

DÉPARTEMENT DE MATHÉMATIQUES ET APPLICATIONS

ENSEIGNEMENT

Introduction à la pensée statistique

Gilles Stoltz (chargé de recherche CNRS)

Paul Baudin (doctorant INRIA)

Table des matières

- Modèle statistique
- Notion d'estimateur
- Construction d'intervalles de confiance pour des lois de Bernoulli
- Cadre des tests d'hypothèses

- Illustrations par des simulations sous Python

- Coupures de presse et bandes dessinées

Passer de l'autre côté du miroir d'Alice :

Du calcul des probabilités au monde merveilleux de la statistique

Gilles Stoltz (chargé de recherche au CNRS)

Paul Baudin (doctorant INRIA)

En probabilités, on part de variables aléatoires X_1, \dots, X_n de loi connue, p.ex. de Bernoulli de paramètre $p_0 \in [0,1]$ connu, et on dit des choses sur des fonctions mesurables de X_1, \dots, X_n . P.ex., l'inégalité de Tchebychev - Markov :

$\forall \varepsilon > 0$,
ie, les valeurs typiques de \bar{X}_n sont proches de p_0 .

$$\mathbb{P}\{ |\bar{X}_n - p_0| > \varepsilon \} \leq \frac{\text{Var}(\bar{X}_n)}{\varepsilon^2} = \frac{p_0(1-p_0)}{n\varepsilon^2};$$

En statistique, on dispose de données x_1, \dots, x_m . Supposons (cf ci-dessous) qu'on puisse les identifier aux réalisations de variables aléatoires X_1, \dots, X_m iid selon une certaine loi, dont on peut connaître la forme mais qui n'est pas entièrement déterminée. P.ex., les X_1, \dots, X_m sont iid selon une certaine loi de Bernoulli de paramètre $p_0 \in [0,1]$ inconnu, à déterminer.

À cause de l'inégalité de Tchebychev - Markov, on se dit que ce p_0 inconnu est sans doute proche de p_0 .

Résumé :

Probabilités : loi connue \rightsquigarrow comportement (d'une fonction) de variables aléatoires

Statistique : comportement de variables aléatoires \rightsquigarrow loi inconnue, à déterminer
(le statisticien est le détective des nombres!)

Exemple : [Les sondages d'opinion]

On interroge au hasard dans la population française un nombre n de sondés à qui l'on pose une question fermée (oui/non). Les réponses sont notées $x_1, \dots, x_n \in \{0,1\}$ (1 si oui, 0 si non).

L'opinion de l'échantillon est résumé par la fréquence $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$.
 Que dire de l'opinion de la population? Une proportion notée p_0 de la population dirait oui: p_0 est inconnue, on ne la connaît que si on interrogeait les 60 millions de Français. (Mais on va voir que p_0 peut être connue typiquement à $\pm 3\%$ rien qu'en interrogeant 1000 Français au hasard!)

Echantillon \bar{x}_n \longleftrightarrow Population p_0
 liens?

Après le tirage: données x_1, \dots, x_n

Avant le tirage: variables aléatoires X_1, \dots, X_n (dont les x_1, \dots, x_n sont les réalisations)

Quitte à effectuer l'approximation d'un tirage sans remise par un tirage avec remise, on peut modéliser les X_1, \dots, X_n comme étant iid selon la loi de Bernoulli de paramètre p_0 .

I. Notion de modèle statistique, d'estimateur.

Comment bien formuler la situation statistique précédente?

→ Notion de modèle statistique: famille de probabilités sur un espace probabilisé $(\Omega, \mathcal{F}, (\mathbb{P}_p)_{p \in \mathcal{I}})$ tel que sous \mathbb{P}_p , X_1, X_2, \dots sont iid $\sim \text{Ber}(p)$

Rq: une telle construction est possible avec $\Omega = \{0,1\}^{\mathbb{N}}$ et le théorème d'extension de Kolmogorov.

Def: Un estimateur T_n de p est toute fonction mesurable de X_1, \dots, X_n .

Ex: $T_n = X_1$, $T_n = 1/2$, $T_n = \bar{X}_n$.

↳ Qu'est-ce qu'un bon estimateur?

Def: T_n est sans biais si $\forall p \in [q]$, $E_p [T_n] = p$

Def: (T_n) est une suite d'estimateurs convergents si
 $\forall p \in [q]$, $T_n \xrightarrow{\mathbb{P}_p\text{-prob.}} p$, i.e.: $\forall \epsilon > 0$,
 $\mathbb{P}_p \{ |T_n - p| > \epsilon \} \rightarrow 0$

Ex:
* $T_n = X_1$ et $T_n = X_n$ sont sans biais
* $(T_n) = (\bar{X}_n)$ est une suite d'estimateurs convergents, cf. loi des grands nombres.

Estimation ponctuelle

Estimation par intervalle

II. Estimation par intervalle.

Def: Soit U_n et V_n deux estimateurs, fondés sur les observations X_1, \dots, X_n (que l'on a construits tels que $U_n \leq V_n$ sur Ω).

