

# Statistiques pour les sciences (MAT-4681)

Arthur Charpentier

# 13 - Loi multinomiale et tableaux croisés

été 2022

# Un peu de formalisme...

## Tableau de comptage

$X$  peut prendre les modalités  $\{x_1, \dots, x_J\}$ . On appelle **tableau de comptage** le vecteur  $\mathbf{n}$  de taille  $J$   $\mathbf{n} = [n_j] = (n_1, \dots, n_J)$  où  $n_j$  est le nombre d'individus dont la modalité est  $x_j$ .

**Example** Considérons l'exemple où  $X$  désigne la couleur des yeux, de la base `HairEyeColor`,

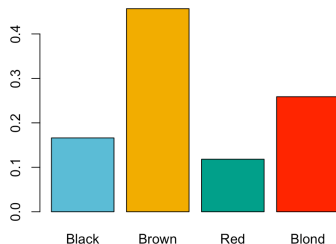
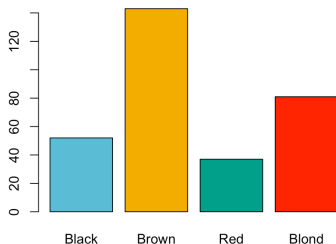
```
1 > data(HairEyeColor)
2 > n = apply(HairEyeColor[, , Sex="Female"], 1, sum)
3 > n
4 Black Brown Red Blond
5 52 143 37 81
```

# Un peu de formalisme...

Si  $n$  est l'effectif total,  $n = \sum_{i=1}^n \mathbf{1}_{j \neq i}$

La fréquence est  $\mathbf{f} = \frac{1}{n} \mathbf{n} = \left[ \frac{n_j}{n} \right]$

```
1 > barplot(n)
2 > f = n/sum(n)
3 > barplot(f)
```



# Loi multinomiale

On suppose que  $\{X_1, \dots, X_n\}$  est une collection de variables catégorielles indépendantes, de loi  $\mathbf{p} = (p_1, \dots, p_J)$

La variable  $Y_{j:i} = \mathbf{1}_j(X_i)$  suit une loi de Bernoulli  $\mathcal{B}(p_j)$ , où

$$p_j = \mathbb{E}[Y_j] = \mathbb{E}(\mathbf{1}_j(X)) = \mathbb{P}[X = x_j]$$

La variable  $N_j = \sum_{i=1}^n \mathbf{1}_j(X_i) = \sum_{i=1}^n Y_{j:i}$  suit une loi binomiale  $\mathcal{B}(n, p_j)$

Espérance, variance et covariance

$$\mathbb{E}(N_i) = np_i \quad \text{var}(N_i) = np_i(1 - p_i)$$

$$\text{cov}(N_i, N_j) = -np_i p_j$$

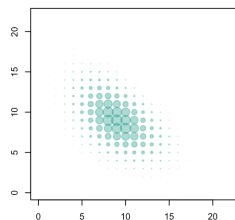
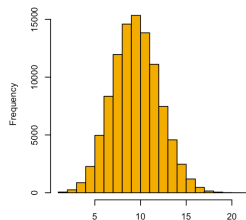
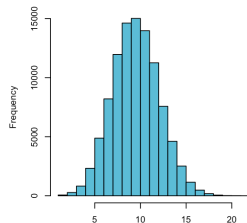
(admis)

# Loi multinomiale

Les variables  $N_i$  et  $N_j$  ne sont pas indépendantes, car  $\sum_{j=1}^J N_j = n$

On peut simuler une loi prenant les valeurs  $\{1, 2, 3\}$ , uniforme ( $\mathbf{p} = (1/3, 1/3, 1/3)$ ),  $n = 30$  fois.

```
1 > X = sample(1:3, size=n, replace=TRUE)
2 > N = table(X)[as.character(1:3)]
```



# Loi multinomiale

On peut montrer que

## Loi multinomiale

$$\mathbb{P}(N_1 = n_1, \dots, N_J = n_J) = \frac{n!}{n_1! \dots n_J!} p_1^{n_1} \dots p_J^{n_J}$$

pour tout  $\mathbf{n} = (n_1, \dots, n_J)$  tel que  $n_1 + \dots + n_J = n$ .

