2
- ▶ 帰無仮説 $H_0: p_1 = p_2$
 - ▶ 差 ($\Delta = p_1 - p_2$) で書き直せば、 $H_0: \Delta = 0$
- ▶ 対立仮説 $H_1: p_1 < p_2$
 - ▶ 差 ($\Delta = p_1 - p_2$) で書き直せば、 $H_1: \Delta < 0$

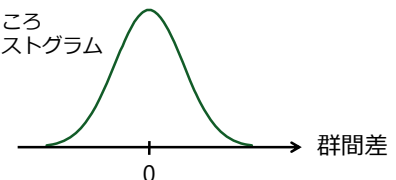
20

帰無仮説が正しい下で

21

- ▶ 同様の研究を繰り返したならば、観察される群間差は、ゼロを中心として、左右対称に分布
- ▶ サンプルサイズ大きくなるにつれ、正規分布に近づく

研究を繰り返したところ
得られた群間差のヒストグラム



21

群間差が従う正規分布

22

- ▶ 帰無仮説が正しい下で、群間差について
 - ▶ 平均は0
 - ▶ 分散は $\frac{t}{N} \cdot \frac{N-t}{N} \cdot \left(\frac{1}{n} + \frac{1}{m}\right)$

比較群	疾病発生		合計
	あり	なし	
試験群	a	b	n
対照群	c	d	m
合計	t	$N - t$	N

22

割合の差の検定

23

- ▶ 帰無仮説の下で、以下の検定統計量 Z が標準正規分布に従う

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n} + \frac{1}{m}\right)}} \sim N(0, 1^2)$$

$$\text{ただし、} \hat{p}_1 = \frac{a}{n}, \hat{p}_2 = \frac{b}{m}, \hat{p} = \frac{t}{N}$$

23

カイ二乗検定

24

- ▶ 検定統計量 Z を2乗して整理すると

$$Z^2 = \chi^2 = \frac{N(ad - bc)^2}{nmt(N - t)} \sim \chi_1^2$$

- ▶ 自由度1のカイ二乗分布
- ▶ 割合の差の検定と全く同じ
- ▶ 2×2分割表以外の分割表にも拡張可能

24

脚気論争の例

25

- ▶ 割合の差の検定

$$\text{▶ } Z = \frac{\left(\frac{14}{333} - \frac{169}{376}\right) - 0}{\sqrt{\frac{183}{709} \cdot \frac{526}{709} \left(\frac{1}{333} + \frac{1}{376}\right)}} = -375.74 < -1.96$$

- ▶ 片側2.5%有意水準で有意差あり
 - ▶ 「洋食は米食より脚気発生割合が低い」

25

この例では？

26

比較群	疾病発生		合計
	あり	なし	
試験群	10 (50%)	10	20
対照群	4 (20%)	16	20
合計	14	26	40

- ▶ 割合の差の検定 : $Z = 1.66 < 1.96$
 - ▶ 有意差なし
- ▶ 「群間に差がなかった」とはいえない

26

仮説検定での2つのエラー

27

- ▶ α エラー (type-I エラー, 第一種の過誤)
 - ▶ 本当は差がないのに有意差ありという誤り
 - ▶ 消費者リスク
 - ▶ (Awatenbo)あわてんぼうさんの間違い
- ▶ β エラー (type-II エラー, 第二種の過誤)
 - ▶ 本当は差があるのに有意差を出せない誤り
 - ▶ 生産者リスク
 - ▶ (Bonyari)ぼんやり者の間違い

27

 α エラーと β エラー

28

検定結果	母集団 (真)	
	差がない	差がある
有意差なし	正しい	第2種の過誤 β エラー
有意差あり	第1種の過誤 α エラー	正しい

有意水準で大きさを決める

サンプルサイズを増やすことで軽減

28

有意差なし \neq 差がない

29

- ▶ 本当に差がなかった
- ▶ β エラーによって、本当は差があるのに有意差なしかも
 - ▶ どちらが正しいかはデータから判定不能
 - ▶ 背理法の理屈からも、帰無仮説が正しいとはいえない