$[U_n, V_n]$ est un intervalle de confiance pour p au niveau $1-\alpha$, où $\alpha \in [q]$, si

$$\forall p \in [q], \quad \mathbb{P}_p \{ p \in [U_n, V_n] \} \geq 1-\alpha.$$

Interprétations: - Pour avoir $1-\alpha = 100\%$, il faut prendre $[U_n, V_n] = [q]$.

- $[U_n, V_n]$ a une probabilité au moins $1-\alpha$ de contenir le paramètre de population p .

↳ si on note $[u_n, v_n]$ la réalisation de $[U_n, V_n]$, alors on est confiant (au niveau $1-\alpha$) que p est dans $[u_n, v_n]$... mais en réalité, p y est ou n'y est pas (il n'y a plus rien d'aléatoire).

- Souvent $[U_n, V_n] = [T_n \pm E_n]$ où T_n et E_n sont des estimateurs

↳ E_n est la marge d'incertitude de l'estimateur T_n .

Constructions d'intervalles de confiance:

1/ Par Tchebychev - Markov:

$$\forall p, \forall \varepsilon > 0, \mathbb{P}_p \{ |\bar{X}_n - p| > \varepsilon \} \leq \frac{p(1-p)}{n\varepsilon^2}$$

$$\forall p, \mathbb{P}_p \left\{ |\bar{X}_n - p| \leq \sqrt{\frac{p(1-p)}{n\alpha}} \right\} \geq 1 - \alpha$$

$$= \mathbb{P}_p \left\{ p \in \left[\bar{X}_n \pm \sqrt{\frac{p(1-p)}{n\alpha}} \right] \right\}$$

version probabiliste

en prenant ε tq. $\frac{p(1-p)}{n\varepsilon^2} = \alpha$

version statistique

Hélas: $\left[\bar{X}_n \pm \sqrt{\frac{p(1-p)}{n\alpha}} \right]$ n'est pas un intervalle de confiance, il faudrait que la marge d'incertitude ne dépende pas de p .

$$\text{On pose } E_n = \max_{p \in [0,1]} \sqrt{\frac{p(1-p)}{n\alpha}} = \frac{1}{2\sqrt{n\alpha}}$$

$$\forall p, \mathbb{P}_p \left\{ p \in \left[\bar{X}_n \pm \frac{1}{2\sqrt{n\alpha}} \right] \right\} \geq 1 - \alpha$$

$$\text{Typiquement } 1 - \alpha = 95\%, \quad \frac{1}{2\sqrt{n\alpha}} = \frac{\sqrt{0.25}}{\sqrt{0.05}} \leq 2.24 / \sqrt{n}$$

2/ Par inégalité de Hoeffding:

$$\forall p, \forall \varepsilon > 0, \mathbb{P}_p \{ |\bar{X}_n - p| > \varepsilon \} \leq 2e^{-2n\varepsilon^2}$$

$$\forall p, \mathbb{P}_p \left\{ |\bar{X}_n - p| \leq \sqrt{\frac{1}{2n} \ln \frac{2}{\alpha}} \right\} \geq 1 - \alpha$$

version probabiliste

version statistique

(en prenant ε tq.
 $2e^{-2n\varepsilon^2} = \alpha$)

c'est-à-dire:

$$\forall p, \mathbb{P}_p \left\{ p \in \left[\bar{X}_n \pm \sqrt{\frac{1}{2n} \ln \frac{2}{\alpha}} \right] \right\} \geq 1 - \alpha$$

Avec $1 - \alpha = 95\%$:

$$\sqrt{\frac{1}{2n} \ln \frac{2}{\alpha}} \leq \frac{1.36}{\sqrt{n}}$$

3/ Méthode la plus standard : Théorème limite central + lemme de Slutsky

On admet (cf. simulations sous Python) que

$$\forall p \quad \frac{\sqrt{n}}{\sqrt{\bar{x}_n(1-\bar{x}_n)}} (\bar{X}_n - p) \xrightarrow[\text{sous } \mathbb{P}_p]{d} \mathcal{N}(0,1) \quad \text{la normale standard}$$

de sorte que $\forall p, \quad \mathbb{P}_p \left\{ -1.96 \leq \frac{\sqrt{n}}{\sqrt{\bar{x}_n(1-\bar{x}_n)}} (\bar{X}_n - p) \leq 1.96 \right\}$

$$= \mathbb{P}_p \left\{ p \in \left[\bar{x}_n \pm 1.96 \frac{\sqrt{\bar{x}_n(1-\bar{x}_n)}}{\sqrt{n}} \right] \right\} \xrightarrow{n \rightarrow \infty} \mathbb{P} \left\{ -1.96 \leq \mathcal{N}(0,1) \leq 1.96 \right\} = 95\%$$

$\left[\bar{x}_n \pm 1.96 \frac{\sqrt{\bar{x}_n(1-\bar{x}_n)}}{\sqrt{n}} \right]$ est un intervalle de confiance dit asymptotique (de niveau 95%)

Avantage \rightarrow Marge d'erreur plus faible $1.96 \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}} \leq 1/\sqrt{n}$
(et parfois bien plus petit, quand \bar{x}_n est éloigné de $1/2$)

Inconvénient \rightarrow Caractère asymptotique.