En particulier si  $J = 2$ , on retrouve la loi binomiale,

$$\mathbb{P}(N_1 = n_1, N_2 = n_2) = \frac{n!}{n_1! n_2!} p_1^{n_1} p_2^{n_2}$$

pour  $n_1$  et  $n_2$  tels que  $n_1 + n_2 = n$ , ou

$$\mathbb{P}(N_1 = n_1, N_2 = n - n_1) = \frac{n!}{n_1! (n - n_1)!} p_1^{n_1} (1 - p_1)^{n - n_1}$$

# Loi multinomiale

On peut montrer que

## Loi multinomiale, approximation

Si  $\{x_1, \dots, x_n\}$  est une collection de variables catégorielles indépendantes, de probabilités  $\mathbf{p} = (p_1, \dots, p_J)$ , et si  $n_j$  est le nombre d'observations de la modalité  $j$ ,

$$\frac{N_j - np_j}{\sqrt{np_j(1 - p_j)}} \approx \mathcal{N}(0, 1)$$

et

$$\sum_{j=1}^J \frac{(N_j - np_j)^2}{np_j} \approx \chi^2(J - 1)$$

(le résultat sera admis ici)

Cette dernière propriété permet de proposer un test de fréquence

Loi multinomiale, test  $H_0 : \mathbf{p} = \mathbf{p}_0$  contre  $H_1 : \mathbf{p} \neq \mathbf{p}_0$

Si  $\{x_1, \dots, x_n\}$  est une collection de variables catégorielles indépendantes, de probabilités  $\mathbf{p} = (p_1, \dots, p_J)$ , pour tester  $H_0 : \mathbf{p} = \mathbf{p}_0$  contre  $H_1 : \mathbf{p} \neq \mathbf{p}_0$  la statistique de test est

$$Q = \sum_{j=1}^J \frac{(N\hat{p}_j - Np_{0,j})^2}{Np_{0,j}} = \sum_{j=1}^J \frac{(n_j - Np_{0,j})^2}{Np_{0,j}}$$

Si  $H_0 : \mathbf{p} = \mathbf{p}_0$  est vraie,  $Q \sim \chi^2(J-1)$ . Et donc

► on rejette  $H_0$  si  $q > Q_{J-1}^{-1}(1 - \alpha)$

où  $Q_\nu$  est la fonction de répartition de la loi du chi-deux,  $\chi^2(\nu)$ .



## Loi multinomiale, test

On a lancé  $n = 600$  fois un dé, est-il biaisé ?

```
1 > table(X1)
2      1      2      3      4      5      6
3     88    109    107     94    105     97
```

$$q = \frac{(88 - 100)^2}{100} + \frac{(109 - 100)^2}{100} + \frac{(107 - 100)^2}{100} + \frac{(94 - 100)^2}{100} + \frac{(105 - 100)^2}{100} + \frac{(97 - 100)^2}{100}$$

```
1 > sum((table(X1)-100)^2/100)
2 [1] 3.44
```

or le quantile à 95% d'une loi  $\chi^2(6 - 1)$  est 11.07

```
1 > qchisq(.95,6-1)
2 [1] 11.0705
```

et la  $p$ -value vaut 36.7%

```
1 > 1-pchisq(3.44,6-1)
2 [1] 0.6324852
```

## Loi multinomiale, test

On a lancé  $n = 600$  fois un (autre) dé, est-il biaisé ?

```
1 > table(X2)
2      1      2      3      4      5      6
3     89    131     93     92    104     91
```

$$q = \frac{(89 - 100)^2}{100} + \frac{(131 - 100)^2}{100} + \frac{(93 - 100)^2}{100} + \frac{(92 - 100)^2}{100} + \frac{(104 - 100)^2}{100} + \frac{(91 - 100)^2}{100}$$

```
1 > sum((table(X2) - 100)^2 / 100)
2 [1] 12.92
```

qui dépasse le quantile à 95% d'une loi  $\chi^2(6 - 1)$  est 11.07

```
1 > qchisq(.95, 6 - 1)
2 [1] 11.0705
```

et la  $p$ -value vaut 36.7%

```
1 > 1 - pchisq(12.92, 6 - 1)
2 [1] 0.0241401
```