29

統計的仮説検定のポイント

30

- ▶ 二者択一の意味決定
 - ▶ 差がある or 差があるか分からない
- ▶ 有意差がある場合、差があることを主張する強力な方法
 - ▶ 背景には背理法の理屈
- ▶ 有意差がない場合、解釈が難しい
 - ▶ 差がない、とってはならない

30

効果や影響の大きさ

31

- ▶ 検定では、差があることしか分からない
- ▶ 脚気割合の場合は、信頼区間を求めた

$$\frac{\hat{p} - \mu}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} = \pm 1.96$$

- ▶ 群間差でも信頼区間を求めたい
 - ▶ 推定値の精度を議論できる

31

脚気発生割合の場合

32

信頼区間

p値の計算

▶ $\frac{\hat{p}-\mu}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} = \pm 1.96$
を μ について解いた

▶ $Z = \frac{\hat{p}-p_0}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim N(0,1^2)$
▶ 事前に定めた p_0 が正しいか

信頼区間は、検定で有意にならない範囲

32

群間差の95%信頼区間

33

信頼係数という

- ▶ 割合の差の検定をベースに、以下の式の μ を解く

$$\frac{(\hat{p}_1 - \hat{p}_2) - \mu}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n} + \frac{1}{m}\right)}} = \pm 1.96$$

33

脚気論争の例

34

$$(\hat{p}_1 - \hat{p}_2) \pm 1.96 \sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n} + \frac{1}{m}\right)}$$

$$\left(\frac{14}{333} - \frac{169}{376}\right) \pm 1.96 \sqrt{\frac{183}{709} \cdot \frac{526}{709} \left(\frac{1}{333} + \frac{1}{376}\right)}$$

▶ (-0.343, -0.472)

34

脚気論争例の解釈

35

- ▶ (-0.343, -0.472)
 - ▶ 洋食にすると脚気発生割合が30%下がる！
 - ▶ ウソ
 - ▶ 洋食にすると脚気発生割合が40%下がる！
 - ▶ 否定できない
 - ▶ 洋食にすると脚気発生割合が下がらない！
 - ▶ ウソ
- ▶ 差の信頼区間が0を含まない
 - ⇔ 検定で有意差あり

35

信頼区間の解釈

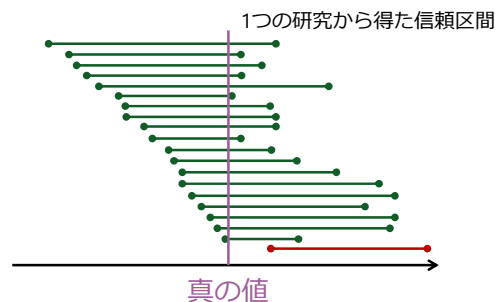
36

- ▶ 厳密な解釈
 - ▶ 同じような研究を繰り返した場合、研究ごとに信頼区間を推定
 - ▶ それら信頼区間100通りあたり、95通りは真の値を含む
- ▶ 実際の解釈
 - ▶ 今回の研究データから効果の大きさとしてありうるであろう範囲

36

仮想的に研究を繰り返し

37

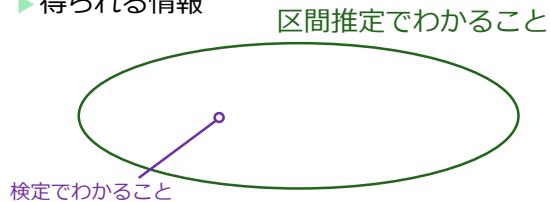


37

区間推定のメリット

38

- ▶ 単に差があることをいうだけでなく、どの程度の効果であるのかを議論できる
 - ▶ 研究の精度もコミの議論
- ▶ 得られる情報



38

Scientists rise up against statistical significance

39

Amrhein V, Greenland S, McShane B. Nature. 2019;567:305-307.