Application \rightarrow Marge d'erreur des sondages en $1/\sqrt{n}$
 \hookrightarrow pour un échantillon de 1000 personnes, précision sur l'estimation à $\pm 3\%$ (seulement)
Les journalistes commencent à intégrer ce fait.

4/ Simulation : Comparaison des trois intervalles de confiance introduits

Pour comparer ds intervalles de confiance, fixer un paramètre p_0 et comparer

- les (espérances des) marges d'erreur $E_{p_0}[E_n]$
- les niveaux $\mathbb{P}_{p_0} \left\{ p_0 \in [\bar{X}_n \pm E_n] \right\} \stackrel{\text{not.}}{=} 1 - \alpha_{n,p}$

Ces quantités sont difficiles (sauf exception) à calculer explicitement : on les estime.

On fixe N , on note $Z = \mathbb{1}_{\{p_0 \in [\bar{x}_n \pm E_n]\}}$
et on simule N copies indépendantes de Z :

$Z_1 \dots Z_N$ dont on fait la moyenne \bar{Z}_N

C'est un estimateur de $1 - \alpha_{n,p_0}$ (avec une marge d'erreur à $\pm 1/\sqrt{N}$).

III. Tests d'hypothèses : Pourquoi deux assertions ^(un peu) contradictoires peuvent simultanément être tenues pour vraies en statistique -

Deux hypothèses (choisies selon le contexte ... par un être humain \rightarrow subjectivité) :

- Une hypothèse de départ H_0 , que l'on conservera sauf si les données la contredisent gravement ;
- Une hypothèse alternative H_1 , vers laquelle on se tournera dans ce dernier cas.

\hookrightarrow Dissymétrie forte entre les rôles des deux hypothèses, H_0 a tendance à être conservée.

Exemple de test pour $H_0 : p_0 = p_{ref}$ contre $H_1 : p_0 \neq p_{ref}$

- Conserve H_0 si $p_{ref} \in [\bar{x}_n \pm E_n]$, intervalle de confiance à 95%
- Rejette H_0 sinon.

Application numérique : $\bar{x}_n = 51\%$ et $n = 1000 \hookrightarrow$ intervalle $[51\% \pm 3\%]$

* Les hypothèses $H_0^{(1)} : p_0 = 49\%$

et $H_0^{(2)} : p_0 = 51\%$

sont toutes deux conservées ! Art de la double négation :

\Leftarrow Au vu de ces données, on ne peut pas affirmer que p_0 ne soit pas

égal à $\begin{cases} 49\% \\ 51\% \end{cases} \rightarrow$ (on énonce une certaine compatibilité, une possibilité).

Il faut faire attention à ne pas effectuer de glissement sémantique ici...

* En revanche l'hypothèse $H_0^{(3)}$: $p_0 = 47\%$ est rejetée
« Ces données montrent que p_0 est différent de 47% »
(on énonce une impossibilité).

