

第8回 F分布・分散分析その1 一元配置

A. F分布

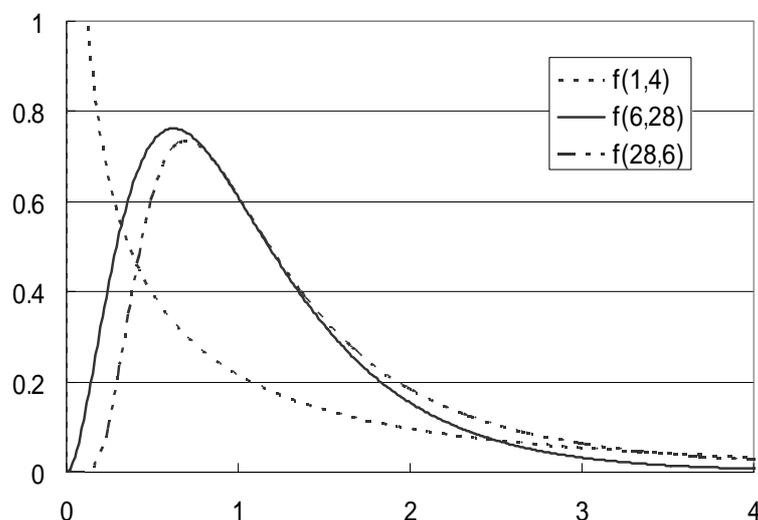
1. F分布

2つの母集団から得られた標本の分散比はF分布で検定する。分散分析という手法はこのF分布によって行うので、F検定は重要な手法である。カイ二乗分布と違い、2つの標本を比較するので、2つの標本それぞれの自由度からF分布が決まる。

F分布： $x_{11}, x_{12}, \dots, x_{1n_1}$ が互いに独立に $N(\mu_1, \sigma_1^2)$ に従い、さらにそれらと独立に

$x_{21}, x_{22}, \dots, x_{2n_2}$ が互いに独立に $N(\mu_2, \sigma_2^2)$ に従うとき、

$F = \frac{V_1/\sigma_1^2}{V_2/\sigma_2^2}$ は自由度 $(n_1 - 1, n_2 - 1)$ のF分布に従う。



F分布もカイ二乗分布と同じく、左右対称ではない。

2. F検定の手順

F分布に基づいて、2つの母集団の分散の比を推定することもできるが、あまり使わないので、ここでは述べない。

F検定は分散分析という手法の骨格をなす方法である。まずここでは基本的なF検定だけを練習しよう。

例：A, B 2種類の飼料を与えて一定期間飼育したハムスターの体重増加量を測定した結果、右のようなデータを得た(単位 g)。飼料による体重増加量のばらつきに差があるかを検定せよ。

飼料A	飼料B
73	71
74	68
70	73
77	71
72	71
72	72
69	69
73	72
72	68
75	

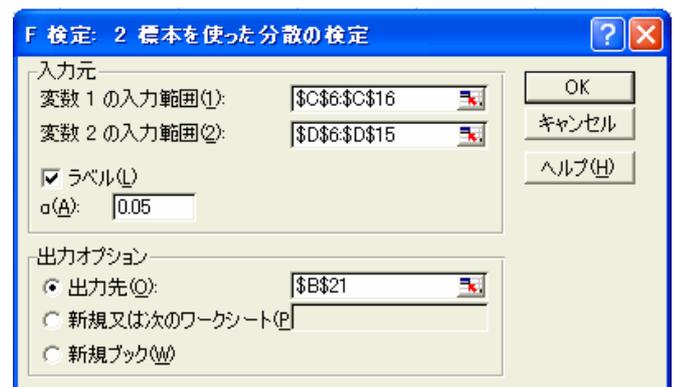
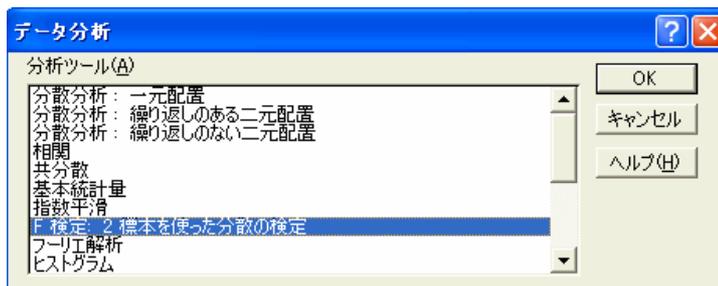
帰無仮説： $\sigma_A^2 = \sigma_B^2$ 2つの標本の母分散は同じである。

