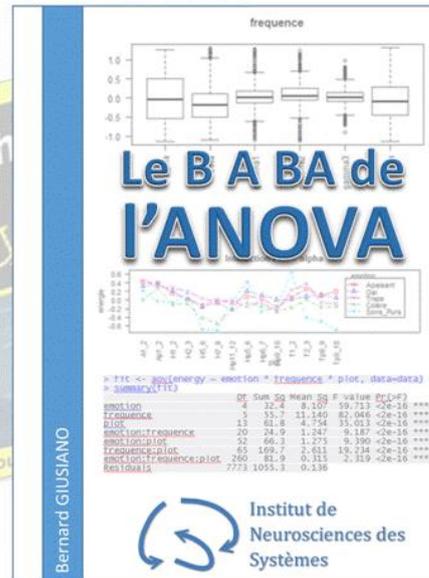




Qu'est-ce que c'est ?
A quoi ça sert ?
Pourquoi ne pas se satisfaire
de simples tests de t 2 par 2 ?



Trois individus (ou objets, ou sujets, ...)



à chacun desquels on affecte une mesure

44

44

52

Bleu = Rouge Rouge ≠ Vert Bleu ≠ Vert

Trois groupes d'individus



dont on calcule la moyenne des mesures

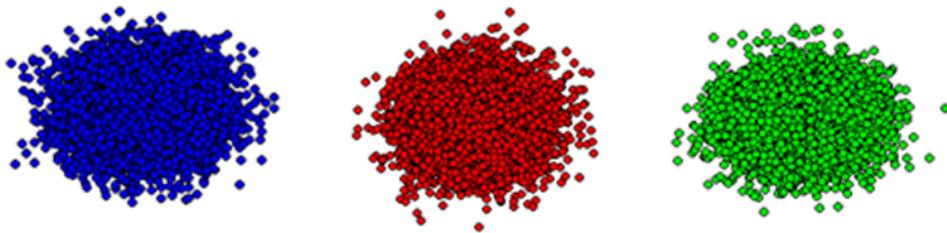
42

42

47

$$m_{\text{bleu}} = m_{\text{rouge}} \quad m_{\text{rouge}} \neq m_{\text{vert}} \quad m_{\text{bleu}} \neq m_{\text{vert}}$$

Trois très grands groupes d'individus (on parle de *populations*)



pour lesquels on ne peut accéder aux
mesures de tous les individus

$$m_{\text{bleu}} \stackrel{?}{=} m_{\text{rouge}} \quad m_{\text{rouge}} \stackrel{?}{\neq} m_{\text{vert}} \quad m_{\text{bleu}} \stackrel{?}{\neq} m_{\text{vert}}$$

On peut démontrer que si on tire un **échantillon** représentatif de la population (*par tirage au hasard de 30 individus, par exemple*), la moyenne de cette échantillon (\bar{x}) est une bonne **estimation** ($\hat{\mu}$) de la moyenne de la population (μ) :

$$\hat{\mu} = \bar{x}$$

Bonne estimation : l'espérance mathématique de $\hat{\mu}$ est égale à μ

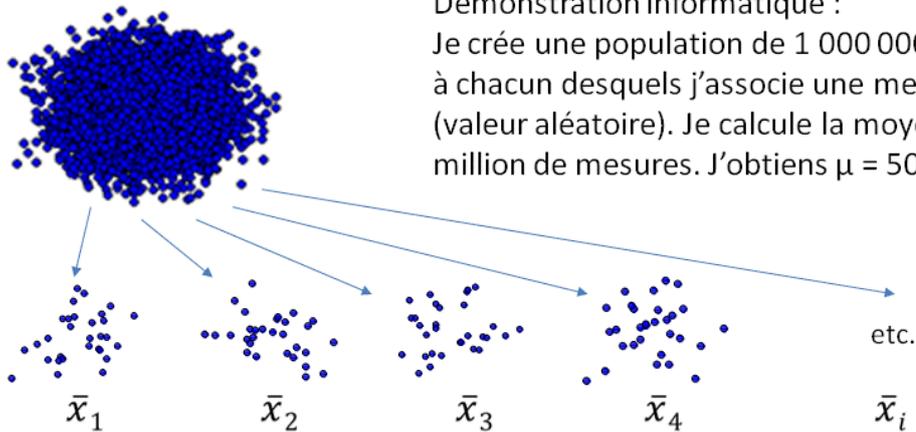
$$E(\hat{\mu}) = \mu$$

Ce qui signifie que si on tire un nombre très grand d'échantillons de même effectif, la moyenne (*espérance*) des moyennes \bar{x} de ces échantillons tend vers μ , moyenne de la population.



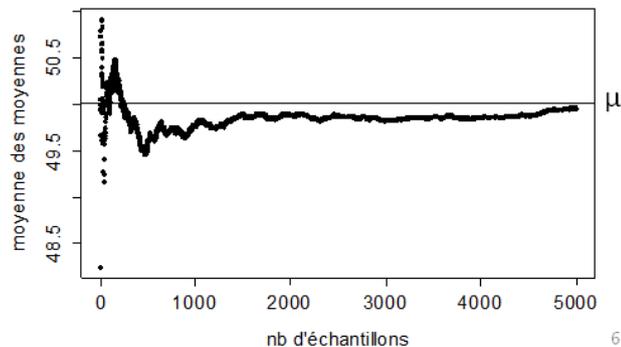
17/11/2015

B A B A de l'ANOVA



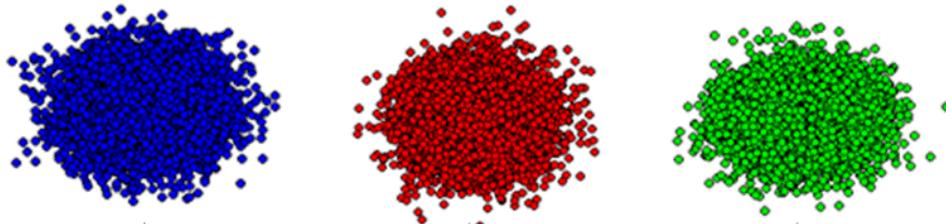
Démonstration informatique :
Je crée une population de 1 000 000 individus à chacun desquels j'associe une mesure (valeur aléatoire). Je calcule la moyenne de ce million de mesures. J'obtiens $\mu = 50.01944$

Je tire i échantillons de 30 individus chacun. Je calcule leur moyenne. Quand i tend vers l'infini, la moyenne des \bar{x}_i tend vers μ .