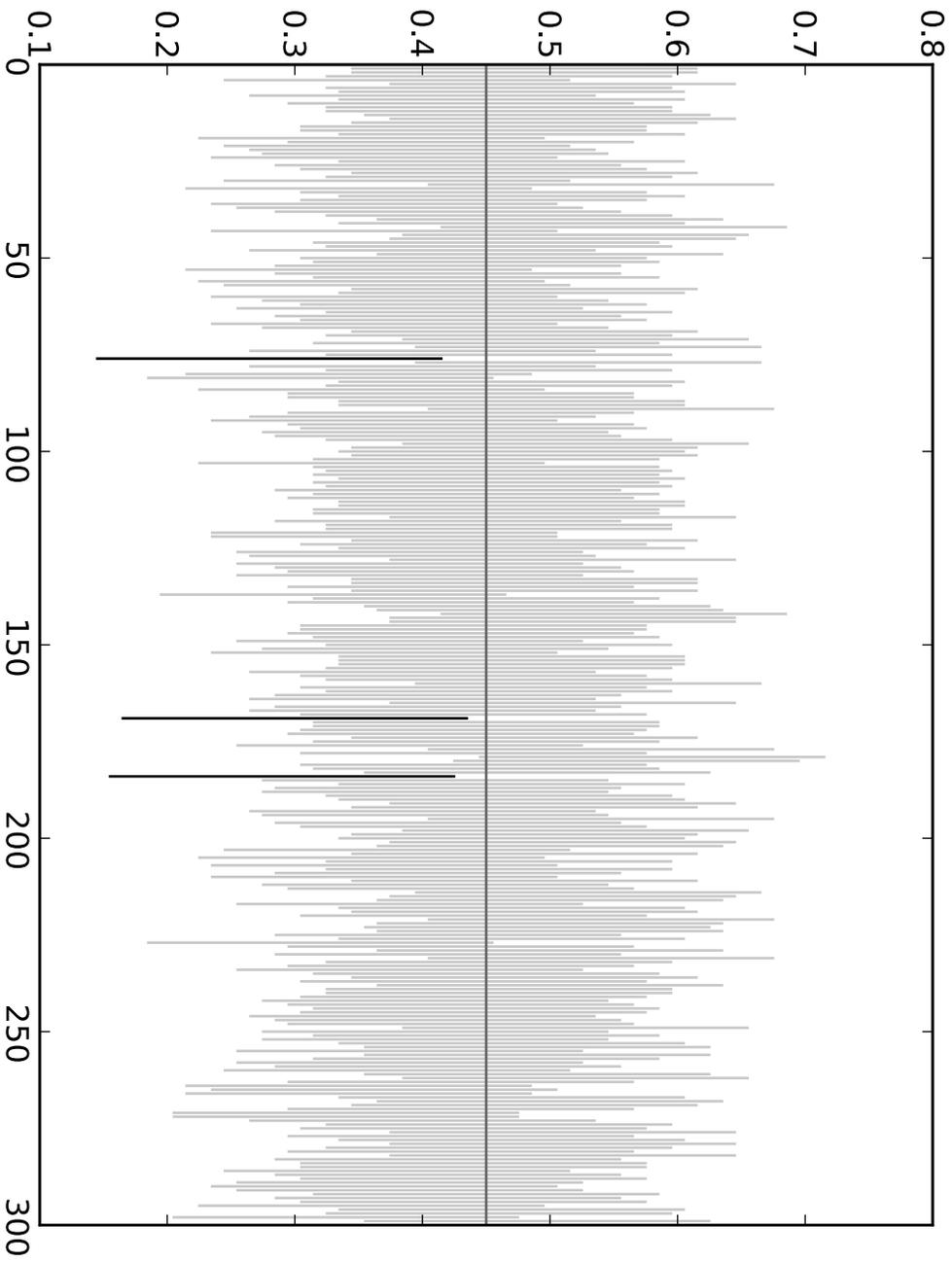
→ En statistique, H_0 est vraie jusqu'à preuve du contraire,
- soit parce qu'elle est vraiment vraie,
- soit parce l'écart du paramètre sous-jacent p_0 à la valeur de référence p_0^* est faible et que notre taille d'échantillon n ne permet pas de la mettre en évidence.

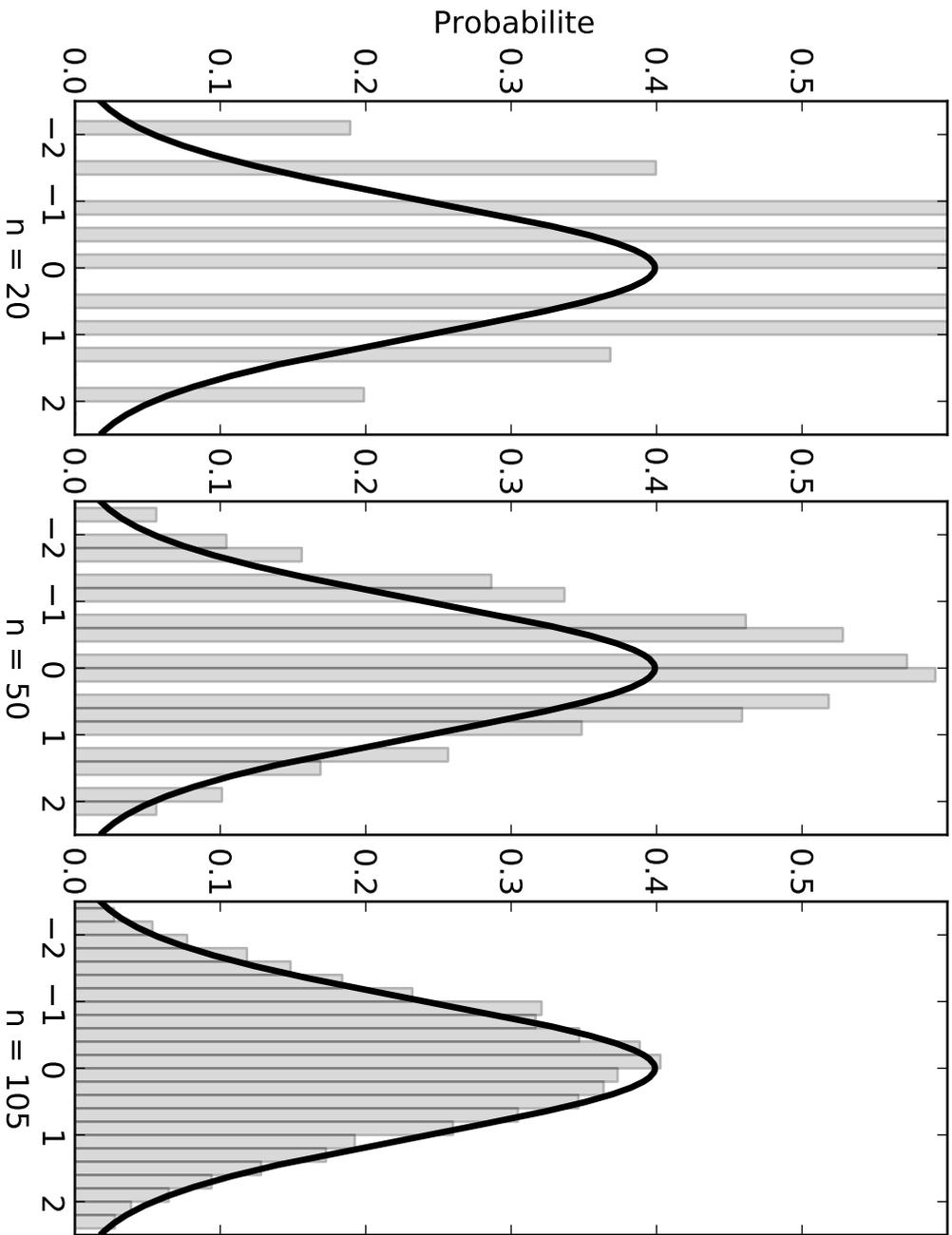
Un test :