飼料Aと飼料Bで飼育したときにハムスターの体重増加量の分散に差がない。

対立仮説： $\sigma_A^2 \neq \sigma_B^2$ 2つの標本の母分散は異なる。

飼料Aと飼料Bで飼育したときにハムスターの体重増加量の分散に差がある。

帰無仮説が成り立つと仮定したときに今回のデータが得られる確率P値はエクセルの分析ツールで以下のように計算できる。



F-検定：2 標本を使った分散の検定		
	飼料A	飼料B
平均	72.7	70.55556
分散	5.344444	3.277778
観測数	10	9
自由度	9	8
観測された分散比	1.630508	
P(F<=f) 両側	0.251207	
F 境界値 両側	3.388124	

検定の結果

P値が0.2512であるから、5%の有意水準で帰無仮説は棄却できない。したがって、5%の有意水準で飼料Aと飼料Bでハムスターの体重増加量のばらつきに差があるとはいえない。

練習 商社Sはタイでタマネギを栽培している。スーパーの基準は厳しく、ほとんど同じ直径(cm)のタマネギでないと納入させてくれない。今まで栽培していた品種Aに代わり、多収で耐病性の強い品種Bを導入したいが直径のばらつきは品種Aと同じではないのかもしれない。無作為に選んだ標本から右のようなデータを得た。品種Bのばらつきは品種Aと異なるのかを有意水準5%で両側検定せよ。

	品種A	品種B
	10.2	11.0
	10.8	11.7
	10.6	10.9
	10.5	11.6
	10.7	11.5
	10.9	12.1
	10.4	11.8
	10.5	12.4
	10.6	11.5
	10.5	

B. 統計的検定についての補足

正規分布を仮定できない場合の検定

t分布, カイ二乗分布, F分布は母集団が正規分布するときに標本の統計量(それぞれ母平均, 母分散, 2つの母分散の比)がどのように分布するかを表したものである。したがって, 母集団が正規分布しないときにこれらの分布を用いた推定・検定をすると誤りをおかす可能性がある。ただし多少正規分布にずれた分布でもこれらの分布で検定してもおおむね大きな間違いはないことがわかっている。しかし, 正規分布を仮定することが危険である, あるいは正規分布に当てはめようのない分布ではノンパラメトリック検定(授業では説明しない)が利用できる。

C. 分散分析とは?

1. 3つ以上の標本平均から母平均が異なるかを同時に比較したい

例: ヤギの成長がよくなるという薬A, B, C, Dのうち, どれが効果があるかあるいはないかを比較したい。したがって, 対照区(コントロール, 薬を与えない区), A, B, C, Dの5つを比較することになる。処理: 薬, 水準: 5つである。

次のような結果を得た。

処理	対照区	A	B	C	D
	100	105	96	100	115
	102	108	97	97	112
	104	110	100	95	100
	105	106	102	104	105
	103	104	99	103	106
平均	102.8	106.6	98.8	99.8	107.6
標準偏差	1.92	2.41	2.39	3.83	5.94

一番値の高いものと次に高いもの, 2番目に高いものと3番目に高いものを次々とt検定すればよいのか?