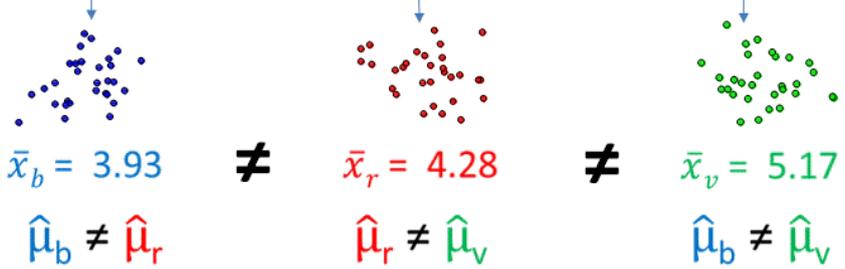


17/11/2015

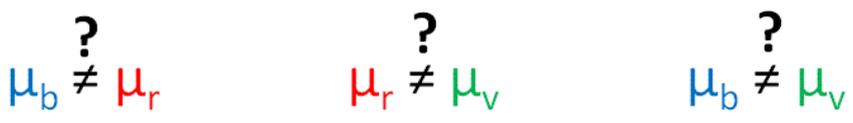
6



Tirage au hasard de 30 individus dans chacune des 3 populations



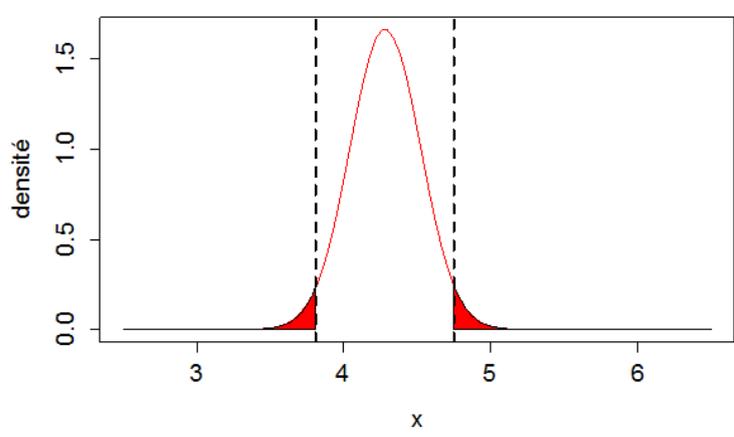
Peut-on en conclure que les moyennes des 3 populations sont différentes ?



NON, car $\hat{\mu}$ n'est pas μ , ce n'en est qu'une **estimation**.

Sa valeur doit être accompagnée de son **intervalle de confiance** qui donne la précision de l'estimation de la moyenne de la population à partir de l'échantillon.

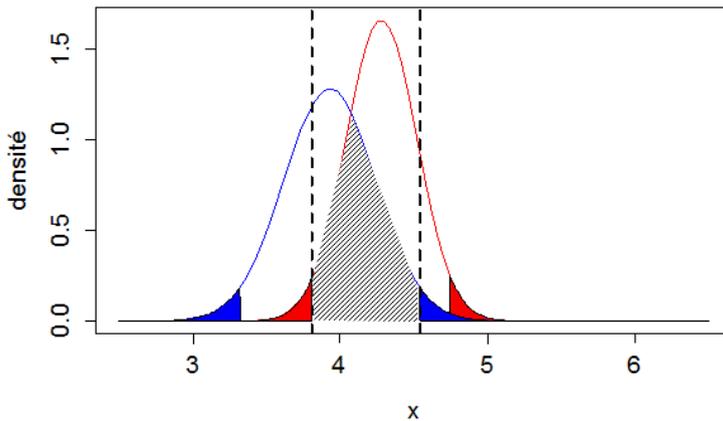
Si on répétait l'estimation un grand nombre de fois, dans $(1-\alpha)$ % des cas l'intervalle de confiance contiendrait la vraie moyenne de la population, μ .



Dans cette estimation à partir d'un échantillon de la **population rouge**, la moyenne μ pourrait bien être dans l'intervalle **[3.81 ; 4.75]** avec un **degré de confiance** de 95%.

$\hat{\mu}_r$ et $\hat{\mu}_b$ ne sont que des **estimations**.

Si la **population rouge** et la **population bleue** ont la même moyenne, μ , celle-ci pourrait bien être dans l'intervalle **[3.81 ; 4.54]** avec un **degré de confiance** équivalent à l'aire hachurée.



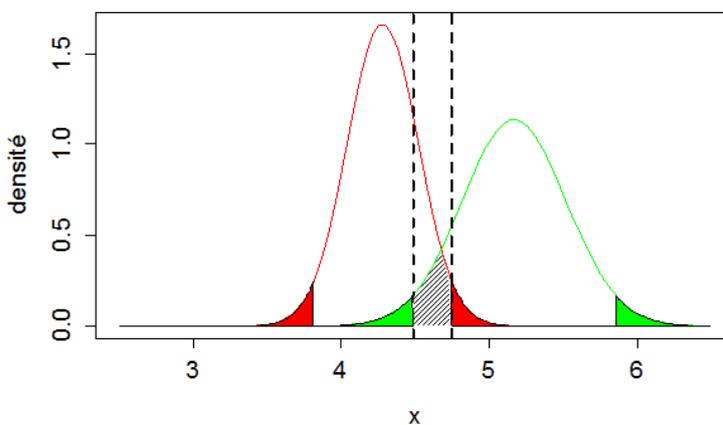
17/11/2015

B A BA de l'ANOVA

9

$\hat{\mu}_r$ et $\hat{\mu}_v$ ne sont que des **estimations**.