- lorsqu'il conserve H_0 dit simplement qu'elle n'est pas impossible ou invraisemblable
- lorsqu'il rejette H_0 fait un progrès négatif: il exclut H_0 (et passe à H_1)

Remarque culturelle :

Les débats sur l'innocence ou la nocivité des OGM semblent contradictoires ; en fait, dans ces dernières, les hypothèses H_0' de nocivité ou H_0'' d'innocence sont simplement conservées (parfois / souvent pour cause de tailles d'échantillon faibles).





Comparaison des performances des trois intervalles de confiance

Pour $p = 0.45$ et $n = 100$, avec $N = 10\,000$ réplifications :

methode Hoeffding :

marge d'erreur 0.136

niveau de confiance 99.40

methode Plug-in :

marge d'erreur 0.097

niveau de confiance 94.51

methode Tchebychev :

marge d'erreur 0.224

niveau de confiance 100.00

Code pour simuler des intervalles fondés sur l'inégalité de Hoeffding

```
## MODULE IMPORTS

from numpy.random import binomial
from numpy import *
import matplotlib.pyplot as plt
import matplotlib.mlab as mlab

## FUNCTION AND PARAMETERS

def histogram_bernoulli(n, p, N):
    t = (random.binomial(n, p, size = (N))/n )
    return t

p = 0.50
n = 100.
N = 300
alpha = .05

## CONFIDENCE INTERVALS FOR BERNOULLI WITH Hoeffding's INEQUALITY

x = histogram_bernoulli(n, p, N)

half_width = sqrt(log(2. / alpha) / (2. * n))
x_IC = [x - half_width, x + half_width]

c = np.arange(N)
list_out = c[((x_IC[0] > p) | (x_IC[1] < p))]
list_in = c[((x_IC[0] < p) & (x_IC[1] > p))]

## FIGURES

fig = plt.figure()
ax = fig.add_subplot(111)

if list_out.size >0:
    ax.vlines(list_out, x_IC[0][list_out], x_IC[1][list_out], color = 'r')

ax.vlines(list_in, x_IC[0][list_in], x_IC[1][list_in], color = 'g')
ax.axhline(y = p)

plt.savefig("Real_confidence.pdf")
```

Code pour comparer les performances des trois intervalles de confiance

```
## IMPORT

from numpy.random import binomial
from numpy import *
import matplotlib.pyplot as plt
import matplotlib.mlab as mlab

## FUNCTION AND PARAMETERS

def histogram_bernoulli(n, p, N):
    t = (binomial(n, p, size = (N)) / n )
    return t

def half_width(n, alpha, method):

    if (method == 'Hoeffding'):
        t = sqrt(log(2. / alpha) / (2. * n))

    elif (method == 'Tchebychev'):
        t = 1. / (2. * sqrt(n * alpha))

    else:
        t = 0

    return t

# Half width of the interval confidence for the CLT and plug-in
# is empirical and would be computed outside the function.

def real_level(x_IC, p):
    N = len(x_IC[0]) * 1.
    c = arange(N)
    list_out = c[((x_IC[0] > p) | (x_IC[1] < p))]
    t = ( 1. - list_out.size / N ) * 100.
    return t

N = 10000.
n = 100.
p = 0.45
alpha = 0.05
```

```

## CONFIDENCE INTERVALS

x = histogram_bernoulli(n, p, N)

shaper = [array([-1.]), array([1.])]
results = {'Hoeffding' : [0.,0.], 'Tchebychev' : [0.,0.]}

# The Bernoulli samples are the same for the three
# interval confidence computations.

for i in results:
    half_width_spe = half_width(n, alpha, i)
    x_IC_spe = x + half_width_spe * (shaper * ones(x.shape))
    real_level_spe = real_level(x_IC_spe, p)

    results[i] = [half_width_spe, real_level_spe]

# TCL Plug in

half_width_plugin = 1.96 * sqrt((x*(1-x)/n))
x_IC_plugin = x + (shaper * half_width_plugin)
real_level_plugin = real_level(x_IC_plugin, p)

results['Plug-in'] = [mean(half_width_plugin), real_level_plugin]

## DISPLAY RESULTS

for i in results:
    print "methode % s : " % i
    print "marge d'erreur %.3f " % results[i][0]
    print "niveau de confiance %.2f\n" % results[i][1]