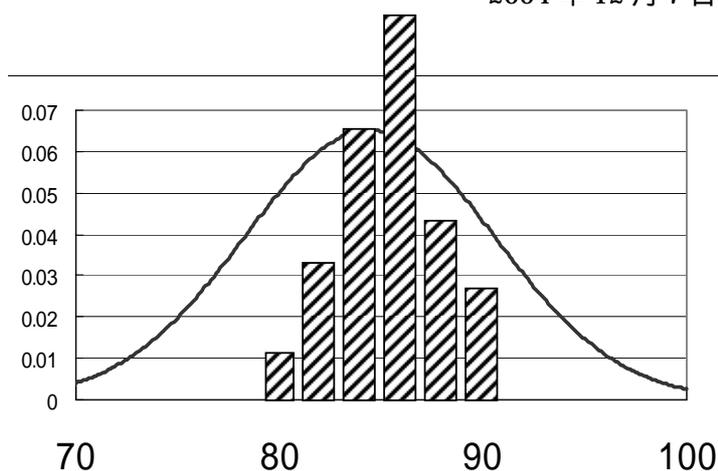
t検定をいくども行う問題点

t検定は決められた一組のデータについて行くと, 決められた有意水準のもとで帰無仮説が棄却できるかを検討できる。上のように5つのデータセットの場合, ${}_5C_2 = 10$ 回のt検定をすることになる。したがって, 全体で見たら有意水準(危険率)は5%以上になってしまう。

同一の正規分布から10個のデータを次々と取って, 一番平均の大きいものと一番平均の小さいものを選び, t検定するとどうなるか? 20回データを取ったときに実験すると・・・

ウズラの雄の体重の模擬実験

1000羽のウズラ(平均体重84.43g, 標準偏差6.12g)のデータから無作為(ランダム)に10羽を選んで, 平均と標準偏差を計算した。一番平均の大きいものと一番平均の小さいものをt検定すると, 5%の有意水準で有意差があった。



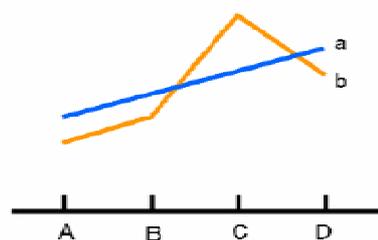
2. 2つ以上の要因を同時に比較したい

例：ヤギの成長がよくなる薬A, B, C, Dはえさと混ぜると効果が高いかもしれない。そう考
えて、麦わら, 稲わら, 濃厚飼料の3種類と組み合わせで試験をすることを考えた。処理は2
つ、薬と飼料であり、薬は5水準、飼料は3水準ある。このとき比較すべき試験は、薬A +
麦わら, 薬A + 稲わら, 薬A + 濃厚飼料, 薬B + 麦わら, 薬B + 稲わら, 薬B + 濃厚飼料, 薬
C + 麦わら, 薬C + 稲わら, 薬C + 濃厚飼料, 薬D + 麦わら, 薬D + 稲わら, 薬D + 濃厚飼料,
対照区 + 麦わら, 対照区 + 稲わら, 対照区 + 濃厚飼料となる。

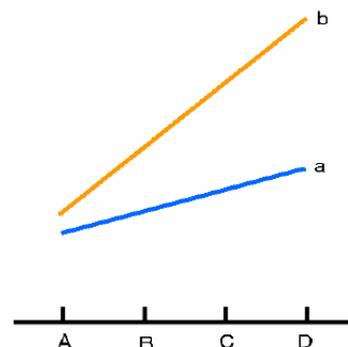
薬Aでは稲わらが一番よいえさだったが、薬Bでは麦わらが一番よいえさだったという結果が
出るかもしれない。t検定では2つの要因が絡み合った結果を解析できない。2つ以上の要因が
絡んだ結果(交互作用)はt検定では解析できない。