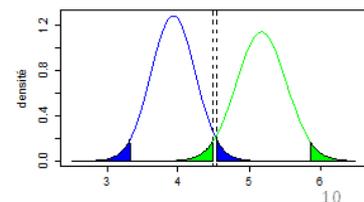
Si la **population rouge** et la **population verte** ont la même moyenne, μ , celle-ci pourrait bien être dans l'intervalle **[4.48 ; 4.75]** avec un **degré de confiance** équivalent à l'aire hachurée.



17/11/2015

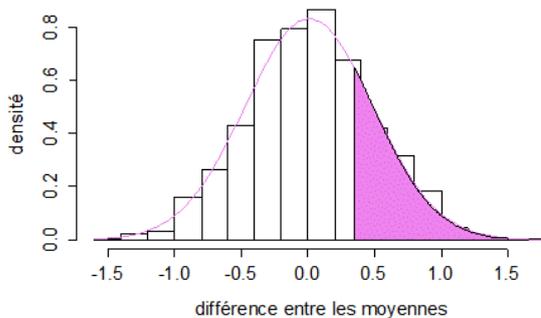
B A BA de l'ANOVA

On comprend que si le degré de confiance est « *trop petit* », on choisira l'**hypothèse** la plus vraisemblable : les deux populations n'ont pas la même moyenne.



Faisons l'**hypothèse** que $\mu_r = \mu_b$ alors $\mu_r - \mu_b = 0$

Si nous répétons un grand nombre de fois le tirage d'un échantillon dans chacune de deux populations de même moyenne $\mu = \mu_r = \mu_b$, que nous calculons chaque fois la différence entre les moyennes des deux échantillons, **cette différence suit une loi normale de moyenne $\mu = 0$.**



Reprenons nos échantillons

rouge et **bleu** :

$\bar{x}_r = 4.28$ et $\bar{x}_b = 3.93$

donc $\bar{x}_r - \bar{x}_b = 0.35$

La probabilité que la différence de moyennes entre les deux échantillons soit égale ou supérieure à 0.35, alors que les populations dont ils sont tirés ont des moyennes égales (**hypothèse nulle**), est donnée par l'aire colorée en violet. C'est le « **petit p** ».

17/11/2015

B A BA de l'ANOVA



11

Conclusion de ce rappel

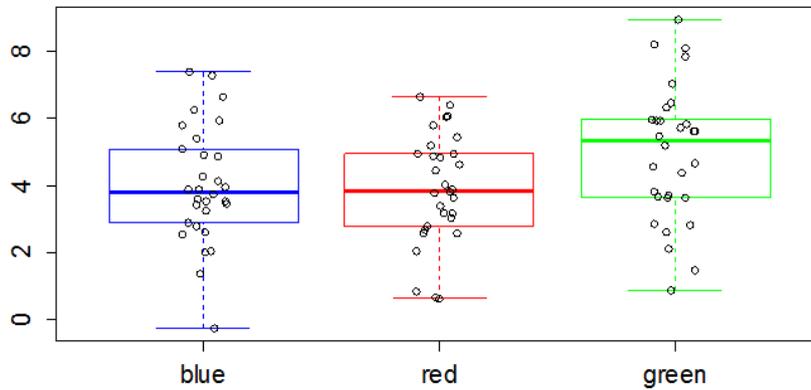
- Si la question est :
« **les moyennes des échantillons** sont-elles différentes ? »
nous n'avons pas besoin d'un test d'hypothèse.
Exemple : Est-ce que la mesure moyenne est différente entre mes 6 patients et mes 8 sujets contrôle ?
- Si la question est :
« **au vu des échantillons, est-ce que les moyennes des populations** dont ils sont tirés sont différentes ? »
nous avons besoin d'un **test d'hypothèse**...
... et nous prenons des risques
liés à l'incertitude sur les estimations.
Exemple : A partir de mes 6 patients et de mes 8 sujets contrôle, est-il raisonnable de dire qu'en général la mesure moyenne des patients est différente de celle des sujets contrôle ?
- L'incertitude est contrôlée par les **conditions d'application** des tests.

17/11/2015

B A BA de l'ANOVA

12

ANOVA à 1 facteur (ayant 3 modalités) - 1 way ANOVA (3 levels)



Nombre de groupes : $k = 3$ {bleu, rouge, vert}

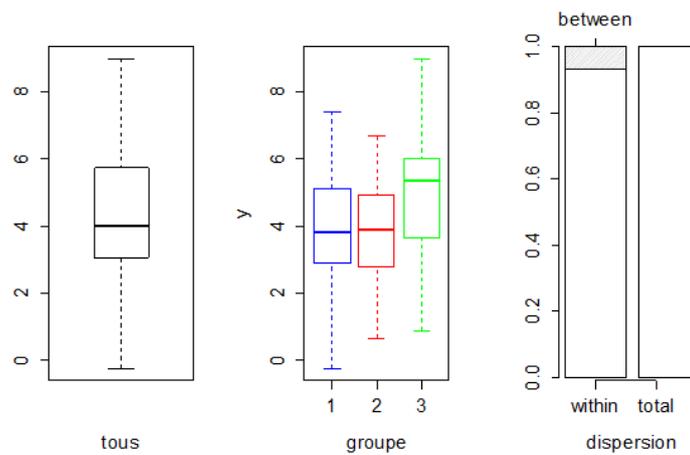
Effectif par groupe : $n_{\text{bleu}} = n_{\text{rouge}} = n_{\text{vert}} = 30$