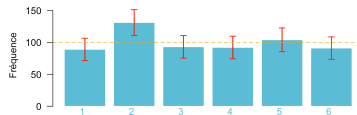
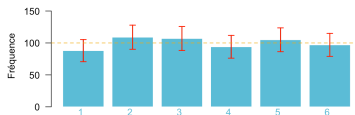
# Tests multiples ★★★

On a ponctuellement des intervalles de confiance, sur les probabilités

$$\left[ \hat{p}_j \pm u_{1-\alpha} \cdot \sqrt{\frac{\hat{p}_j(1 - \hat{p}_j)}{n}} \right] \text{ où } \hat{p}_j = \frac{n_j}{n},$$

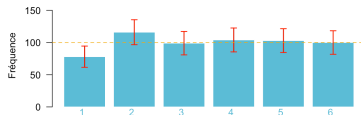
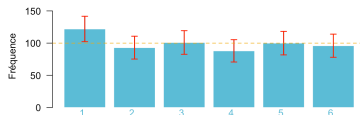
et sur les fréquences

$$\left[ n_j \pm u_{1-\alpha} \cdot \sqrt{\frac{n_j(n - n_j)}{n}} \right]$$



# Tests multiples ★★★

Ces intervalles de confiance sont associés à 6 tests simples

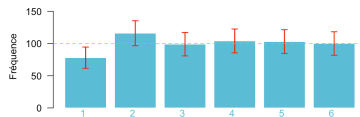
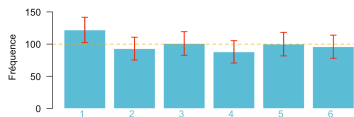


On peut refuser un test simple (un sur les six)

```
1 > table(X)
2   1   2   3   4   5   6
3  78 116  99 104 103 100
4 > prop.test(table(X)[1],600,1/6)
5
6 1-sample proportions test with continuity correction
7
8 data:  table(X)[1] out of 600, null probability 1/6
9 X-squared = 5.547, df = 1, p-value = 0.01851
10 alternative hypothesis: true p is not equal to
    0.1666667
11 95 percent confidence interval:
12  0.1046716 0.1601808
```

# Tests multiples ★★★

et on peut accepter le test multiple ( $p$ -value de 17.6%)



Car ici, on regarde un test multiple,  $H_0 : p_1 = \dots = p_6$ .

```
1 > 1-pchisq(sum((table(X)-100)^2/100),6-1)
2 [1] 0.175996
```

i.e.

```
1 > chisq.test(table(X), p = rep(1/6,6))
2
3   Chi-squared test for given probabilities
4
5 data:  table(X)
6 X-squared = 7.66, df = 5, p-value = 0.176
```

## Un peu de formalisme...

La formule de base repose sur

$$Z_j = \frac{\text{comptage observé} - \text{comptage attendu sous } H_0}{\sqrt{\text{comptage attendu}}} = \frac{O_j - E_j}{\sqrt{E_j}}$$

Si on a assez d'observations,  $Z_j \approx \mathcal{N}(0, 1)$  et

$$Q = Z_1^2 + Z_2^2 + \dots + Z_{J-1}^2 + Z_J^2 \approx \chi^2(J-1)$$

que l'on notera aussi

$$Q = \sum_{j=1}^J \frac{(O_j - E_j)^2}{E_j} \approx \chi^2(J-1)$$

## Un peu de formalisme...

De manière générale, la statistique de test est

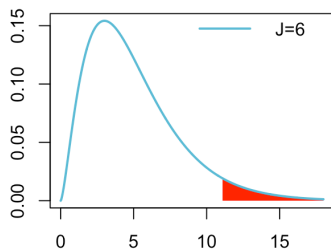
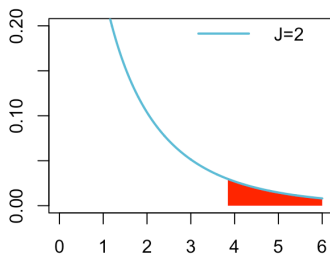
$$q = \sum_{j=1}^J \frac{(N\hat{p}_{0,j} - Np_{0,j})^2}{Np_{0,j}} = \sum_{j=1}^J \frac{(n_j - Np_{0,j})^2}{Np_{0,j}}$$

la  $p$ -value est

$$p = \mathbb{P}[Q > q | Q \sim \chi^2(J-1)]$$

mais on peut passer par la région de rejet

- ▶ si  $q > Q_{J-1}^{-1}(1 - \alpha)$  on rejette  $H_0$
- ▶ si  $q < Q_{J-1}^{-1}(1 - \alpha)$  on ne rejette pas  $H_0$