39

P値の誤用入門①

40

- ▶ 死亡群と生存群では年齢が異なった
- ▶ $P > 0.05$ だったから両群は同じ
- ▶ 年齢の違いでは $P = 0.04$ 、性別の違いでは $P = 0.02$ ゆえ、性別のほうが与える影響が大きい

40

P値の誤用入門②

41

- ▶ 有意差がつかないと論文にならないから Fisher 検定をカイ二乗検定に変えた
- ▶ 有意差があったから、グループ間には間違いなく差があるんだ

41

P値廃止運動

42

- ▶ あまりに科学界で検定を濫用している
- ▶ いっそ、P値を示すことを禁止しては
- ▶ 医学研究では臨床試験におけるアウトカム比較でのみ検定を使用すべき

42

結果変数の型による分類

43

目的	連続尺度	分類尺度	時間イベント尺度
分布の記述	ヒストグラム、箱ヒゲ図、散布図	ヒストグラム、分割表	生存曲線 (Kaplan-Meier法)
要約統計量	平均、分散、中央値、パーセント点、相関係数	頻度、一致度、相関係数	x年生存確率、中央生存期間
検定 (単純)	t検定、分散分析、Wilcoxon検定	χ^2 検定、Fisher正確検定	ログランク検定
検定 (層別)	共分散分析	Mantel-Haenszel検定	層別ログランク検定
回帰モデル	重回帰分析	ロジスティック回帰分析	Cox回帰分析

43

血圧を下げる飲料！？

44

- ▶ 収縮期血圧が140mmHg以上 (平均は144mmHg)であった 大学生男性20名に飲料を摂取させた
- ▶ 30分後に血圧を測定したところ、 平均収縮期血圧は125mmHgになった
- ▶ この飲料は収縮期血圧を下げる！

あなたは血圧の高い人にこの飲料をすすめますか？

44

のんだ、なおった、きいた？

45

- ▶ 同じ条件でこの飲料をのまなかったら？
 - ▶ 一息ついてリラックスしただけ？
- ▶ そもそも血圧の高い大学生男性って？
 - ▶ 直前に運動をしていた？？
- ▶ 平均は120mmHgだが、そのバラツキは？
- ▶ 血圧の測定器はどのようなもの？
 - ▶ 同じ機種を用いたか？
 - ▶ 人が聴診して測定したか？

45

コントロール(対照)の設定

46

- ▶ 飲料を摂取しないだけで、 その他は同じ条件にしたグループをおく
- ▶ 血圧の変化量を 飲料摂取群と非摂取群で比較

46

平均への回帰 regression to the mean

47

- ▶ ある基準以上という人だけ選んでみると、 次の測定までに何もなされていなくても、 ランダムなバラツキがゆえに 平均値に結果が近づく現象
 - ▶ たまたま血圧の高かった人が対象者になった ため、飲料とは関係なく30分後の血圧が 下がっただけ？

47

対象集団

48

- ▶ 第1相試験では、健康な男性大学生が対象
 - ▶ それはヒトへの安全性をみるためでしょう
- ▶ この飲料を売るターゲットは中高年
 - ▶ 同じ高血圧でも大学生のそれとはわけが違う
 - ▶ 大学生と同じように効くとは限らない

48

測定の信頼性

49

- ▶ 血圧計の較正は難しい
 - ▶ 全く同じものを測定しても、 血圧計によって値が異なる
- ▶ ましてや人が目で見ていたら・・・
- ▶ そもそも血圧は個人内変動が激しい
- ▶ 測定回数は全員1回だけ？
 - ▶ 2回測定した平均をとった人もいた！？
- ▶ 1mmHg単位での信頼性はない！

49

コントロール群をおいた研究

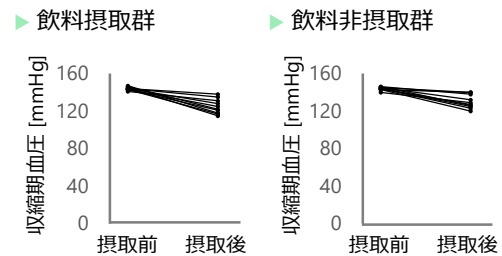
50

- ▶ 対象
 - ▶ 高血圧症と診断された、50歳以上の男女
- ▶ 方法：ランダム化
 - ▶ 飲料摂取群10名
 - ▶ 飲料非摂取群10名
 - ▶ 有効成分を除いた飲料（プラセボ）を摂取
- ▶ 測定
 - ▶ 同じ施設・測定器・実施者により
安静時血圧を摂取前と1カ月の継続摂取後に
各々3回測定し、その中央値を測定値に採用