Effectif total : $n = 90$

L'analyse de la variance Analysis Of VAriance (ANOVA)

Elle repose sur la décomposition de la variance totale :

$\text{variance}_{\text{totale}} \equiv \text{variance}_{\text{entre les groupes}} + \text{variance}_{\text{à l'intérieur des groupes}}$



Mais il s'agit bien de comparer des moyennes :

variance_{totale} : différence entre valeurs et moyenne globale

$$SCE_{totale}^* = \sum_i^n (x_i - \bar{x})^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$$

variance_{entre les groupes} : différence entre moyennes des groupes et moyenne globale

$$SCE_{inter\ groupe} = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2 = \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_j - \bar{x})^2$$

variance_{à l'intérieur des groupes} : différence entre valeurs et moyennes des groupes

$$SCE_{intra\ groupe} = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$$

*SCE : Somme des Carrés des Ecartés à la moyenne

17/11/2015

B A BA de l'ANOVA



15

La variance est la moyenne des carrés des écarts :

$$\text{estimation de la variance}_{totale} = S_{totale}^2 = \frac{SCE_{totale}}{ddl_{total}^*} = \frac{SCE_{totale}}{n-1}$$

$$\text{estimation de la variance}_{entre les groupes} = S_{intra}^2 = \frac{SCE_{inter}}{ddl_{inter}} = \frac{SCE_{inter}}{k-1}$$

$$\text{estimation de la variance}_{à l'intérieur des groupes} = S_{inter}^2 = \frac{SCE_{intra}}{ddl_{intra}} = \frac{SCE_{intra}}{n-k}$$

$$SCE_{totale} = SCE_{inter} + SCE_{intra}$$

*ddl : nombre de degrés de liberté

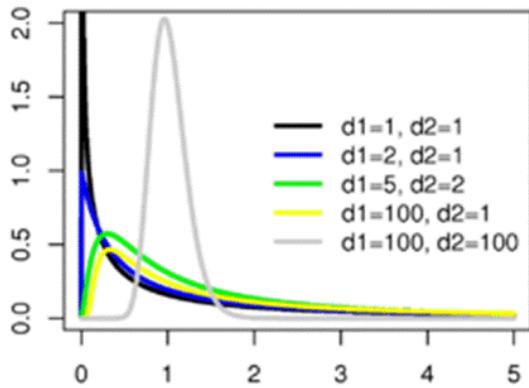
17/11/2015

B A BA de l'ANOVA

16

Dans l'hypothèse nulle (H_0) où il n'y a pas de différence entre les moyennes des 3 populations d'où sont tirés les 3 échantillons,

le rapport entre la variance inter et la variance intra suit une loi de distribution connue, la loi du F de Fisher-Snedecor, à $v_1 = k - 1$ et $v_2 = n - k$ degrés de liberté.



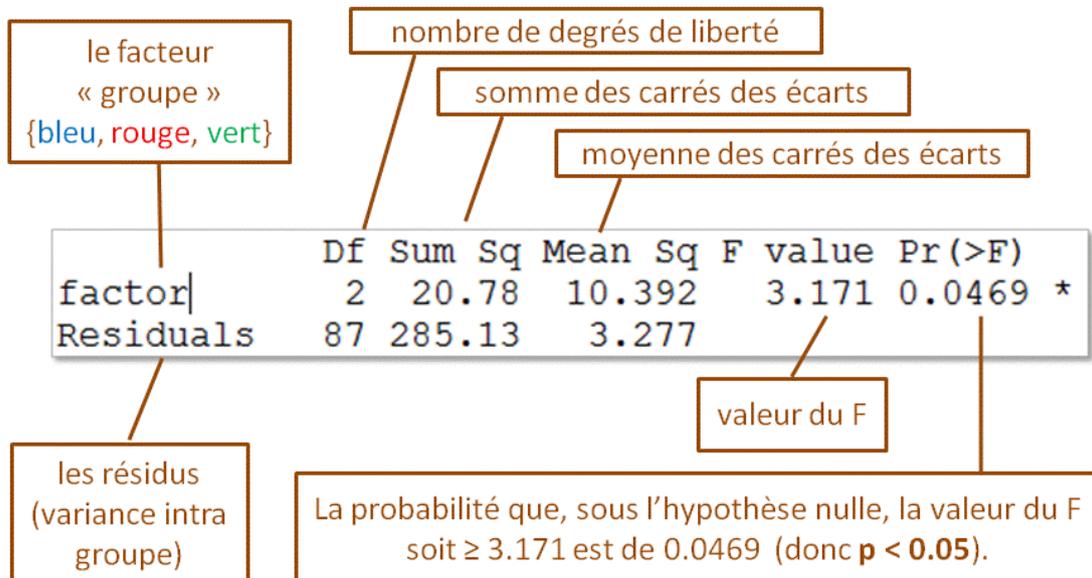
$$F = \frac{S_{inter}^2}{S_{intra}^2} = \frac{S_{factorielle}^2}{S_{résiduelle}^2}$$

