

# Introduction à la théorie des sondages - Cours 1

Gaspar Massiot  
gaspar.massiot@ined.fr



2024-2025

# Introduction

- Note
  - Devoir Maison (TD) Coefficient 1
  - Devoir sur table (CM) Coefficient 2 : Durée 1h, **Calculatrice et 1 Feuille RV autorisée**
- Intervenants
  - Cours : Gaspar Massiot ([gaspar.massiot@ined.fr](mailto:gaspar.massiot@ined.fr))
  - TD : Ulysse Lebec ([ulebec@mediametrie.fr](mailto:ulebec@mediametrie.fr)) et Nicolas Salaün ([nsalaun@mediametrie.fr](mailto:nsalaun@mediametrie.fr))
  - TP : Tony Bissonnier ([tbissonnier@mediametrie.fr](mailto:tbissonnier@mediametrie.fr))
- Supports de cours inspirés de ceux de Thomas Merly-Alpa, Paul Cochet (Ined), Antoine Rebecq et Martin Chevalier (Insee)

# Sommaire I

## 1 Pourquoi le sondage ?

- Concept
- Utilisations
- Un échantillon “représentatif” ?
- Pondération

## 2 Notion d'estimateur

- Vocabulaire
- Retour sur l'estimateur naïf
- Les probabilités d'inclusion

# Chapitre 1

## Pourquoi le sondage ?

## Partie 1

### Concept

# Concept

Qu'est-ce que l'échantillonnage / l'estimation par sondage ?

- Une population de grande taille
- Compter ou interroger est coûteux
- On sélectionne quelques individus qui répondent " pour tout le monde"

Idee cruciale : sélectionner **aléatoirement** ces individus.

# Historique

Historiquement et conceptuellement, rien d'évident !

- Laplace (1785) : recensement par une sous-partie de la population

# Historique

Historiquement et conceptuellement, rien d'évident !

- Laplace (1785) : recensement par une sous-partie de la population
- Kiaer (1895) : échantillon "représentatif"



# Historique

Historiquement et conceptuellement, rien d'évident !

- Laplace (1785) : recensement par une sous-partie de la population
- Kiaer (1895) : échantillon "représentatif"  
... puis 1925 : acceptation de l'échantillonnage aléatoire

# Historique

Historiquement et conceptuellement, rien d'évident !

- Laplace (1785) : recensement par une sous-partie de la population
- Kiaer (1895) : échantillon "représentatif"  
... puis 1925 : acceptation de l'échantillonnage aléatoire
- Gallup (1936) : élections américaines

# Élections américaines de 1936

- Duel entre Alfred Landon (Républicain) et Franklin Roosevelt (Démocrate)
- Un magazine interroge ses 2 millions de lectorices : victoire de Landon
- Gallup fait un sondage sur 50 000 personnes : il prédit la victoire de Roosevelt

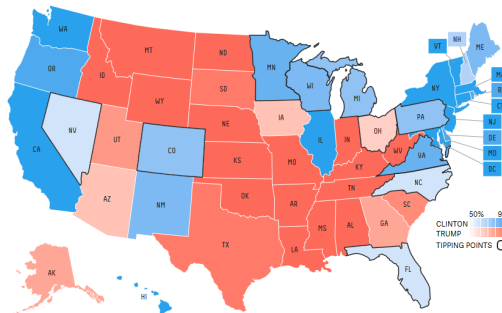
The map shows the following electoral college distribution:

State	Electoral College Votes	Winner
Alaska	3	Roosevelt
Alabama	9	Roosevelt
Arizona	3	Roosevelt
Arkansas	3	Roosevelt
California	22	Roosevelt
Colorado	4	Roosevelt
Connecticut	8	Roosevelt
Delaware	3	Roosevelt
District of Columbia	3	Roosevelt
Florida	7	Roosevelt
Georgia	8	Roosevelt
Idaho	4	Roosevelt
Illinois	12	Roosevelt
Indiana	11	Roosevelt
Iowa	7	Roosevelt
Kansas	6	Roosevelt
Kentucky	12	Roosevelt
Louisiana	9	Roosevelt
Maine	4	Roosevelt
Massachusetts	17	Roosevelt
Michigan	19	Roosevelt
Minnesota	11	Roosevelt
Mississippi	9	Roosevelt
Missouri	11	Roosevelt
Montana	4	Roosevelt
Nebraska	4	Roosevelt
Nevada	3	Roosevelt
New Hampshire	4	Roosevelt
New Jersey	16	Roosevelt
New Mexico	4	Roosevelt
New York	47	Roosevelt
North Carolina	12	Roosevelt
North Dakota	4	Roosevelt
Ohio	26	Roosevelt
Oklahoma	9	Roosevelt
Oregon	5	Roosevelt
Pennsylvania	36	Roosevelt
Rhode Island	4	Roosevelt
South Carolina	8	Roosevelt
South Dakota	4	Roosevelt
Tennessee	11	Roosevelt
Texas	23	Roosevelt
Vermont	3	Roosevelt
Virginia	11	Roosevelt
Washington	8	Roosevelt
West Virginia	4	Roosevelt
Wisconsin	11	Roosevelt
Wyoming	4	Roosevelt
New England	5	Landon

# Élections américaines de 2016

Nate Silver, <http://fivethirtyeight.com> :

Chance of winning



# Élections américaines de 2016

Élection de Donald Trump, Brexit... Pourquoi les sondages ont-ils eu tout faux ?

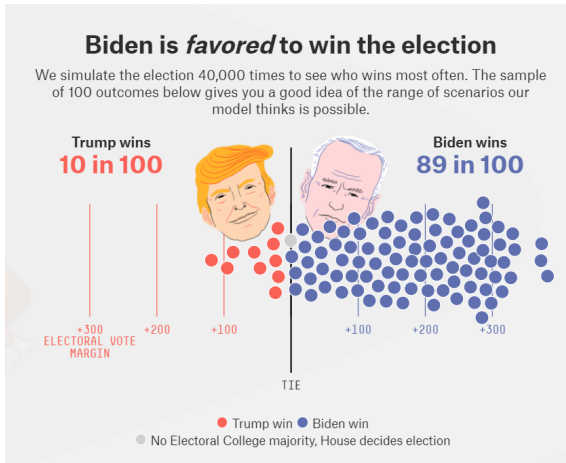
# Élections américaines de 2016

Élection de Donald Trump, Brexit... Pourquoi les sondages ont-ils eu tout faux ?

- Marge d'erreur et précision : 1 000 personnes ?
- Temporalité : changement d'avis (exemple : hausse du vote Fillon à la primaire LR 2017)
- Mensonge ou camouflage des intentions (exemple : traitement du vote FN depuis 2002)
- Ne souhaitent pas répondre : à suivre dans le cours
- Enquêtes Internet ou par téléphone : meilleures méthodes ?

# Élections américaines de 2020

En 2020, une meilleure prédiction :

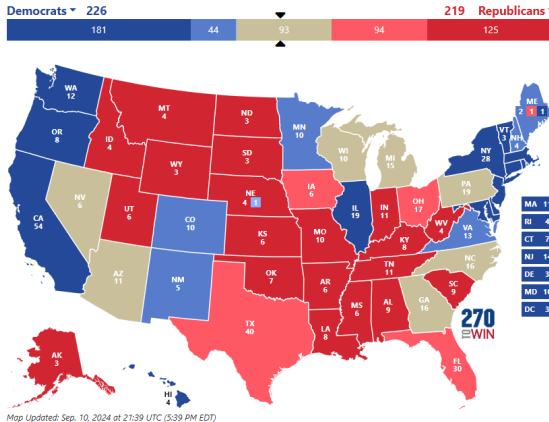




Pourquoi le sondage ?  
Notion d'estimateur

Concept  
Utilisations  
Pourquoi faire une enquête ?  
Un échantillon "représentatif" ?  
Pondération

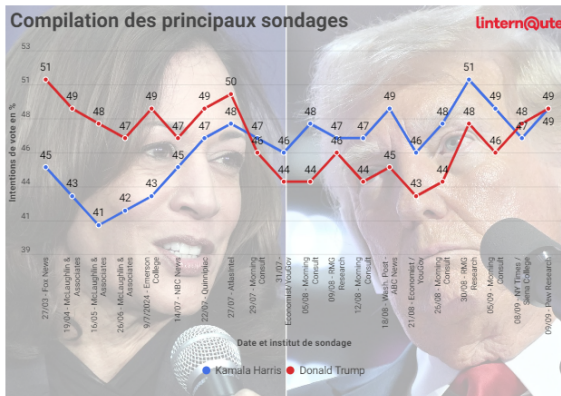
# Et pour 2024 ?