# Un peu de formalisme...

## Tableau de contingence

$X$  peut prendre les modalités  $\{x_1, \dots, x_I\}$  et  $Y$  les modalités  $\{y_1, \dots, y_J\}$ . On appelle **tableau de contingence** la matrice  $N$ ,  $I \times J$ ,  $N = [n_{ij}]$  où  $n_{ij}$  est le nombre d'individus dont les modalités sont  $x_i$  et  $y_j$ . On parle parfois aussi de **tri-croisé**.

**Example** Considérons l'exemple où  $X$  désigne la couleur des cheveux, et  $Y$  la couleur des yeux, de la base HairEyeColor,

```
1 > data(HairEyeColor)
2 > HairEyeColor[, , Sex="Female"]
3      Eye
4 Hair   Brown Blue Hazel Green
5  Black   36    9     5     2
6  Brown   66   34    29    14
7   Red   16    7     7     7
8  Blond    4   64     5     8
```



# Un peu de formalisme...

## Effets marginaux

Les **effets marginaux** sont notés

$$n_{i,\cdot} = \sum_j n_{i,j} \text{ et } n_{\cdot,j} = \sum_i n_{i,j}$$

L'effectif total de la population est alors

$$n = \sum_i n_{i,\cdot} = \sum_j n_{\cdot,j} = \sum_{i,j} n_{i,j}$$

```
1 > apply(HairEyeColor[, , Sex="Female"], 2, sum)
2 Brown   Blue  Hazel  Green
3    122    114     46    31
4 > apply(HairEyeColor[, , Sex="Female"], 1, sum)
5 Black Brown   Red  Blond
6     52    143    37    81
```

## Un peu de formalisme...

On pose alors  $F = \frac{1}{n}N = [f_{i,j}]$ , où  $f_{i,j} = \frac{n_{i,j}}{n}$ .

```
1 > HairEyeColor[, , Sex="Female"] / sum(HairEyeColor[, , Sex  
    ="Female"])  
2      Eye  
3 Hair      Brown      Blue      Hazel      Green  
4   Black 0.11501597 0.02875399 0.01597444 0.006389776  
5   Brown 0.21086262 0.10862620 0.09265176 0.044728435  
6    Red  0.05111821 0.02236422 0.02236422 0.022364217  
7   Blond 0.01277955 0.20447284 0.01597444 0.025559105
```

De la même manière, on peut définir les effets marginaux

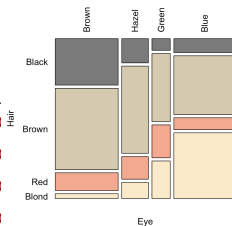
$$f_{i,\cdot} = \sum_j f_{i,j} \text{ et } f_{\cdot,j} = \sum_i f_{i,j}$$

# Probabilités conditionnelles

	brown	hazel	green	blue	
black	63.0%	13.9%	4.6%	18.5%	100.0%
brown	41.6%	18.9%	10.1%	29.4%	100.0%
red	36.6%	19.7%	19.7%	23.9%	100.0%
blond	5.5%	7.9%	12.6%	74.0%	100.0%
	37.2%	15.7%	10.8%	36.3%	



	brown	hazel	green	blue	
black	30.9%	16.1%	7.8%	9.3%	18.2%
brown	54.1%	58.1%	45.3%	39.1%	48.3%
red	11.8%	15.1%	21.9%	7.9%	12.0%
blond	3.2%	10.8%	25.0%	43.7%	21.5%
	100.0%	100.0%	100.0%	100.0%	