Table de Fisher-Snedecor, $\alpha = 5\%$ (95^e centile)

v_1 (numérateur)	v_2 (dénominateur)	1	2	3	4	5	6	7	8	9	10	20	30	40	50	60	80	100	
1	1	161.45	199.50	215.71	224.58	230.16	233.99	236.77	238.88	240.54	241.88	243.00	243.95	244.75	245.45	246.07	246.63	247.14	247.61
2	1	18.51	19.00	19.15	19.25	19.32	19.37	19.41	19.44	19.46	19.47	19.48	19.49	19.50	19.50	19.51	19.51	19.52	19.52
3	1	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85	8.81	8.79	8.76	8.75	8.74	8.73	8.73	8.73	8.73	8.73
4	1	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	6.00	5.96	5.90	5.85	5.82	5.80	5.79	5.78	5.78	5.78
5	1	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.66	4.60	4.56	4.54	4.53	4.53	4.53	4.53
6	1	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.10	4.06	3.97	3.91	3.87	3.85	3.84	3.84	3.84	3.84
7	1	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.68	3.64	3.54	3.48	3.44	3.42	3.41	3.41	3.41	3.41
8	1	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.39	3.35	3.25	3.18	3.14	3.12	3.11	3.11	3.11	3.11
9	1	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.18	3.14	3.04	2.96	2.92	2.90	2.89	2.89	2.89	2.89
10	1	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.87	2.79	2.75	2.73	2.72	2.72	2.72	2.72
20	1	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.24	2.16	2.12	2.10	2.09	2.09	2.09	2.09
30	1	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.21	2.16	2.05	1.97	1.93	1.91	1.90	1.90	1.90	1.90
40	1	4.08	3.23	2.84	2.61	2.45	2.34	2.25	2.19	2.13	2.08	1.97	1.89	1.85	1.83	1.82	1.82	1.82	1.82
50	1	4.03	3.18	2.79	2.56	2.40	2.29	2.20	2.13	2.07	2.02	1.91	1.83	1.79	1.77	1.76	1.76	1.76	1.76
60	1	4.00	3.15	2.76	2.53	2.37	2.26	2.17	2.10	2.04	1.99	1.87	1.79	1.75	1.73	1.72	1.72	1.72	1.72
70	1	3.98	3.13	2.74	2.51	2.35	2.24	2.15	2.07	2.01	1.96	1.84	1.76	1.72	1.70	1.69	1.69	1.69	1.69
80	1	3.96	3.11	2.72	2.49	2.33	2.22	2.13	2.06	2.00	1.95	1.83	1.75	1.71	1.69	1.68	1.68	1.68	1.68
90	1	3.95	3.10	2.71	2.47	2.31	2.20	2.11	2.04	1.98	1.93	1.81	1.73	1.69	1.67	1.66	1.66	1.66	1.66
100	1	3.94	3.09	2.70	2.46	2.31	2.19	2.10	2.03	1.97	1.92	1.80	1.72	1.68	1.66	1.65	1.65	1.65	1.65
100	100	3.93	3.08	2.69	2.45	2.29	2.18	2.09	2.02	1.96	1.91	1.79	1.71	1.67	1.65	1.64	1.64	1.64	1.64

Densité de probabilité (ou fonction de masse)

ANOVA des données de notre exemple :



Conclusion : rejet de l'hypothèse nulle (au risque α), trop improbable. Les moyennes des trois groupes sont significativement différentes.

Conditions d'application de l'ANOVA

1. Indépendance des échantillons

- La meilleure solution : tirer les échantillons au hasard.

2. Normalité des populations.

- Il faut pouvoir considérer que les populations d'où sont tirés les échantillons sont distribuées selon la loi normale.

3. Egalité des variances (**homoscédasticité**)

- Le facteur étudié peut influencer les moyennes, mais pas les variances.

Des méthodes existent pour vérifier ces 3 conditions d'application.

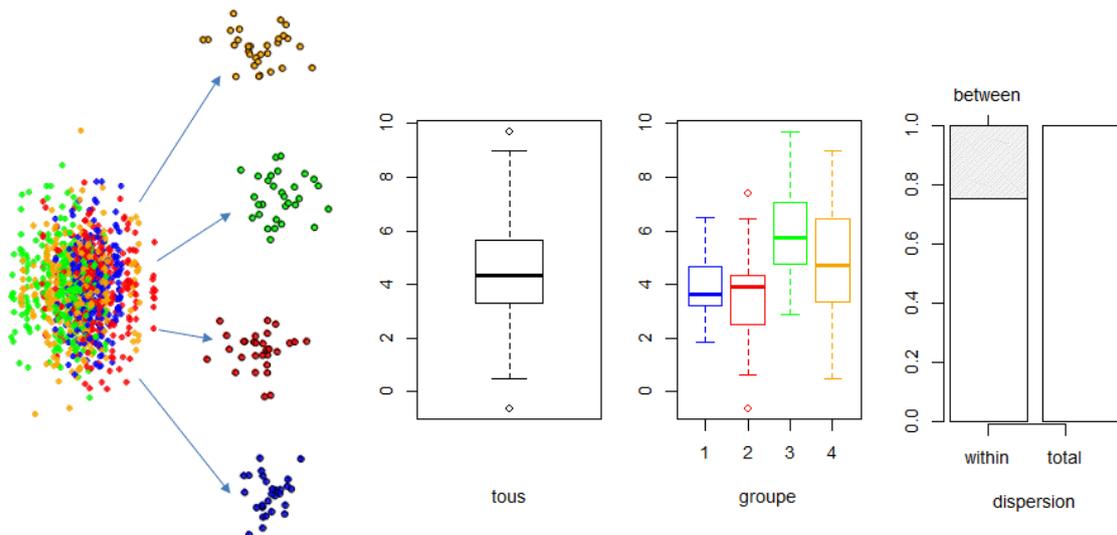
L'indépendance est la condition la plus critique.

L'ANOVA est robuste vis-à-vis de la normalité stricte et de l'homoscédasticité si les échantillons sont assez grands (> 30).

Si ces conditions ne peuvent pas être satisfaites, on peut utiliser le **test de Kruskal-Wallis** comme équivalent non paramétrique de l'analyse de variance à un facteur.

Mais quel groupe est différent de quel autre groupe ?

On utilise un autre exemple : une population dont on tire quatre échantillons indépendants et de même effectif ($n_j = 30$), un pour chaque couleur.