Pourquoi le sondage ?  
Notion d'estimateur

Concept  
Utilisations  
Pourquoi faire une enquête ?  
Un échantillon "représentatif" ?  
Pondération

Et pour 2024 ?



## Partie 2

### Utilisations

# Statistique publique

- Enquêtes auprès des ménages : le moral des ménages, le taux de chômage
- Enquêtes auprès des entreprises - ESA (Enquête Sectorielle Annuelle) : Chiffre d'affaire par secteur, chiffres d'investissement, ...

# Statistique publique

Et d'autres sujets :

- Epicov : Enquête rapide sur le covid-19 pendant les confinements (Insee-Inserm-Drees) ;
- Panel ELIPSS : Panel de sciences sociales (Sciences Po) ;
- EMP : Enquête Mobilité des Personnes (INSEE-SDES) ;
- Familles et Employeurs (Ined) ...

## Autres exemples

- Biologie : dénombrement d'espèces

## Autres exemples

- Biologie : dénombrement d'espèces
- Politique

## Autres exemples

- Biologie : dénombrement d'espèces
- Politique
- Marketing



## Autres exemples

- Biologie : dénombrement d'espèces
- Politique
- Marketing



# Pour aller plus loin sur l'utilisation des sondages et leurs limites

BLAST, Sondages d'opinion : l'Overdose



## Partie 3

### Pourquoi faire une enquête ?

# Conception

Une enquête peut être coûteuse (en budget - 2 millions pour une enquête INSEE, mais aussi en temps des enquêté.e.s). Il faut donc s'assurer que le sujet est :

- Pertinent (contraintes européennes, demandes d'études, sujet actuel)
- Non couvert (autres enquêtes, autres données)
- Réalisable (pas trop complexe, légalité, anonymisation)

# Questionnaire

Une fois les objectifs identifiés, il faut réaliser un questionnaire :

- Qui colle aux concepts
- Mais compréhensible par l'enquêté : ni équivoque, ni flou
- Qui permette de la comparabilité avec d'autres sources

⇒ Une étape cruciale mais difficile !

## Partie 4

### Un échantillon "représentatif" ?

# Échantillon représentatif

Un "échantillon représentatif" : "Village" de 100 habitants

- Est-ce que le concept d'échantillon représentatif est toujours pertinent ?
- Si on veut connaître le secteur automobile en France, quelle est la bonne stratégie ?

"Sondage" devrait toujours aller de pair avec **"objectif"** (même si les objectifs pour un même échantillon peuvent être nombreux).

## Secteur automobile

Quelle est le chiffre d'affaires moyen d'une entreprise du secteur automobile ?

- On a intérêt à bien interroger Renault et Peugeot.



## Secteur automobile

Quelle est le chiffre d'affaires moyen d'une entreprise du secteur automobile ?

- On a intérêt à bien interroger Renault et Peugeot.
- On doit aussi interroger au hasard des garages.

## Secteur automobile

Quelle est le chiffre d'affaires moyen d'une entreprise du secteur automobile ?

- On a intérêt à bien interroger Renault et Peugeot.
- On doit aussi interroger au hasard des garages.
- Ce n'est pas utile d'interroger trop de garages, car ils se ressemblent.

## Secteur automobile

Quelle est le chiffre d'affaires moyen d'une entreprise du secteur automobile ?

- On a intérêt à bien interroger Renault et Peugeot.
- On doit aussi interroger au hasard des garages.
- Ce n'est pas utile d'interroger trop de garages, car ils se ressemblent.

Renault			50 Md€
Peugeot			40 Md€
Garage 1			300 k€
Garage 2			200 k€

# L'estimation naïve

Pour l'estimation du total et de la moyenne d'une variable  $Y$ , l'estimateur « naïf » est :

- Pour le total, la somme des valeurs  $Y$  des individus de l'échantillon.
- Pour la moyenne, la moyenne des valeurs  $Y$  des individus de l'échantillon.

En général, l'estimation naïve est fausse (*biaisée*), surtout quand l'échantillon est choisi de façon complexe.

## Secteur automobile

Quelle est le chiffre d'affaires moyen d'une entreprise du secteur automobile ?