t検定ではわからないこと 交互作用

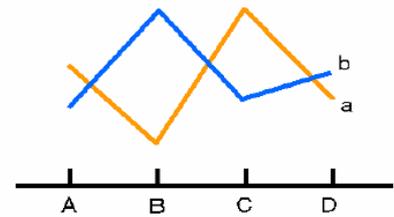
組み合わせの妙



相乗効果



打ち消しあう場合



交互作用のあるときは1つの要因（処理）だけを見て結論づけるわけには行かない．要因（処理）の組み合わせを考える必要が出てくる．

さらに要因を一つ一つ順番に調べるやり方では要因を組み合わせた場合の最適な条件にたどり着けるとは限らない．

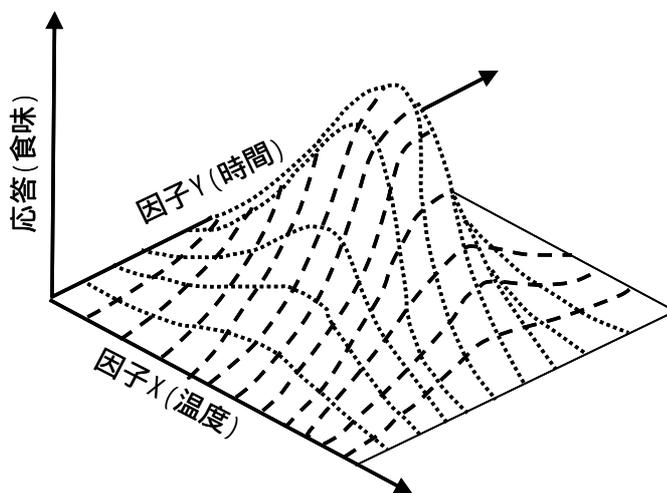


図 2 因子に対する応答曲面

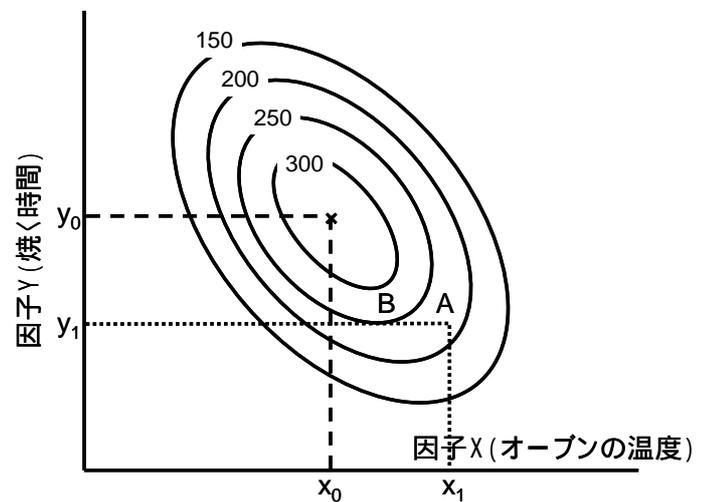


図 2 因子応答曲面の等高線図

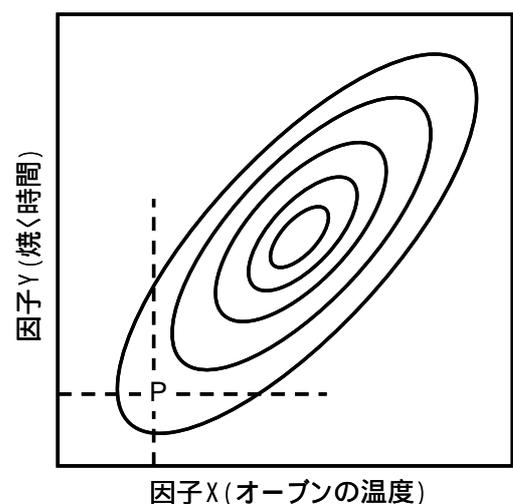


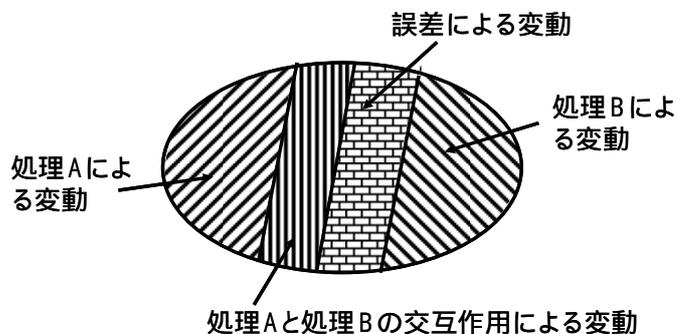
図 最適な水準に到達しない実験例
上の例では因子Xを固定したまま, 因子Yの水準を変えても応答は下がる. しかし, 真の最適水準の組み合わせは遙かに右上にある.