50

血圧の変化

51



研究仮説：
飲料摂取群は非摂取群より収縮期血圧が下がるか？

51

エンドポイント（評価項目）

52

- ▶ そのままの値
 - ▶ 1カ月後の測定値
- ▶ 絶対的変化量
 - ▶ (1カ月後の測定値) - (摂取前の測定値)
- ▶ 相対的変化量
 - ▶ $\frac{(1カ月後の測定値) - (摂取前の測定値)}{(摂取前の測定値)}$

52

絶対的変化量を採用

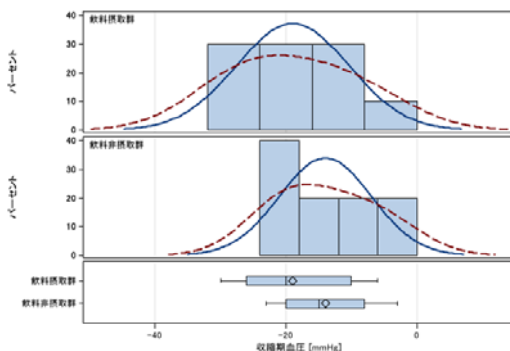
53

- ▶ 飲料摂取群（10名）
 - ▶ -6, -9, -10, -16, -18, -22, -24, -26, -29, -30
 - ▶ 平均[SD]：-19.0 [8.6]
 - ▶ 中央値[四分位範囲]：-20 [-10, -26]
- ▶ 飲料非摂取群（10名）
 - ▶ -3, -5, -8, -11, -13, -17, -19, -20, -21, -23
 - ▶ 平均[SD]：-14.0 [7.1]
 - ▶ 中央値[四分位範囲]：-15 [-8, -20]

53

絶対的変化量の分布

54



54

分布の違い？

55

- ▶ モーメントを用いた比較
 - ▶ 平均：位置
 - ▶ 分散：分布の広がり
 - ▶ 歪度：分布の歪み
 - ▶ 尖度：分布のすその重さ
- ▶ 中央値や分布関数を用いた比較
 - ▶ 中央値：位置
 - ▶ 分布関数：分布の広がり

55

位置の比較

56

- ▶ t 検定：平均値の比較
- ▶ Wilcoxonの順位和検定：中央値の比較
- ▶ 並べ替え検定：どちらも可能
- ▶ 分布の広がり、歪み、すその重さは、比較群間で同様と仮定
 - ▶ ランダム化によって介入前の分布は同じに
 - ▶ 介入によって分布の形状が変わることは考えづらい

56

ランダム化

57

- ▶ 20名を同じ確からしさに
撮取群か非撮取群に割付
 - ▶ 合計10名ずつになるように
- ▶ ランダム化した結果、
対象者の割り付けパターン数は・・・？
 - ▶ ${}_{20}C_{10}=184,756$ 通り

57

並べ替え検定

58

- ▶ 帰無仮説：どちらの群でも変化量が同じ
 - ▶ 全割付パターンを計算すれば、
帰無仮説の下、観察されうる結果がすべて判明
 - ▶ しかも逆のパターンが必ずあるので、
「差がない」場合の結果の分布が分かる
 - ▶ AABABBというパターンがあれば、
BBABAAというパターンも必ずある
- ▶ 例えば、変化量平均値の群間差を
各割付パターンで求めてみよう

58

割付を並べ替え

59

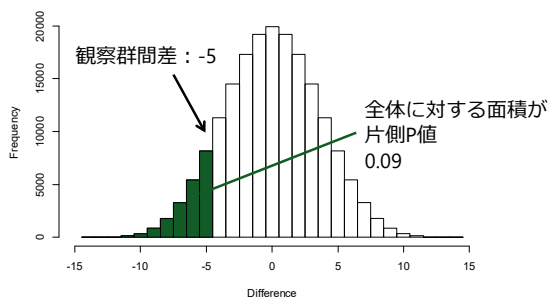
ID	変化量	実際の割付	パターン①	パターン②
1	-6	撮取群	撮取群	撮取群
2	-9	撮取群	非撮取群	非撮取群
⋮	⋮	⋮	⋮	⋮
10	-30	撮取群	撮取群	非撮取群
11	-3	非撮取群	非撮取群	非撮取群
12	-5	非撮取群	撮取群	撮取群
⋮	⋮	⋮	⋮	⋮
20	-23	非撮取群	非撮取群	撮取群