Renault			50 Md€
Peugeot			40 Md€
Garage 1			300 k€
Garage 2			200 k€

Estimateur naïf :  $(50 + 40 \text{ Md} + 300 + 200 \text{ k}) / 4 \approx 22 \text{ Md } \text{€}$

## Partie 5

### Pondération

## Pondérer ?

Pour éviter d'utiliser l'estimateur naïf, on utilise généralement ce qu'on appelle des poids, qu'on note  $w$  (pour *weight* en anglais).

Le poids d'un individu correspond au nombre d'individus que l'individu de l'échantillon représente dans la population. Si l'on interroge 1 individu sur 100, le poids est alors de 100.

L'estimateur pondéré du total est alors la somme des  $w_i y_i$  sur l'échantillon.

## Retour sur l'exemple

Retour sur le secteur automobile. S'il n'y a qu'un Renault et qu'un Peugeot, il existe en fait près de 80 000 garages.  
Les deux garages enquêtés en représentent donc 80 000 : leur poids  $w$  est de :

$$w_i = \frac{80000}{2} = 40000$$



## Retour sur l'exemple

Renault	Dans l'échantillon		50 Md€
Peugeot	Dans l'échantillon		40 Md€
Garage 1	Dans l'échantillon		300 k€
Garage 2	Dans l'échantillon		200 k€
Garage 3			?
...			?
Garage 80 000			?

## Retour sur l'exemple

On introduit la pondération :

Renault	Dans l'échantillon	1	50 Md€
Peugeot	Dans l'échantillon	1	40 Md€
Garage 1	Dans l'échantillon	40 000	300 k€
Garage 2	Dans l'échantillon	40 000	200 k€

Estimateur naïf :  $(50 + 40 \text{ Md} + 300 + 200 \text{ k}) / 4 \approx 22 \text{ Md } \text{€}$

## Retour sur l'exemple

On introduit la pondération :

Renault	Dans l'échantillon	1	50 Md€
Peugeot	Dans l'échantillon	1	40 Md€
Garage 1	Dans l'échantillon	40 000	300 k€
Garage 2	Dans l'échantillon	40 000	200 k€

Estimateur naïf :  $(50 + 40 \text{ Md} + 300 + 200 \text{ k}) / 4 \approx 22 \text{ Md } \text{€}$

Estimateur pondéré :

$$(1*50 + 1*40\text{Md} + 40\,000*300 + 40\,000*200\text{k}) / (1 + 1 + 40\,000 + 40\,000)$$

soit environ 1,4 millions d'€

# À retenir

- On construit notre sondage et donc notre échantillon dans un but précis.
- On utilise les résultats obtenus en se rappelant de notre méthode de sondage, via la pondération.

## Chapitre 2

# Notion d'estimateur

## Partie 1

### Vocabulaire

# Notations - Définitions

- Population  $\mathcal{U} = \{u_1, \dots, u_k, \dots, u_N\}$
- L'individu  $u_k \in \mathcal{U}$  est repéré sans ambiguïté par son identifiant  $k$ .
- Variable d'intérêt  $Y$ , qui prend la valeur  $y_k$  pour l'individu  $k$
- Objectif du sondage : Mesurer  $\Phi(Y)$ , une fonction dépendant de  $Y$ .

# Notations - Définitions

$Y$  peut être

- quantitative (exemple : revenu). Dans ce cas  $\Phi$  peut être le total, la moyenne, etc.
- qualitative, c'est-à-dire prendre un nombre fini de valeurs (exemple : sexe). Dans ce cas,  $\Phi$  peut être la répartition dans la population.



# Notations - Définitions

La **base de sondage** donne les moyens d'identifier et de joindre les unités d'échantillonnage, souvent il s'agit des individus mais cela peut aussi être un *proxy*.

# Notations - Définitions

- Échantillon  $s \subset \mathcal{U}$
- Si  $s = \mathcal{U}$ , recensement
- Chaque individu  $u_k, k \in s$  est interrogé, et on relève  $y_k$
- Les  $y_k, k \in s$  sont utilisés pour construire un **estimateur**  $\hat{\Phi}$  de  $\Phi$

# Plan de sondage - définition

On note  $\mathcal{S}$  l'ensemble des parties de  $\mathcal{U}$ .