3. ばらつき（分散）を偶然誤差と意味のあるものに分解する

処理の間の違い，交互作用は本当に意味のあるものなのか？（有意であるか？）

実験データのばらつきは処理によって生じたもの（意味のある部分）と誤差によるばらつき（意味のない部分）に分けることができる．

分散の加法性



データの総変動 = 処理による変動 + 誤差変動

= 個々の処理による変動（主効果）+ 交互作用 + 誤差変動

2つの標本の分散の検定 F検定

処理による変動が誤差による変動に比べて，十分に大きいのか（正確には処理による変動と誤差による変動は等しいという帰無仮説が棄却できるのか）を検討する．

D. 分散分析の理論（一元配置の場合）

例：ハムスターをひまわり，大豆，人工餌の3種類のどれで育てるのが一番よいかを実験した．

実験結果に全く誤差がなく，餌の効果だけが現れたらどうなるか？

全く差がない場合，
実験結果は下のようになる

餌の効果に差があるなら，
（効果の合計=0）

餌の効果に差があるなら，
（効果の合計=0）

	ひまわり	大豆	人工餌		ひまわり	大豆	人工餌		ひまわり	大豆	人工餌
1	15	15	15	1	0	-3	3	1			
2	15	15	15	2	0	-3	3	2			
3	15	15	15	3	0	-3	3	3			
4	15	15	15	4	0	-3	3	4			
5	15	15	15	5	0	-3	3	5			

実験結果に誤差がランダムに適当に混ざっているならばどうなるか？

餌の効果が誤差なく発揮されると
実験結果は下のようになる

	ひまわり	大豆	人工餌
1			
2			
3			
4			
5			

誤差があるなら、
(誤差の合計=0)