59

変化量平均値の群間差の分布

60

- ▶ 184,756通り計算



60

並べ替え検定に基づくP値

61

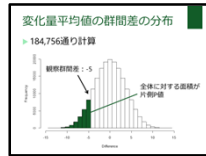
- ▶ ありうるパターンの何番目にいるか
 - ▶ 正確なP値
 - ▶ 二値(成功/失敗)なら手計算で可能(Fisher検定)
 - ▶ 変化量 (アウトカム) が連続量だと、
通り数が増えすぎて大変
 - ▶ 10例 vs 10例では18万通り以上計算せねば
 - ▶ 100例 vs 100例では約 10^{29} 通り
- ▶ 変化量の違いをみる指標には
平均値のほか、中央値や分位点も可能
 - ▶ 検討したい仮説に応じて選べる

61

平均値の違いに注目

62

- ▶ 帰無仮説 $H_0: \mu_X = \mu_Y$
 - ▶ 摂取群平均値 μ_X と非摂取群の平均値 μ_Y が同じ
- ▶ 対立仮説 $H_1: \mu_X < \mu_Y$
 - ▶ 摂取群平均値 μ_X のほうが、より減少している
- ▶ 平均値の群間差の分布
さえ分かれば・・・



62

t分布

63

- ▶ 正規分布より少しずそが重い分布
 - ▶ 自由度が $+\infty$ の場合は正規分布に一致
- ▶ データが独立に同一の正規分布に従う場合、平均の差に従う正確な分布
 - ▶ データが正規分布に従っている必要はない
 - ▶ 別にデータが正規分布に従っていなくても、ランダム化試験のように分散が同じと仮定できるなら検定は妥当(valid)

63

Studentのt検定

64

- ▶ t が自由度 $n_X + n_Y - 2$ の t 分布に従う
 - ▶ n_X, n_Y : 各群の人数
 - ▶ \bar{x}, \bar{y} : 各群の平均値

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{(1/n_X + 1/n_Y)s^2}}$$

$$s^2 = \frac{\sum_i^{n_X} (x_i - \bar{x})^2 + \sum_i^{n_Y} (y_i - \bar{y})^2}{(n_X - 1) + (n_Y - 1)}$$

64

自由度 Degrees of Freedom

65

- ▶ 偏差平方和(データのバラツキの表現方法)
 - ▶ n 個のデータ x_1, \dots, x_n のバラツキは、それぞれ平均値からの偏差 $(x_i - \bar{x})$ の二乗和

$$\sum_i^n (x_i - \bar{x})^2$$
- ▶ 自由度
 - ▶ 偏差平方和が何個の独立な二乗和から成るか
 - ▶ 何個の二乗和が特定されれば、偏差平方和が判明するか
 - ▶ 平均が分かっているため、 $n - 1$ 個が独立な個数

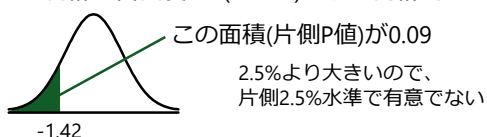
65

飲料の例 (t検定)

66

- ▶ 帰無仮説の下で統計量を計算

$$t = \frac{-19 - (-14)}{\sqrt{\left(\frac{1}{10} + \frac{1}{10}\right) \frac{664 + 448}{(10 - 1) + (10 - 1)}}} = -1.42$$
- ▶ 自由度18の t 分布から P 値を計算
 - ▶ 各群の自由度は $9 (= 10 - 1)$ より、両群では18