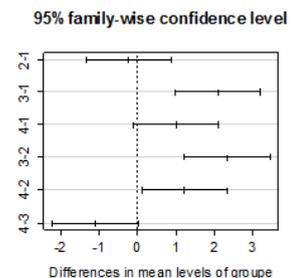


ANOVA de ce nouvel exemple :

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
factor	3	102.8	34.28	12.5	3.86e-07 ***
Residuals	116	318.2	2.74		

Tests post-hoc de Tukey (Tukey Honest Significant Differences) :

	diff	lwr	upr	p adj
2-1	-0.2301661	-1.3448419	0.88450977	0.9495035
3-1	2.1078624	0.9931865	3.22253825	0.0000164
4-1	1.0041868	-0.1104890	2.11886266	0.0932113
3-2	2.3380285	1.2233526	3.45270433	0.0000016
4-2	1.2343529	0.1196770	2.34902874	0.0237150
4-3	-1.1036756	-2.2183514	0.01100026	0.0533531



La méthode de Tukey consiste à déterminer la différence minimum entre deux groupes qui peut être considérée comme significative.

L'ajustement de p pour les comparaisons multiples (ici au nombre de 6) est assuré par la référence à la distribution des intervalles de Student.

cf. Wikipedia : Tukey's range test

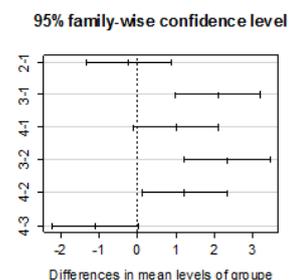
Quelle différence avec des test de t deux à deux ?

Si on réalise les 6 tests de comparaisons 2 à 2 avec le test du t de Student, se pose le problème des comparaisons multiples :

$$\alpha_{\text{global}} \approx \alpha_{\text{chaque test}} \times \text{nb de comparaisons} \Rightarrow \alpha_{\text{global}} \approx 0.05 \times 6 = 0.30$$

Traitons-le avec la correction de Bonferroni : $\alpha_{\text{corrigé}} = \alpha / \text{nb de comparaisons}$

	Tukey	test t	test t
difference	p adj	p	p adj
2-1	0.9495035	0.5417	1.0
3-1	0.0000164	0.0000006	0.0000039
4-1	0.0932113	0.02006	0.1203465
3-2	0.0000016	0.0000016	0.0000098
4-2	0.0237150	0.01237	0.0742211
4-3	0.0533531	0.0237	0.1422164

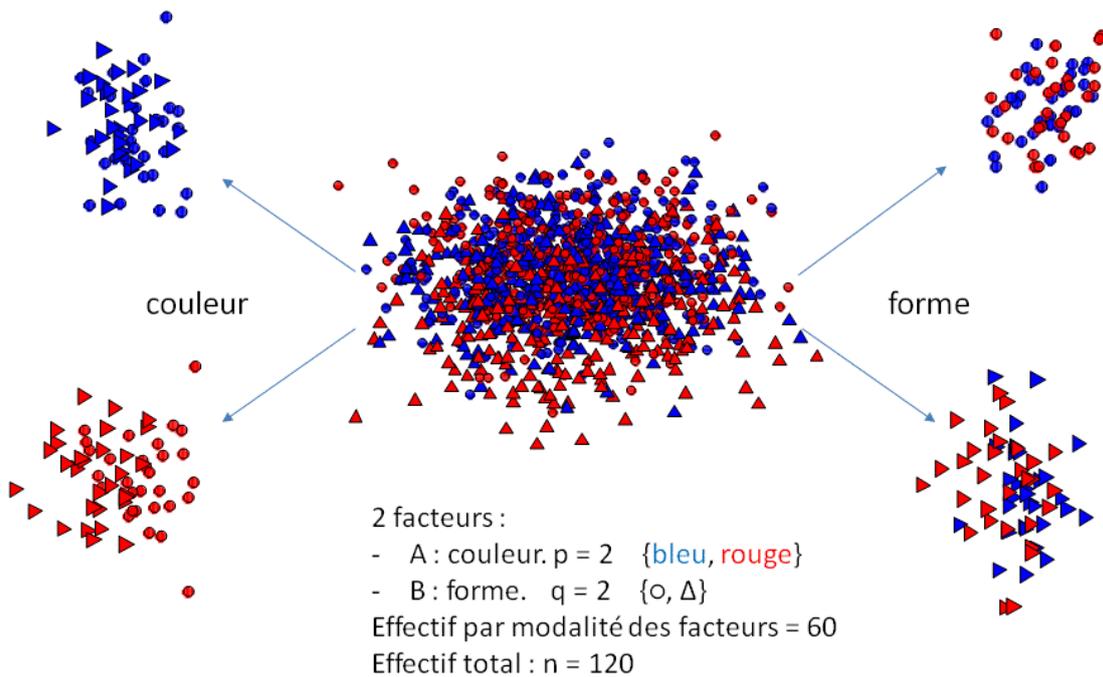


Le test de t n'a pas mis en évidence la différence entre 4 et 2.

Le test de Tukey est donc plus **puissant** que le test de t corrigé.

Il est aussi plus **robuste** car il utilise toutes les données de l'expérience chaque fois, et pas seulement les données des deux groupes testés.