```

Code pour justifier l'approximation par une loi normale pour construire le troisième intervalle

```
## MODULE IMPORTS

from numpy.random import binomial
from numpy import *
import matplotlib.pyplot as plt
import matplotlib.mlab as mlab

## FUNCTIONS AND PARAMETERS
random.seed(100)

def data_bernoulli(n, p, N):
    mean_trial = binomial(n, p, size = (N))/n
    t = sqrt(n) * (mean_trial - p) / sqrt(mean_trial * (1. - mean_trial))
    return t

p = 0.45

mu = 0
sigma = 1

n = [20., 50., 105.]
N = 1e4 * ones(3)

x = [data_bernoulli(n[i], p, N[i]) for i in range(len(n)) ]

## HISTOGRAMS

num_bins = arange(-2.6, 2.6, .2)
fig = plt.figure()

for i in range(3):
    print i
    ax = fig.add_subplot(131 + i)
    ax.hist(x[i], num_bins, normed=1, facecolor='gray', alpha = 0.30)
    bincenters = arange(-2.5, 2.5, 0.01)
    y = mlab.normpdf( bincenters, mu, sigma)
    l = ax.plot(bincenters, y, 'k-', linewidth=2.5)
```

```
if i==0:
    ax.set_ylabel('Probabilite')
    ax.set_xlabel('n = 20')
if i==1:
    ax.set_xlabel('n = 50')
if i==2:
    ax.set_xlabel('n = 105')
ax.set_xlim(-2.5, 2.5)
ax.set_ylim(0, 0.6)
ax.grid(False)

plt.show(patches)
plt.savefig("Histogram_plug-in.pdf")
```



FIGURE 1 – Pan sur le bec du Canard! Les variations à $\pm 3\%$ des cotes de popularité ne sont pas statistiquement significatives.



FIGURE 2 – Il est facile de briller à peu de frais... Un peu d'efforts (cf. formules vues ensemble) et vous saurez faire mieux que ce mauvais détective!