66

95%信頼区間

67

- ▶ 効果の大きさを含めた議論
- ▶ 真の平均値の差 δ

$$t = \frac{(\bar{x} - \bar{y}) - \delta}{\sqrt{(1/n_X + 1/n_Y)s^2}}$$
- ▶ 両側5%水準で有意とならない δ の範囲

$$(\bar{x} - \bar{y}) \pm 2.10 \sqrt{(1/n_X + 1/n_Y)s^2}$$

自由度18の t 分布における片側2.5%点

67

飲料の例（平均値の区間推定）

68

$$-19 - (-14) \pm 2.10 \sqrt{\left(\frac{1}{10} + \frac{1}{10}\right) \frac{664 + 448}{18}}$$

$$= (-12.4, 2.4)$$

- ▶ 95%信頼区間に0を含んでいる
 - ▶ すなわち、片側2.5%(両側5%)で有意でない

68

信頼区間に注目

69

- ▶ 平均の差は-5
 - ▶ その信頼区間は(-12.4, 2.4)
- ▶ -10mmHgくらい下がるならば臨床的な意義もあるかも
- ▶ 十分なサンプルサイズだった?

69

医学的有意差と統計的有意差

70

- ▶ 対象者数が非常に多い場合
 - ▶ わずかな治療効果(差)であっても統計的に有意
 - ▶ わずかな差が臨床的に重要?
- ▶ 対象者数が少ない場合
 - ▶ (非常に)大きな差でも統計的有意とはならない
 - ▶ 医学的に重要な差であれば、無視すべきものでなく、さらに検討すべき

70

(Studentの)t検定の仮定

71

- ▶ 分散が同じという仮定が必要
 - ▶ ランダム化試験では、リーズナブル
- ▶ 分散が多少違っていても、比較群間で人数が大体揃っていれば、検定はほぼ妥当(使ってよい)
 - ▶ α エラー率が、名目水準を上回らないこと
 - ▶ 観察研究では、めったに人数は揃わない

71

Welchのt検定

72

- ▶ 分散が同じという仮定が不要
 - ▶ 分散が異なる場合にも対応
 - ▶ s_x^2, s_y^2 : 各群で求めた不偏分散

$$t_w = \frac{\bar{x} - \bar{y}}{\sqrt{s_x^2/n_x + s_y^2/n_y}}$$

- ▶ Satterthwaiteの近似(式は省略)によって求めた自由度のt分布に従う

72

予備検定方式

73

- ▶ ものの本には、~~等分散性の検定を行い、有意差があればWelch型、なければStudent型~~とある
 - ▶ 予め検定を行い、次に行う検定方法を決める
- ▶ 研究デザインに応じて、検定方法は決めるべき

73

StudentかWelchか

74

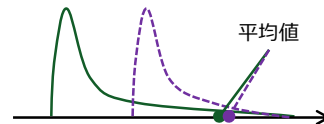
- ▶ ランダム化比較試験であればStudent型
 - ▶ 割付比が1:1の場合
 - ▶ 検出力(1-βエラー)はStudent型が高い
 - ▶ 本当は差がある場合に有意差ありという確率
 - ▶ ランダム化する人(サンプルサイズ)を減らせる
- ▶ 観察研究はWelch型
 - ▶ 検定はいつでもよく、信頼区間を表示することを最優先
 - ▶ 少しでもマシな信頼区間を出すために、Welch型に基づく信頼区間を示すべき

74

分布形によらない検定

75

- ▶ t検定は元のデータが正規分布に従うとき最も効率の良い(=有意になりやすい)検定
 - ▶ せめて左右対称な分布形であるとよい
- ▶ 例えば極端に歪んだ分布
 - ▶ 平均値の解釈がわかりづらい
 - ▶ t検定では効率が悪い(=有意差を出しづらい)



75

ノンパラメトリックな検定

76

- ▶ データの確率分布を仮定しない
 - ▶ 対義語: パラメトリックな検定
- ▶ データの分布に注目 $F_X(u) = F_Y(u + \Delta)$
 - ▶ 帰無仮説 $H_0: \Delta = 0$
 - ▶ 分布が重なっている
 - ▶ 対立仮説 $H_1: \Delta \neq 0, (\Delta > 0, \Delta < 0)$
 - ▶ 分布がΔだけずれている
- ▶ 解釈は中央値の比較とすればよい