Le plan de sondage  $p$  est une loi de probabilité sur  $\mathcal{S}$  i.e. :

$$\forall s \in \mathcal{S}, p(s) \geq 0$$

$$\sum_{s \in \mathcal{S}} p(s) = 1$$

## Plan de sondage - exemple

Soit  $\mathcal{U} = \{1, 2, 3\}$ . On a alors :  
 $\mathcal{S} =$

## Plan de sondage - exemple

Soit  $\mathcal{U} = \{1, 2, 3\}$ . On a alors :

$$\mathcal{S} = \{\{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$$

## Plan de sondage - exemple

Soit  $\mathcal{U} = \{1, 2, 3\}$ . On a alors :

$$\mathcal{S} = \{\{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}\}$$

On peut définir un plan de sondage  $p$  par :

$$p(\{1\}) = 0 \quad p(\{1, 2\}) = \frac{1}{2} \quad p(\{1, 2, 3\}) = 0$$

$$p(\{2\}) = 0 \quad p(\{1, 3\}) = \frac{1}{3}$$

$$p(\{3\}) = 0 \quad p(\{2, 3\}) = \frac{1}{6}$$

# Paramètre d'intérêt

$Y$  est la **variable d'intérêt** et  $\Phi(Y)$  est le **paramètre d'intérêt**.

Attention,  $Y$  n'est **pas aléatoire** !

L'aléatoire repose entièrement sur l'échantillonnage décrit par le plan de sondage  $p$ .

# Estimateur

Une fois l'échantillon  $s$  tiré, on **estime**  $\Phi(Y)$  à l'aide d'une fonction, notée  $\hat{\Phi}$ , qui dépend de l'échantillon.

$\hat{\Phi}$  est appelé un **estimateur** de  $\Phi(Y)$ .



# Espérance

$$\mathbb{E}(\hat{\Phi}) = \sum_s p(s) \cdot \hat{\Phi}(s)$$

C'est la valeur moyenne de  $\hat{\Phi}$  obtenue avec le plan de sondage considéré **sur tous les échantillons possibles**.

# Biais

$$B(\hat{\phi}) = \mathbb{E}(\hat{\phi}) - \phi$$

Si  $B(\hat{\phi}) = 0$ , alors on parle **d'estimateur sans biais**.

## Variance / Précision

$$\text{Var}(\hat{\Phi}) = \sum_s p(s) \cdot \left[ \mathbb{E}(\hat{\Phi}) - \hat{\Phi}(s) \right]^2$$

C'est une mesure de la dispersion des valeurs  $\hat{\Phi}(s)$  autour de leur moyenne.

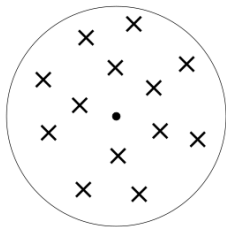
## Variance / Précision

Quantités liées :

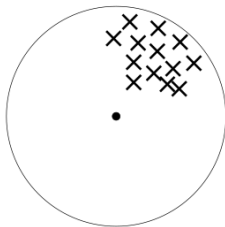
$$\sigma(\hat{\Phi}) = \sqrt{\text{Var}(\hat{\Phi})}, \text{écart-type}$$

$$CV(\hat{\Phi}) = \frac{\sigma(\hat{\Phi})}{\mathbb{E}(\hat{\Phi})}, \text{coefficient de variation}$$

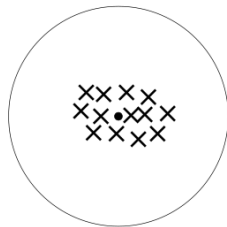
# Schéma



Cas 1



Cas 2



Cas 3

# Erreur quadratique moyenne

$$\begin{aligned}EQM(\hat{\Phi}) &= \sum_s p(s) \cdot [\Phi - \hat{\Phi}(s)]^2 \\ &= \text{Var}(\hat{\Phi}) + B(\hat{\Phi})^2\end{aligned}$$

Entre deux estimateurs sans biais, celui qui a la plus petite variance est de meilleure qualité.

# Construction d'un intervalle de confiance

La **vraie variance**  $\text{Var}(\hat{\Phi})$  n'est pas connue (il faudrait pour cela pouvoir tirer tous les échantillons).

Il faudra donc estimer la variance à partir des données de l'échantillon. L'estimateur sera noté  $\hat{V}(\hat{\Phi})$  ou  $\widehat{\text{Var}}(\hat{\Phi})$ .