	ひまわり	大豆	人工餌
1	2	3	1
2	0	0	-4
3	3	-2	-1
4	-3	0	-1
5	1	-2	3

左の2つを足すと
実験結果は下のようになる

	ひまわり	大豆	人工餌
1			
2			
3			
4			
5			

誤差によって、データの処理による違いがわかりにくくなった。
餌の効果が±3に対して、誤差も±4もあるからである。

得られた実験結果から誤差を分離してみる

実際に得られた実験結果は
下の通りである

	ひまわり	大豆	人工餌
1			
2			
3			
4			
5			

列の合計
列の平均
列の効果

左の実験結果から列の平均を
それぞれ引くと誤差を分離できる

	ひまわり	大豆	人工餌
1			
2			
3			
4			
5			

誤差の合計=0

列の効果を判定する：誤差に比べて十分に
大きいのか？

ばらつきのうち、誤差と効果によるものを分
けて、比較してみよう（右の表）。

誤差と効果を比較する

	繰り返し	ひまわり	大豆	人工餌
効果	1			
	2			
	3			
	4			
	5			
誤差	1			
	2			
	3			
	4			
	5			

分散分析を行う：効果による変動が誤差による変動に比べて十分に大きいのか？をF検定で検定する。

帰無仮説：効果による変動と誤差による変動には差がない。

誤差変動よりも処理の変動の方が大きいとかがえてよいかから片側検定となる。

上の帰無仮説のいうことは下の図のように読み替えることもできる

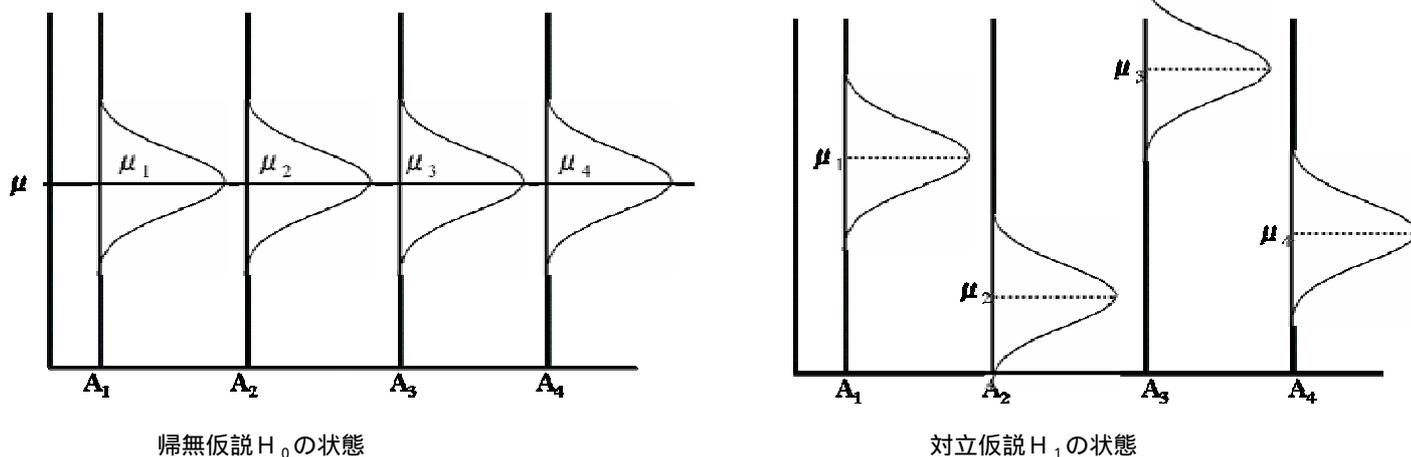


図 分散分析における帰無仮説と対立仮説

帰無仮説： $\mu_1 = \mu_2 = \mu_3 = \mu_4$ どの水準でも母平均は同じである

対立仮説：水準（処理）間の母平均のどれか一つは異なる

次の分散比Fを計算し、この分散比が得られる確率P値を計算する（実際はエクセルなどのソフトがデータを入力しただけでP値を計算する）

$$F = \frac{\text{効果のばらつき的大小} = \text{効果の分散}}{\text{誤差のばらつき的大小} = \text{誤差の分散}} = \frac{V_1}{V_2}$$

ここで 分散 $V = \frac{\sum (x_i - \bar{x})^2}{\phi}$

は自由度：データの数から使用した平均値の数を引いたもの

効果には0.6, -3.2, 2.6の3つしかなく、しかもこれらを作り出すために全体の平均を使っているので、自由度は $= 3 - 1 = 2$ となる。

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
0.6	0.6	0.36
0.6	0.6	0.36
0.6	0.6	0.36
0.6	0.6	0.36
0.6	0.6	0.36
<hr/>		
-3.2	-3.2	10.24
-3.2	-3.2	10.24
-3.2	-3.2	10.24
-3.2	-3.2	10.24
-3.2	-3.2	10.24
<hr/>		
2.6	2.6	6.76
2.6	2.6	6.76
2.6	2.6	6.76
2.6	2.6	6.76
2.6	2.6	6.76

したがって、右の表のように効果の分散を計算でき、 $V_1 = 43.4$ であることがわかった。

$$S = \sum (x_i - \bar{x})^2 = 86.80$$

$$V_1 = \frac{86.80}{2} = 43.4$$

誤差のデータは15個ある。使用した平均は各列にあるので、15.6, 11.8, 17.6の3つであり、自由度は
 $= 15 - 3 = 12$ となる。

したがって、右の表のように誤差の分散を計算でき、
 $V_2 = 5.43$ であることがわかった。

以上からF値は $F = \frac{V_1}{V_2} = \frac{43.4}{5.43} = 7.99$

このようなF値が得られる確率P値を計算すると、
 0.006229であることから（片側検定）

有意水準1%で帰無仮説は棄却される。

処理と誤差のばらつきには有意水準1%で有意な差がある。
 すなわちハムスターの成長は餌によって変化すると結論
 できる。

x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$
1.4	1.4	1.96
-0.6	-0.6	0.36
2.4	2.4	5.76
-3.6	-3.6	12.96
0.4	0.4	0.16

3.2	3.2	10.24
0.2	0.2	0.04
-1.8	-1.8	3.24
0.2	0.2	0.04
-1.8	-1.8	3.24

1.4	1.4	1.96
-3.6	-3.6	12.96
-0.6	-0.6	0.36
-0.6	-0.6	0.36
3.4	3.4	11.56

$$S = \sum (x_i - \bar{x})^2 = 65.20$$

$$V_2 = \frac{65.20}{12} = 5.43$$

分散分析では誤差のばらつきより処理のばらつきの方が大きいと判断してよいから、片側検定である。

E. 分散分析の実際（一元配置の場合）

1. 分散分析の結果の表現方法

分散分析の結果はふつう下のような表に書いて示す。

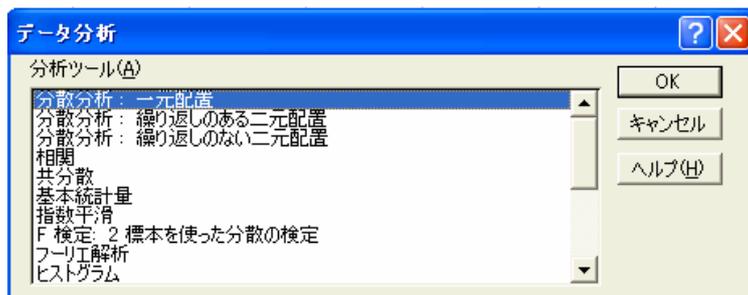
変動因	自由度	平方和 S	平均平方 (分散) V	分散比 F
処理間	A	S_A	V_A	F_0
誤差	E	S_E	V_E	
全体	T	S_T		