Autre intérêt de l'ANOVA : l'analyse multifactorielle

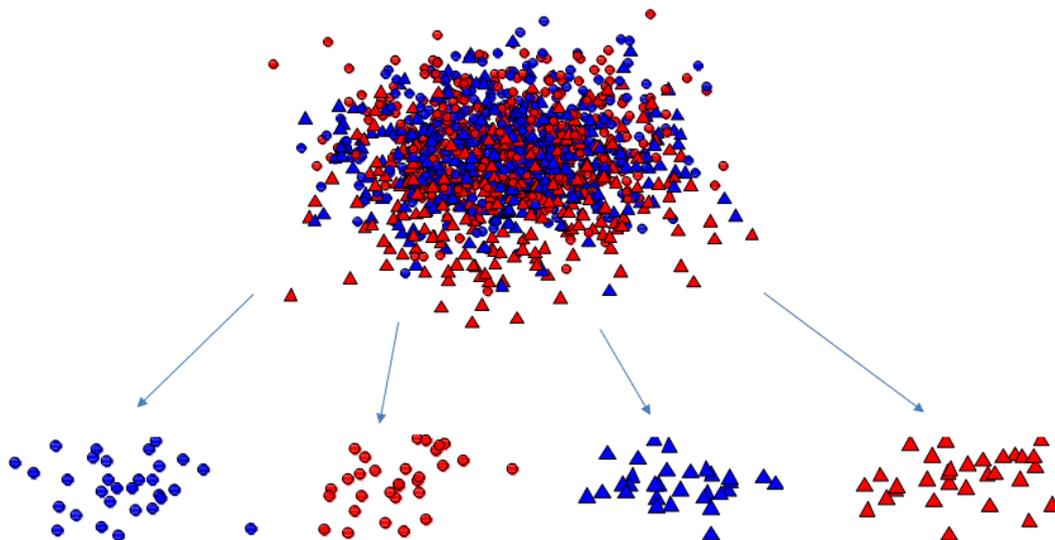


17/11/2015

B A BA de l'ANOVA

23

ANOVA à 2 facteurs (2 modalités) - 2 ways ANOVA (2 levels)



17/11/2015

B A BA de l'ANOVA

24

ANOVA à 2 facteurs (2 modalités) - 2 ways ANOVA (2 levels)

$$SCE_{totale} = SCE_{inter} + SCE_{intra}$$

$$SCE_{Totale} = SCE_{Factorielle} + SCE_{Résiduelle}$$

$$SCE_{Factorielle} = SCE_{Facteur A} + SCE_{Facteur B} + SCE_{Interaction AB}$$

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
A	1	2.10	2.10	2.537	0.11390
B	1	32.46	32.46	39.243	6.56e-09 ***
A:B	1	6.85	6.85	8.284	0.00476 **
Residuals	116	95.96	0.83		

ANOVA à 2 facteurs (2 modalités) - 2 ways ANOVA (2 levels)

Error Decomposition

$$\underbrace{\sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^r (Y_{ijk} - \bar{Y}_{...})^2}_{SS_{Total}} = \underbrace{r \cdot b \cdot \sum_{i=1}^a (\bar{Y}_{i..} - \bar{Y}_{...})^2}_{SS_A} + \underbrace{r \cdot a \cdot \sum_{j=1}^b (\bar{Y}_{.j.} - \bar{Y}_{...})^2}_{SS_B} + \underbrace{r \times \sum_{i=1}^a \sum_{j=1}^b (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...})^2}_{SS_{A \times B}} + \underbrace{\sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^r (Y_{ijk} - \bar{Y}_{ij.})^2}_{SS_{within}}$$

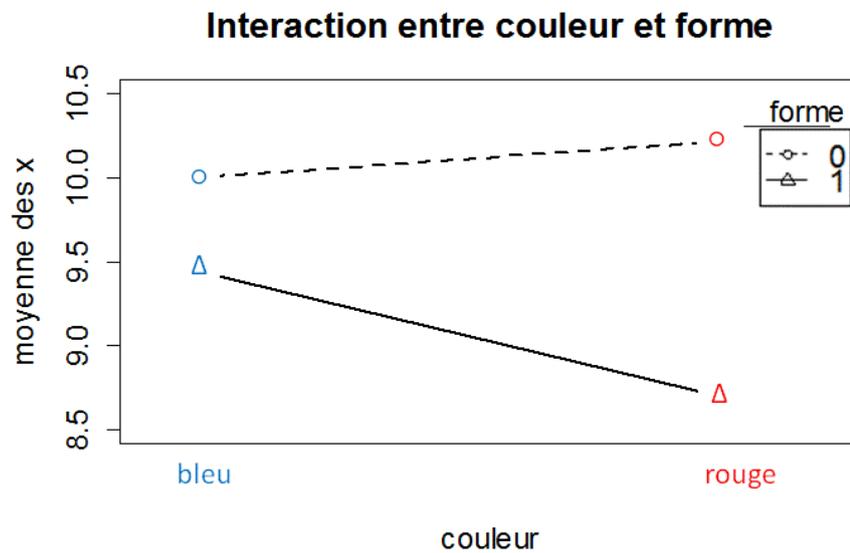
ANOVA Table

Source	Degrees of Freedom	SS	MS	F
A	a-1	SS_A	MS_A	MS_A / MS_{within}
B	b-1	SS_B	MS_B	MS_B / MS_{within}
A × B	(a-1)(b-1)	$SS_{A \times B}$	$MS_{A \times B}$	$MS_{A \times B} / MS_{within}$
Within	ab(r-1)	SS_{within}	MS_{within}	
Total	abr-1	SS_{Total}		