76

Wilcoxonの順位和検定

77

- ▶ Mann-WhitneyのU検定 と等価
- ▶ 全データの順位を割当
- ▶ 群の順位和 U_X を計算
 - ▶ U_X の期待値 $n_X(n_X + n_Y + 1)/2$
 - ▶ U_X の分散 $n_X n_Y (n_X + n_Y + 1)/12$
- ▶ 以下の検定統計量を計算

$$\frac{U_X - n_X(n_X + n_Y + 1)/2}{\sqrt{n_X n_Y (n_X + n_Y + 1)/12}}$$

77

飲料の例 (Wilcoxonの順位和検定)

78

- ▶ 各データの順位を計算
 - ▶ 摂取群
 - ▶ -30, -29, -26, -24, -22, -18, -16, -10, -9, -6
 - ▶ 1, 2, 3, 4, 6, 10, 12, 15, 16, 18 (順位)
 - ▶ 非摂取群
 - ▶ -23, -21, -20, -19, -17, -13, -11, -8, -5, -3
 - ▶ 5, 7, 8, 9, 11, 13, 14, 17, 19, 20 (順位)
- ▶ 摂取群での順位和を計算
 - ▶ $U_X = 1 + 2 + \dots + 18 = 87$

78

飲料の例 (Wilcoxonの順位和検定)

79

- ▶ 検定統計量を計算

$$\frac{87 - 10(10 + 10 + 1)/2}{\sqrt{10 \cdot 10 (10 + 10 + 1)/12}} = -1.36$$
- ▶ 検定統計量が標準正規分布に従うことを利用し、P値を計算
 - ▶ $-1.36 > -1.96$ より片側2.5%水準で有意差なし

79

t検定かWilcoxon検定か①

80

- ▶ 統計解析の医学的な解釈
 - ▶ 多くの状況では平均値の比較が解釈容易
 - ▶ 平均値に差がでた or 中央値に差がでた
 - ▶ 平均値の信頼区間 or 中央値の信頼区間
 - ▶ 平均値を議論することが無意味な場合
 - ▶ 分布が大きく歪んでいる
 - ▶ 平均値は外れ値の影響を受けやすい
 - ▶ 一部の患者のみ特異的に大きな値をとる状況
 - ▶ 閾値を定め、閾値を以上/未満を議論すべき

80

t検定かWilcoxon検定か②

81

- ▶ ものの本には、
 - ▶ ~~サンプルサイズが少ない時に、Wilcoxon検定~~
- ▶ どちらも分布の位置の違いを評価
 - ▶ サンプルサイズが多いときには同等 (t検定のロバストネス)
 - ▶ サンプルサイズが極端に少ない時には?
 - ▶ 仮説検定の問題か? (データの提示で十分?)
 - ▶ Wilcoxon検定は検出力が低い

81

2群の比較といっても...

82

- ▶ 同一患者さんに、時期をずらして、2つの治療法を施す臨床試験
 - ▶ クロスオーバー試験
- ▶ 治療しなければ回復が見込めない下での治療前後での患者さんの状態を比較
 - ▶ 手術による効果の検討
- ▶ 術式による手術成績を比較するため、背景のリスク要因が似た対象者同士をマッチングして比較

82

中心角膜厚データ (一部)

83

- ▶ 術後24週と2年で違いがあるか

ID	術後24週	術後2年
1	511	532
2	525	538
3	540	546
4	640	710
5	509	529
6	505	525
7	626	550
8	489	503
9	595	543
10	539	523
11	561	572

Kinoshita S, et al. N Engl J Med. 2018; 995-1003.