# Test d'indépendance I

## Indépendance $X \perp\!\!\!\perp Y$

Soit  $X$  et  $Y$  deux variables discrètes,  $X$  et  $Y$  sont indépendantes - noté  $X \perp\!\!\!\perp Y$  - si

$$\mathbb{P}[X = x, Y = y] = \mathbb{P}[X = x] \cdot \mathbb{P}[Y = y], \quad \forall x, y.$$

Compte tenu des notations précédentes,

- ▶ on estime  $\mathbb{P}[X = x_i, Y = y_j]$  par  $\hat{p}_{i,j} = f_{i,j} = \frac{n_{i,j}}{n}$
- ▶ on estime  $\mathbb{P}[X = x_i]$  par  $\hat{p}_{i,\cdot} = f_{i,\cdot} = \frac{n_{i,\cdot}}{n}$
- ▶ on estime  $\mathbb{P}[Y = y_j]$  par  $\hat{p}_{\cdot,j} = f_{\cdot,j} = \frac{n_{\cdot,j}}{n}$

# Test d'indépendance II

## Indépendance empirique $\mathbf{x} \perp\!\!\!\perp \mathbf{y}$

Soit  $(\mathbf{x}, \mathbf{y}) = \{(x_1, y_1), \dots, (x_n, y_n)\}$  des couples de variables catégorielles appariées. Si  $[n_{i,j}]$  est le tableau de contingence associé, on dira que  $\mathbf{x}$  et  $\mathbf{y}$  sont empiriquement indépendants - noté  $\mathbf{x} \perp\!\!\!\perp \mathbf{y}$  - si

$$\hat{p}_{i,j} = \hat{p}_{i,\cdot} \hat{p}_{\cdot,j} \text{ ou } \frac{n_{i,j}}{n} = \frac{n_{i,\cdot}}{n} \frac{n_{\cdot,j}}{n}, \quad \forall i \text{ et } j$$

$$\text{ou } n_{i,j} = \frac{n_{i,\cdot} n_{\cdot,j}}{n}, \quad \forall i \text{ et } j$$

On pourra noter  $\hat{p}_{i,j}^\perp = \hat{p}_{i,\cdot} \hat{p}_{\cdot,j}$ , et on aura indépendance si  $\mathbf{p} = \mathbf{p}^\perp$

Un test naturel sera un test du chi-deux.

## Test d'indépendance III

Test  $H_0 : X \perp\!\!\!\perp Y$  contre  $H_1 : X \not\perp\!\!\!\perp Y$

Soit  $(\mathbf{x}, \mathbf{y}) = \{(x_1, y_1), \dots, (x_n, y_n)\}$  des couples de variables catégorielles appariées. Si  $[n_{i,j}]$  est le tableau de contingence associé, pour tester  $H_0 : X \perp\!\!\!\perp Y$  contre  $H_1 : X \not\perp\!\!\!\perp Y$  la statistique de test est

$$Q = \sum_{i=1}^I \sum_{j=1}^J \frac{(n\hat{p}_{i,j} - n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j})^2}{n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}} = \sum_{i=1}^I \sum_{j=1}^J \frac{(n_{i,j} - n_{i,\cdot}n_{\cdot,j}/n)^2}{n_{i,\cdot}n_{\cdot,j}/n}$$

Si  $H_0 : X \perp\!\!\!\perp Y$  est vraie,  $Q \sim \chi^2((I-1)(J-1))$ . Et donc

► on rejette  $H_0$  si  $q > Q_{(I-1)(J-1)}^{-1}(1 - \alpha)$

où  $Q_\nu$  est la fonction de répartition de la loi du chi-deux,  $\chi^2(\nu)$ .

# Test d'indépendance IV

Notons qu'on peut aussi écrire la statistique de test sur les probabilités, et pas les comptages

$$Q = \sum_{i=1}^I \sum_{j=1}^J \frac{(n\hat{p}_{i,j} - n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j})^2}{n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}} = n \sum_{i=1}^I \sum_{j=1}^J \frac{(\hat{p}_{i,j} - \hat{p}_{i,\cdot}\hat{p}_{\cdot,j})^2}{\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}}$$

On peut également écrire

$$Q = \sum_{i=1}^I \sum_{j=1}^J \epsilon_{i,j}^2 \text{ où } \epsilon_{i,j} = \frac{(n\hat{p}_{i,j} - n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j})}{\sqrt{n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}}}$$

où, si  $H_0 : X \perp\!\!\!\perp Y$  est vraie,  $\epsilon_{i,j} \approx \mathcal{N}(0, 1)$ .