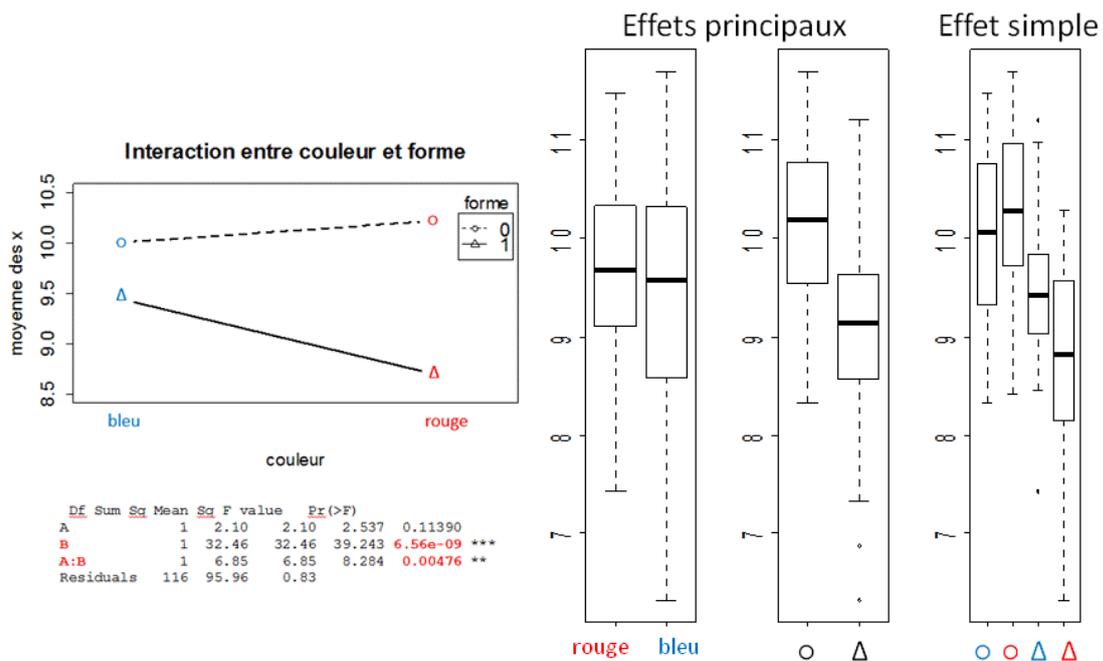


Interaction : l'effet d'un facteur dépend du niveau de l'autre facteur.

Du point de vue de leur moyenne, $\circ < \circ$ alors que $\Delta > \Delta$

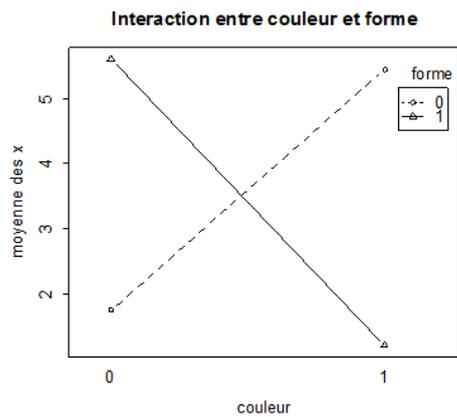


Interaction : effet simple et effets principaux

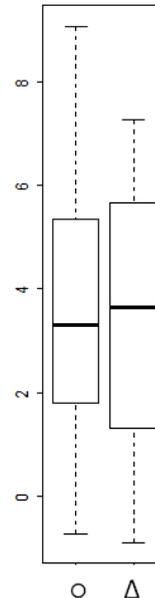
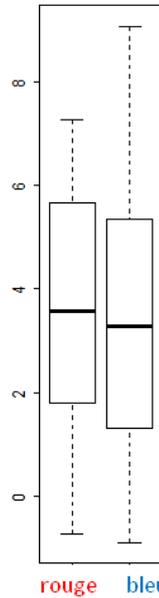


L'absence d'effets principaux peut cacher des effets simples !

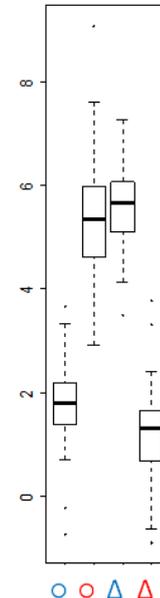
	Df	Sum Sq	Mean Sq	F value	Pr(>F)
A	1	0.0	0.0	0.019	0.892
B	1	0.0	0.0	0.033	0.856
A:B	1	470.0	470.0	505.571	<2e-16 ***
Residuals	116	107.8	0.9		



Effets principaux



Effet simple

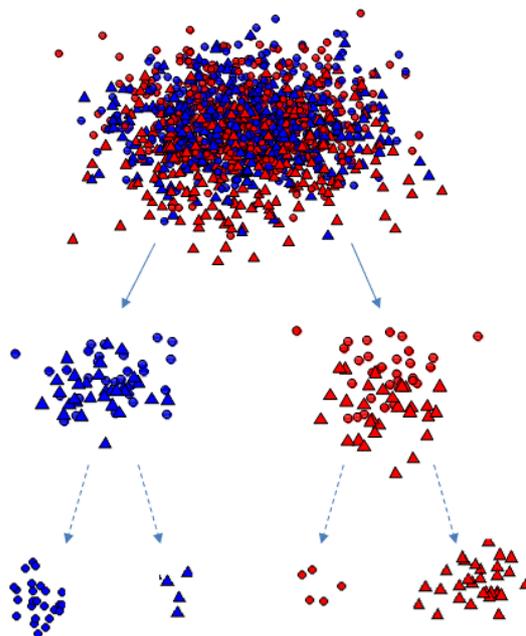


17/11/2015

B A BA de l'ANOVA

29

Facteur fixé vs facteur aléatoire



Je tire 2 échantillons au hasard, un **bleu** et un **rouge**.

Le facteur couleur est un **facteur fixé** : bleu et rouge sont les deux seules modalités qui m'intéressent et je fixe la taille de mes échantillons.

Le facteur forme est un **facteur aléatoire** : je ne sais pas si les deux formes observées sont les seules dans ma population et je ne peux pas fixer le nombre de ronds ni de triangles.

17/11/2015

B A BA de l'ANOVA

30

Mesures répétées

...

Facteurs emboîtés vs croisés

...

Are ANOVA and linear regression twin princesses grown in different castles ?

- L'ANOVA c'est pour les variables explicatives qualitatives, la régression linéaire c'est pour les variables explicatives quantitatives.
- L'ANOVA met l'accent sur la comparaison (test d'hypothèse), le modèle linéaire met l'accent sur la prédiction (modèles).
- L'ANOVA est un cas particulier de la régression linéaire généralisée.
- Tout ce qu'on peut faire avec l'ANOVA, on peut le faire avec le modèle linéaire généralisé... peut-être pas tout à fait tout d'après Gelman.
- Mais moi, j'aime bien l'ANOVA.

Merci de votre attention