さらに有意水準5%, 1%で有意であればそれぞれ, *, **を F_0 の右肩につけるのが慣習となっている。有意差が検出されなかったときは ns をつけることもある。

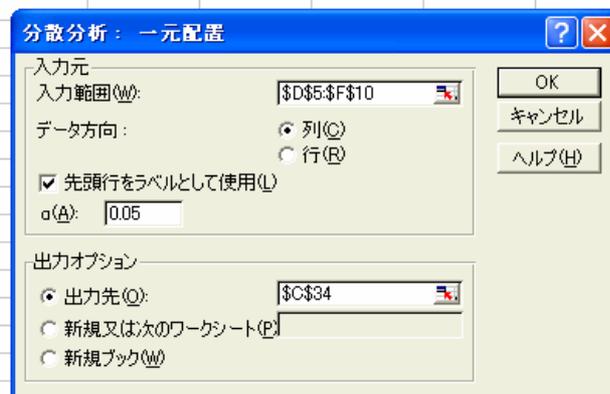
2. エクセルでの分散分析

例: ハムスターをひまわり, 大豆, 人工餌の3種類のどれで育てるのが一番よいかを実験した.

	ひまわり	大豆	人工餌
1	17	15	19
2	15	12	14
3	18	10	17
4	12	12	17
5	16	10	21



	ひまわり	大豆	人工餌
1	17	15	19
2	15	12	14
3	18	10	17
4	12	12	17
5	16	10	21



分散分析: 一元配置						
概要						
グループ	標本数	合計	平均	分散		
ひまわり	5	78	15.6	5.3		
大豆	5	59	11.8	4.2		
人工餌	5	88	17.6	6.8		
分散分析表						
変動要因	変動	自由度	分散	観測された分散比	P-値	F 境界値
グループ間	86.8	2	43.4	7.987730061	0.006229	3.88529
グループ内	65.2	12	5.433333			
合計	152	14				

エクセルで分散分析するときにはF値を調べる必要はない。P-値をみれば、帰無仮説が成り立つ場合、今回えられたようなデータが出現する確率(P値)がわかる。P値が0.05以下であれば有意水準5%で帰無仮説は棄却され、処理間の母平均に差があること、すなわち処理によってハムスターの成長が変わったことが示される。

3. 分散分析の結果の意味

分散分析は処理（水準）間の母平均に差がないという帰無仮説を検定している

帰無仮説： $\mu_1 = \mu_2 = \mu_3 = \mu_4$ どの水準でも母平均は同じである

対立仮説：水準（処理）間の母平均のどれか一つは異なる

したがって、分散分析の結果、有意差があることが分かった場合、その意味するところは、処理（水準）の中で一つは母平均の異なるものがある、すなわち処理によって変わるということである。一般的には一番値の小さいものと一番値の大きいものとの間には有意差があるということになる。それ以外の処理（水準）間に差があるかは分散分析ではわからない。これを調べるのは多重比較法である。多重比較法については後の授業で簡単に説明する。

練習 1 :

例: ヤギの成長がよくなるという薬A, B, C, Dのうち、どれが効果があるかあるいはないか、どれが一番効果があるかを比較したい。したがって、対照区(コントロール, 薬を与えない区), A, B, C, Dの5つを比較することになる。処理: 薬, 水準: 5つである。
次のような結果を得た。

処理	対照区	A	B	C	D
	100	105	96	100	115
	102	108	97	97	112
	104	110	100	95	100
	105	106	102	104	105
	103	104	99	103	106
平均	102.8	106.6	98.8	99.8	107.6
標準偏差	1.92	2.41	2.39	3.83	5.94