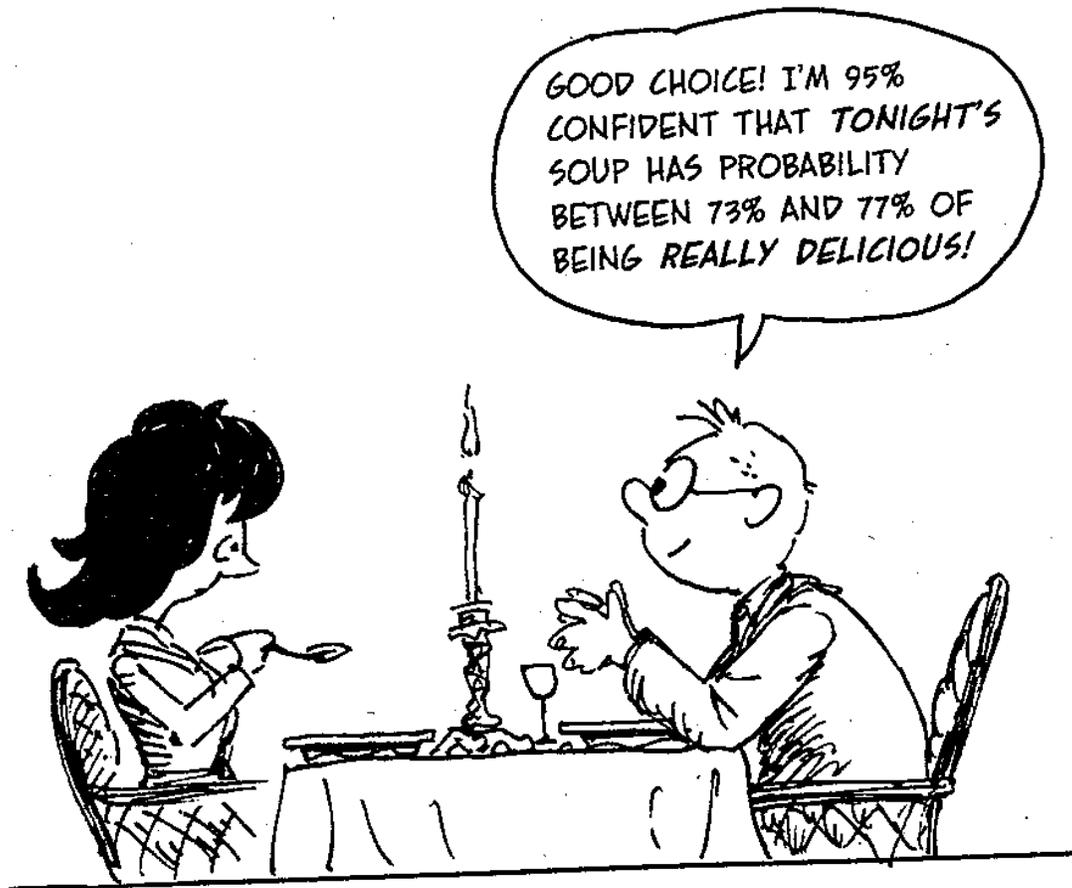
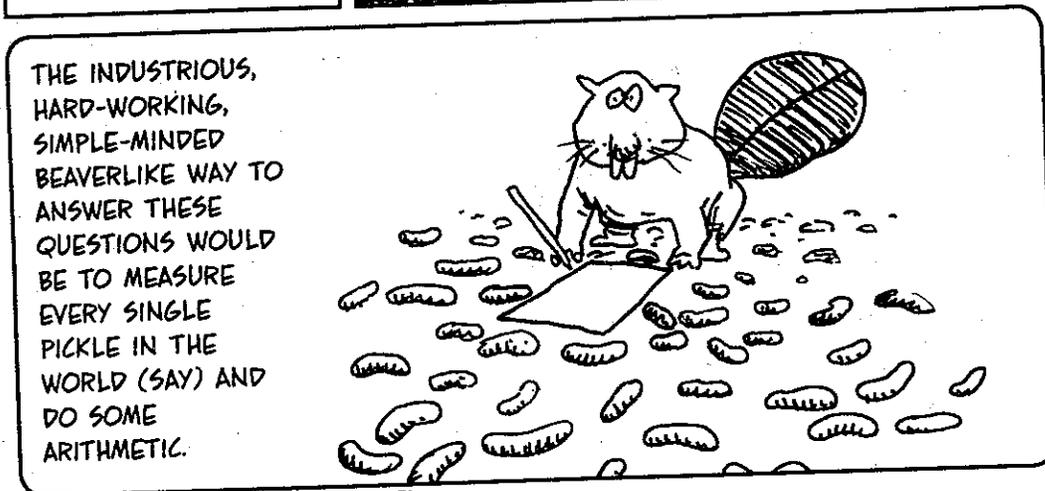
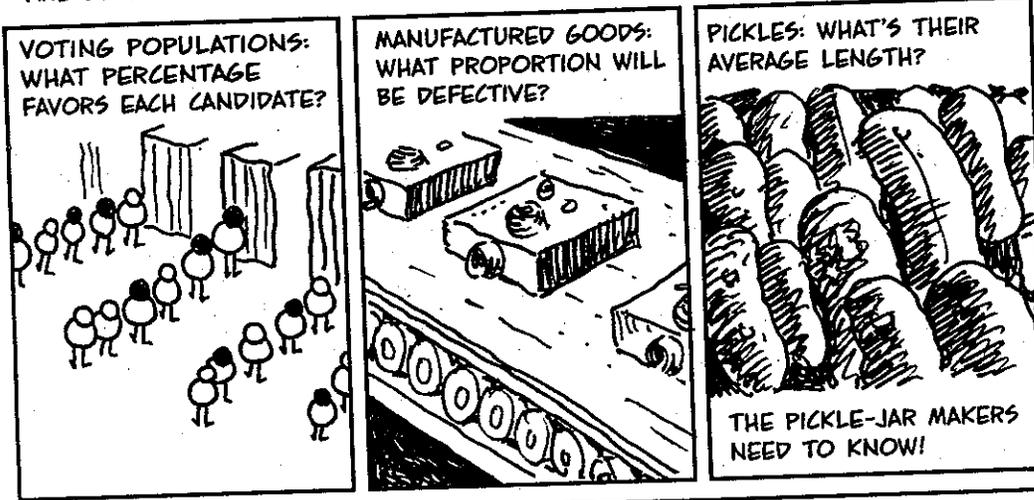


FIGURE 3 – La statistique donne des intervalles de confiance. C'est utile, mais pas nécessairement romantique.

THE PROBLEM WITH THE WORLD IS THAT THE COLLECTIONS OF STUFF IN IT ARE SO LARGE, IT'S HARD TO GET THE INFORMATION WE WANT:



BUT WE AREN'T BEAVERS—WE'RE STATISTICIANS! WE'RE LOOKING FOR THE EASY WAY OUT...