$\epsilon_{i,j}^2$  est appelée **contribution au test du chi-deux**

# Test d'indépendance V

```
1 > N = HairEyeColor[, , Sex="Female"] + HairEyeColor[, , Sex
  = "Male"]
2 > (Q = chisq.test(N))
3
4 Pearson's Chi-squared test
5
6 data:  N
7 X-squared = 138.29, df = 9, p-value < 2.2e-16
8 > Q$observed
9      Eye
10 Hair   Brown Blue Hazel Green
11 Black    68   20    15     5
12 Brown   119   84    54    29
13 Red     26   17    14    14
14 Blond     7   94    10    16
```



# Test d'indépendance VI

```
1 > Q$expected
2      Eye
3 Hair      Brown      Blue      Hazel      Green
4  Black  40.13514  39.22297  16.96622  11.675676
5  Brown 106.28378 103.86824  44.92905  30.918919
6  Red    26.38514  25.78547  11.15372   7.675676
7  Blond  47.19595  46.12331  19.95101  13.729730
```

Comme attendu, on notera que  $n_{\cdot j}^{\perp} = n_{\cdot j}$  pour tout  $j$

```
1 > apply(Q$observed, 2, sum)
2 Brown  Blue  Hazel  Green
3   220   215    93    64
4 > apply(Q$expected, 2, sum)
5 Brown  Blue  Hazel  Green
6   220   215    93    64
```

(et on vérifiera que  $n_{\cdot j}^{\perp} = n_{\cdot j}$  pour tout  $j$ )

# Test d'indépendance VII

```
1 > Q
2
3   Pearson's Chi-squared test
4
5 data:  N
6 X-squared = 138.29, df = 9, p-value < 2.2e-16
```

On rejette ici  $H_0 : X \perp\!\!\!\perp Y$  car  $q$  dépasse le quantile à 95% d'une loi du  $\chi^2(3 \times 3)$ ,

```
1 > qchisq(.95, 3*3)
2 [1] 16.91898
```

avec une  $p$ -value inférieure à  $10^{-16}$ .

On peut aussi calculer les **résidus**  $\epsilon_{i,j} = \frac{(n\hat{p}_{i,j} - n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j})}{\sqrt{n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}}}$

# Test d'indépendance VIII

$$\epsilon_{i,j} = \frac{(n\hat{p}_{i,j} - n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j})}{\sqrt{n\hat{p}_{i,\cdot}\hat{p}_{\cdot,j}}}$$

```
1 > Q$residuals
2      Eye
3 Hair      Brown      Blue      Hazel      Green
4  Black  4.3984 -3.0694 -0.4774 -1.9537
5  Brown  1.2335 -1.9495  1.3533 -0.3451
6   Red   -0.0750 -1.7301  0.8523  2.2827
7  Blond -5.8510  7.0496 -2.2278  0.6127
```

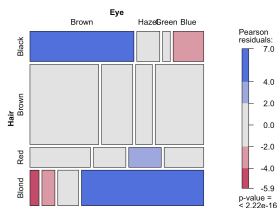
On rejette  $H_0$  car on a

- ▶ trop de personnes aux cheveux Black ayant les yeux Brown
- ▶ trop de personnes aux cheveux Blond ayant les yeux Blue
- ▶ pas assez de personnes Black ayant les yeux Blue
- ▶ pas assez de personnes Blond ayant les yeux Brown

# Test d'indépendance IX

	brown	hazel	green	blue	
black	68	15	5	20	108
brown	119	54	29	84	286
red	26	14	14	17	71
blond	7	10	16	94	127
	220	93	64	215	

	brown	hazel	green	blue	
black	40	17	12	39	108
brown	106	45	31	104	286
red	26	11	8	26	71
blond	47	20	14	46	127
	220	93	64	215	



on compare  $n_{i,j}$  et  $n_{i,j}^\perp$

$$n_{i,j}^\perp = \frac{n_{i,\cdot} \cdot n_{\cdot,j}}{n}$$