練習2：右の例で水稻の1穂穎花数には品種間差があるかを分散分析せよ。

水稻の1穂穎花数の品種間差異を5品種について調査した。					
	コシヒカリ	日本晴	NPT160	IR72	タカナリ
	56	80	188	131	177
	52	84	169	129	207
	64	64	131	106	179
	60	71	233	95	173
	17	81	186	152	109
	98	93	183	152	151
	122	107	167	145	143
	101	104	135	144	162
	128	77	148	153	123
	85	60	126	83	188
	125	96	177	160	96
	60	111	158	115	152
	108	78	127	113	162
	85	80	162	71	111
	83	80	117	76	139
	97	94	144	101	154
	84	85	165	76	132
	108	92	98	68	214
平均	85.16667	85.38889	156.33333	115	154
標準偏差	29.40038	13.9332	31.93007	31.8932	32.76476

帰無仮説：

対立仮説：

検定結果：

4. 分散分析をする上での注意

なるべく反復数はそろえる

今回学んだ一元配置の分散分析では反復数が個々の処理区で異なってもそれほど問題はない。しかし、来週以降に学ぶより複雑な分散分析では、反復数が異なると解析が面倒になるだけでなく、精度も大きく落ちてしまう。実験開始のときは反復数をそろえて実験するのが普通であるが、事故や不注意なので反復数がそろわなくなることもあるかもしれない。しかし、できるだけ反復のそろうように実験することが基本である。なお反復がそろわないからといって、一部のデータを削除するのは間違ったやり方である。

複雑な実験はなるべく避ける

前項とも関連するが、分散分析ではデータが複雑になるほど、解析が面倒かつ間違いやすくなる。特にコンピューターで計算させるときは、データの入力の仕方を間違いやすくなり、自分の目的とする分散分析をするにはデータの構造が複雑（あるいはでたらめ）で、解析不能ということもあり得る。そのうえ、反復がそろわなかったときの影響も大きくなる。必要のない複雑な実験は避けるのはもちろんのこと、必要だとしてもできるだけ簡単な実験計画にならないかをよく検討してから実験するべきである。実験計画を立てた時点で、どういう分散分析をするのかを決めておくのが正しい統計解析方法である。

正規分布するデータが前提条件である

分散分析では比較する母集団それぞれが正規分布すること，母分散が等しいことが理論的には前提条件になる．しかし，分散分析は多少その前提条件からはずれていても，結果が大きく左右されない頑健性をもっている．

すべての水準に対して母分散が等しいことも前提条件である

水準ごとの反復数がみな同じである場合，この前提が多少崩れても影響はあまりない．

フィッシャーの3原則を満たした実験計画のもとで，分散分析を行う

誤差に系統誤差が入ると解析結果の妥当性が失われる危険性がある．系統誤差を除去したり，分散分析の解析の妨害とならない偶然誤差に転化するのがフィッシャーの3原則に述べられた反復，無作為化，局所管理である．詳細は後の授業で学ぶ予定である．

F．宿題

1．第5回の宿題で調査したデータを用いて，卵の重さのばらつきが2つの店の間で同じであることを有意水準5%でF検定（両側検定）せよ．

2．第7回の宿題3．で調べたデータについて分散分析せよ．

3．次回の授業では二元配置の分散分析について学ぶ．2つの処理を組み合わせたときのデータの解析方法である．

産地と品種の違いがリンゴの糖度に及ぼす影響を知りたい．

施肥量と品種の違いがイネの1穂穎花数に及ぼす影響を知りたい．

作期と品種の違いがトマトの酸度に及ぼす影響を知りたい．

二元配置で分散分析できそうなデータを研究室の卒業実験などから手に入れる．あるいは自分で実験してもかまわない．なお，データには反復が2つ以上あり，それぞれの処理で反復の数は同じにすること．データの配置は下の図のようになる．

産地と品種の違いがリンゴの糖度に及ぼす影響を各処理3つずつのリンゴで調査した．

産地：青森，長野，山形， 品種：紅玉，ふじ，むつ

	青森	長野	山形
紅玉	() () ()	() () ()	() () ()
ふじ	() () ()	() () ()	() () ()
むつ	() () ()	() () ()	() () ()

他にも次のような例が考えられる

コンビニ弁当の売れ行き 曜日，天気の2つの要因

イネの1穂穎花数 品種，施肥の2つの要因

小鳥のさえずる回数 気温，太陽の明るさの2つの要因