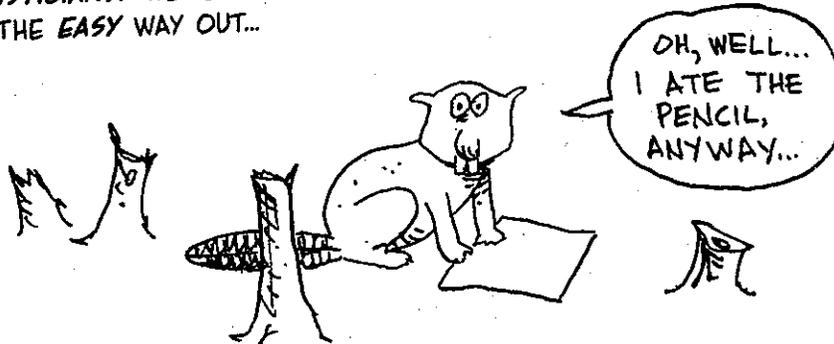


FIGURE 4 – Vous n'êtes pas des petites bêtes industrielles : en conséquence de quoi, vous ferez désormais des sondages plutôt que des décomptes exhaustifs sur des populations immenses.

“opinion way

LE FIGARO LCI RTL

Baromètre “opinionway FIDUCIAL

des élections régionales

Région Alsace

18 mars 2010

Toute publication totale ou partielle doit impérativement utiliser la mention complète suivante : « Baromètre OpinionWay – Fiducial pour Le Figaro / LCI / RTL » et aucune reprise de l'enquête ne pourra être dissociée de cet intitulé.

Méthodologie

- Étude réalisée auprès d'un **échantillon de 801 personnes**, représentatif de la **population d'Alsace, âgées de 18 ans et plus et inscrites sur les listes électorales**.
- L'échantillon a été constitué selon la méthode des quotas, au regard des critères de sexe, d'âge, de catégorie socioprofessionnelle, de catégorie d'agglomération et de département de résidence.
- Mode d'interrogation:** L'échantillon a été interrogé par téléphone au domicile des personnes.
- Dates de terrain:** les interviews ont été réalisées les **16 et 17 mars 2010**.
- OpinionWay rappelle par ailleurs que les résultats de ce sondage doivent être lus en tenant compte des marges d'incertitude : **3 à 4 points au plus** pour un échantillon de 800 répondants.
- La notice de cette enquête est consultable à la commission des sondages.

Le Figaro-LCI - Baromètre OpinionWay - Fiducial des élections régionales - Région Alsace / 18 Mars 2010 page 3

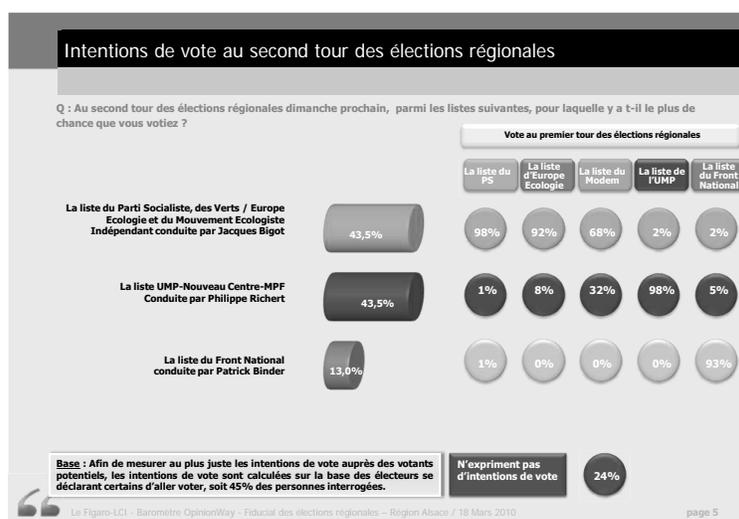


FIGURE 5 – Un exemple de résultat de sondage fourni par un institut avec indication de la marge d'erreur (d'incertitude) possible.

OUR METHOD IS TO TAKE A **SAMPLE**... A RELATIVELY SMALL SUBSET OF THE TOTAL POPULATION, THE WAY POLLSTERS DO AT ELECTION TIME.



AN OBVIOUS QUESTION IS: HOW BIG A SAMPLE DO WE HAVE TO TAKE TO GET MEANINGFUL RESULTS?



AND THE ANSWER, WHICH YOU SHOULD INSCRIBE IN YOUR BRAIN FOREVERMORE, WILL TURN OUT TO BE: IF n IS THE NUMBER OF ITEMS IN THE SAMPLE, THEN EVERYTHING IS GOVERNED BY

$$\frac{1}{\sqrt{n}}$$

GOVERNED BY $\frac{1}{\sqrt{n}}$? DIDN'T EVEN KNOW IT WAS ON THE BALLOT!



FIGURE 6 – Planification de sondages et mauvais jeux de mots associés.