# Construction d'un intervalle de confiance

Estimateurs plug-ins des quantités liées à la variance :

$$\hat{\sigma}(\hat{\Phi}) = \sqrt{\widehat{\text{Var}}(\hat{\Phi})}, \text{ écart-type}$$

$$\widehat{CV}(\hat{\Phi}) = \frac{\hat{\sigma}(\hat{\Phi})}{\hat{\Phi}}, \text{ coefficient de variation}$$



## Construction d'un intervalle de confiance

On fait **l'hypothèse** :  $\hat{\Phi} \sim \mathcal{N}(\Phi, \text{Var}(\Phi))$

L'intervalle de confiance à 95% est défini par :

$$IC_{95\%} = \left[ \hat{\Phi} - 2\sigma(\hat{\Phi}); \hat{\Phi} + 2\sigma(\hat{\Phi}) \right]$$

L'intervalle de confiance **estimé** est défini par :

$$\widehat{IC}_{95\%} = \left[ \hat{\Phi} - 2\hat{\sigma}(\hat{\Phi}); \hat{\Phi} + 2\hat{\sigma}(\hat{\Phi}) \right]$$

## Partie 2

### Retour sur l'estimateur naïf

# L'estimateur naïf

Rappel : pour l'estimation du total et de la moyenne d'une variable  $Y$ , l'estimateur « naïf » s'écrit :

$$\hat{T}(Y)_{naif} = \sum_{k \in s} y_k$$
$$\hat{\bar{y}}_{naif} = \frac{1}{n} \sum_{k \in s} y_k$$

# L'estimateur naïf

En général, l'estimation naïve est biaisée :

$$\mathbb{E}(\hat{\phi}_{naif}) = \sum_s p(s) \cdot \hat{\phi}(s) \\ \neq \phi$$

$\mathbb{E}(\hat{\phi})$  est la valeur moyenne de  $\hat{\phi}$  obtenue avec le plan de sondage considéré **sur tous les échantillons possibles**.

## Partie 3

### Les probabilités d'inclusion

## Probabilité d'inclusion $\pi_k$

Pour améliorer l'estimateur naïf, il faut utiliser une pondération.  
On va calculer celle-ci à l'aide des **probabilités d'inclusion**.

La probabilité d'inclusion simple d'un individu  $k$  est la probabilité que cet individu soit dans l'échantillon. Ainsi, pour  $k \in \mathcal{U}$ ,

$$\pi_k = \mathbb{P}(k \in s) = \mathbb{P}(\delta_k = 1) = \sum_{s \ni k} p(s)$$

où  $\delta_k$  est l'indicatrice d'appartenance de  $k$  à  $\mathcal{S}$ , appelée aussi variable de Cornfield.

## Probabilité d'inclusion $\pi_{kl}$

La probabilité d'inclusion double de deux individus  $k$  et  $l$  est la probabilité que ces deux individus soient ensemble dans l'échantillon. Ainsi, pour  $k, l \in \mathcal{U}$ ,

$$\pi_{kl} = \mathbb{P}(k, l \in s) = \mathbb{P}(\delta_k \delta_l = 1) = \sum_{s \ni k, l} p(s)$$

**Attention :** on n'a pas  $\pi_{kl} = \pi_k \pi_l$  en général ! On note par ailleurs  $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$ .

## Probabilités d'inclusion $\pi_k$ et $\pi_{kl}$ - Propriétés

$$\mathbb{E}(\delta_k) = \pi_k$$

$$\mathbb{E}(\delta_k \delta_l) = \pi_{kl}$$

$$\text{Var}(\delta_k) = \pi_k(1 - \pi_k) \quad \text{Cov}(\delta_k \delta_l) = \Delta_{kl}$$



# Probabilités d'inclusion $\pi_k$ et $\pi_{kl}$ - Propriétés

Pour un plan à **taille fixe**  $n$ , on a :

$$\begin{aligned}\sum_{k \in \mathcal{U}} \pi_k &= n \\ \sum_{\substack{k, l \in \mathcal{U} \\ k \neq l}} \pi_{kl} &= n(n-1) \\ \sum_{\substack{l \in \mathcal{U} \\ l \neq k}} \pi_{kl} &= \pi_k(n-1)\end{aligned}$$