83

対応のあるt検定

84

- ▶ i 番目の患者さんの結果 X_{Ai}, X_{Bi}
 - ▶ 差 $d_i = X_{Ai} - X_{Bi}$
- ▶ 帰無仮説 H_0 : 差の平均がゼロ
 - ▶ 差の平均値とその標準偏差 \bar{d}, s_d
 - ▶ 検定統計量 t が自由度 $n - 1$ の t 分布に従う

$$t = \frac{\bar{d}}{\sqrt{s_d^2/n}}$$

84

中心角膜厚の例

85

- ▶ 差の平均値とその標準偏差を計算
 - ▶ $\bar{d} = 2.82, s_d = 39.13$
- ▶ 検定統計量の計算
 - ▶ $t = \frac{2.82}{\sqrt{39.13^2/11}} = 0.24$
- ▶ 自由度 $10 (= 11 - 1)$ の t 分布から P 値を計算
 - ▶ 片側 P 値: 0.41

85

Wilcoxonの符号付き順位検定

86

- ▶ 対応のあるデータに対するノンパラ検定
- ▶ 帰無仮説 H_0 : 分布の中央がゼロ
 - ▶ ゼロを境に、正の値と負の値は同じバラツキ
- ▶ 差の絶対値 $|d_i|$ について、順位 R_i を考える
 - ▶ ゼロは順位をつけない

86

順位和の計算

87

- ▶ $T_+ = \sum R_i$ (正の値をとったデータ)
- ▶ $T_- = \sum R_i$ (負の値をとったデータ)
- ▶ 帰無仮説の下では、 $T_+ = T_-$
 - ▶ $E(T_+) = E(T_-) = n(n+1)/4$
 - ▶ $V(T_+) = V(T_-) = n(n+1)(2n+1)/24$
- ▶ 検定統計量が標準正規分布に従う

$$\frac{T_+ - n(n+1)/4}{\sqrt{n(n+1)(2n+1)/24}}$$

87

順位和の計算

88

ID	術後24週	術後2年	差	絶対値の順位	正符号の順位
1	511	532	21	8	8
2	525	538	13	3	3
3	540	546	6	1	1
4	640	710	70	10	10
5	509	529	20	6.5	6.5
6	505	525	20	6.5	6.5
7	626	550	-76	11	
8	489	503	14	4	4
9	595	543	-52	9	
10	539	523	-16	5	
11	561	572	11	2	2

- ▶ 検定統計量は、 $\frac{(8+3+1+10+6.5+6.5+4+2)-11 \cdot (11+1)/4}{\sqrt{11 \cdot (11+1) \cdot (2 \cdot 11+1)/24}} = 0.71$
- ▶ 片側P値=0.24

88

対応のない／ある群間比較

89

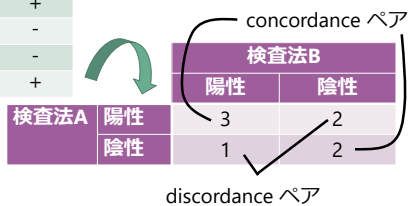
変数の種類	対応のないデータ	対応のあるデータ
連続量	<ul style="list-style-type: none"> • (2標本) t検定 • Wilcoxon順位和検定 	<ul style="list-style-type: none"> • (1標本) t検定 • Wilcoxon符号付き順位検定
カテゴリカル	<ul style="list-style-type: none"> • χ^2検定 • Fisherの直接検定 	<ul style="list-style-type: none"> • McNamer検定

89

対応のある2値アウトカム

90

患者ID	検査法A	検査法B
1	+(陽性)	-(陰性)
2	+	+
3	+	+
4	-	-
5	-	+
6	+	-
7	-	-
8	+	+



90

McNamer検定

91

- ▶ Discordanceペアだけに注目

		検査法B	
		陽性	陰性
検査法A	陽性	a	b
	陰性	c	d

- ▶ 自由度1のカイ二乗検定を利用

$$\chi^2 = \frac{(b - c)^2}{b + c}$$

91

まとめ①

92

- ▶ 検定と推定
 - ▶ 検定は差があることをいうためには強力
 - ▶ 推定は効果の大きさを含めたより多くの議論
- ▶ カイ二乗検定
 - ▶ カテゴリカルデータの群間比較

92

まとめ②

93

- ▶ 並べ替え検定、t検定、Wilcoxon順位和検定
 - ▶ 平均の比較をするなら、たいていt検定
 - ▶ ランダム化比較試験ではStudent型
観察研究ではWelch型に基づく信頼区間
- ▶ 対応のあるデータに対する検定
 - ▶ 同一対象者、マッチされたペアのように、
背景情報を似せている（同じにする）場合

93