

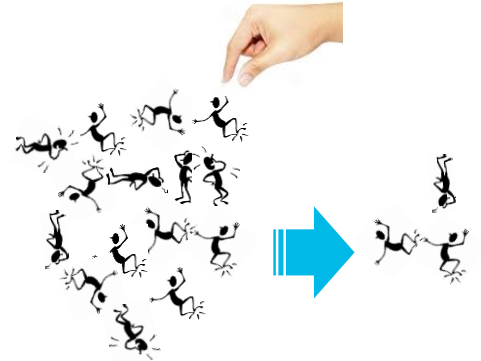
# Meeting 16:

## Sampling issues (without and with a decision-analytic approach)

# Sampling issues in classical inference

Sampling in a general sense can be of two kinds:

I. Sampling without replacement from a finite population



II. Sampling from an infinite population/from a process/(with replacement from a finite population)



Case II also covers making several measurements on a specific object.

*For simple random sampling ...*

## How many units should be sampled?

Depends to the objective of sampling:

- to estimate the value of a parameter?
  - the number of units should be chosen from a desired bound on the point estimate (desired length of a confidence interval)
  - requires prior knowledge of the population variance

$$n_0 \geq \frac{4 \cdot (z_{\alpha/2} \cdot \sigma)^2}{D^2}$$

$$n = \frac{n_0}{1 + (n_0 - 1)/N}$$

$D$  is the desired length of a confidence interval,  $N$  is the size of the population

- to be able to reject a (range of) value(s) of a parameter with a high probability when the true parameter value is at a certain distance from that (range of) value(s)?
  - the number of units should be chosen so that the power function of a hypothesis test at such a distance is at least as high at that probability
  - requires prior knowledge of the population variance

## *For stratified sampling and cluster sampling ...*

How many units should be sampled, and if the population is stratified, how should they be allocated over strata?

- to estimate the value of a parameter?
  - the number of units should be chosen from a desired bound on the point estimate (desired length of a confidence interval)
  - requires prior knowledge of the population variance, stratum variances and sizes, numbers and sizes of clusters

*For sampling from several populations (and planning of experiments) ...*

How many units should be sampled from each population? How many measurements should be made for each experimental setup?

- to be able to detect with a certain probability a difference in a parameter between two populations? to be able to reject the hypothesis of no differences in effect between the different set-ups of an experiment?
  - the number of units should be chosen so that the power function of the hypothesis test used at such a distance is at least as high at that probability
  - requires prior knowledge of the population variances

# The Bayesian approach to sampling

How many units should be sampled to be able to state with a certain probability  $p_0$  that

- a parameter is at least/most a certain value or within a specific range?
- the difference in a parameter between two populations is at least/at most a certain value or within a specific range?
- ...

$$P(\theta \geq \theta_0 | n, y_1, \dots, y_n) = \int_{\theta \geq \theta_0} q(\theta | n, y_1, \dots, y_n; \psi) d\theta$$
$$q(\theta | n, y_1, \dots, y_n; \psi) = \frac{f(y_1, \dots, y_n | \theta; \psi) \cdot p(\theta | \psi)}{\int_{u \in \Theta} f(y_1, \dots, y_n | u; \psi) \cdot p(u | \psi) du}$$

Solve for  $n$   $P(\theta \geq \theta_0 | n, y_1, \dots, y_n) \geq p_0$

## *Examples from drugs sampling*

Consignments of drugs in forms of pills, capsules, ampoules or plastic bags can be very extensive (e.g. thousands of pills in big sacs)



Analysis must by legal reasons be made “unit-wise” and is time-consuming  $\Rightarrow$  as small sample sizes as possible are desired.

However, drug seizures are usually homogeneous with respect to the active substance in each unit.

For so-called *identifying analysis* what is of interest is whether a unit contains the active (illicit) substance or not – percentages of that substance or presence of other substances is of minor importance.

Hence, the sampling scheme (hypergeometric or approx. binomial) is expected to render homogeneous samples, i.e.  $x = n$  (all units are illicit ones) or  $x = 0$  (no unit is illicit).

# 1. Homogeneity expected from visual inspection and experience

Consider a case with a seizure of 5000 pills, all of the same colour (blue), form (circular) and printing (e.g. the Mitsubishi trade mark)



The forensic scientist would say “this is a seizure of Ecstasy pills”.

Consider historical cases with blue pills

Group the cases into  $M$  clusters with respect to another parameter, e.g. the print on the pill.

Find an estimate of the *prior distribution* for the proportion  $\theta$  of Ecstasy pills among blue pills.

Nordgaard A. (2006) Quantifying experience in sample size determination for drug analysis of seized drugs. *Law, Probability and Risk* **4**: 217-225



Cluster	Accumulated size of seizure	Accumulated size of sample	Number of Ecstasy pills	Number of Non-Ecstasy pills
1	$N_1$	$n_1$	$x_1$	$n_1 - x_1$
2	$N_2$	$n_2$	$x_2$	$n_2 - x_2$
...	...	...	...	...
$M$	$N_M$	$n_M$	$x_M$	$n_M - x_M$

Use a generic *beta prior* for the proportion  $\theta$  of Ecstasy pills in the current seizure:

$$p(\theta | v_1, v_2) = \frac{\theta^{v_1-1} \cdot (1-\theta)^{v_2-1}}{B(v_1, v_2)}; 0 \leq \theta \leq 1$$

$$p(\theta | \nu_1, \nu_2) = \frac{\theta^{\nu_1-1} \cdot (1-\theta)^{\nu_2-1}}{B(\nu_1, \nu_2)}$$

Use the grouped data to estimate the parameters  $\nu_1$  and  $\nu_2$  of this beta prior.

This can be done by the *maximum likelihood method* using that the probability of obtaining  $x_i$  Ecstasy pills in cluster  $i$  is

$$P(\tilde{x}_i = x_i | \theta, n_i) \approx \frac{\binom{\lfloor N_i \cdot \theta \rfloor}{x_i} \cdot \binom{\lfloor N_i \cdot (1-\theta) \rfloor}{n_i - x_i}}{\binom{N_i}{n_i}}$$

*Hypergeometric  
distribution*

where “ $\lfloor \cdot \rfloor$ ” stands for rounding downwards to nearest integer

The likelihood function of  $\nu_1$  and  $\nu_2$  in light of observed numbers in all clusters ( $\mathbf{x} = (x_1, \dots, x_M)$ ) then becomes

$$L(\nu_1, \nu_2 | \mathbf{x}) = \prod_{i=1}^M P(\tilde{x}_i = x_i | \nu_1, \nu_2, n_i) = \prod_{i=1}^M \int_0^1 P(\tilde{x}_i = x_i | \theta, n_i) \cdot p(\theta | \nu_1, \nu_2) d\theta$$

The obtained point estimates of  $\nu_1$  and  $\nu_2$  can be assessed with respect to *bias* and *variance* using *bootstrap resampling*.

In Nordgaard (2006) original point estimates of  $\nu_1$  and  $\nu_2$  for historical cases of blue pills at SKL (*now* NFC) are

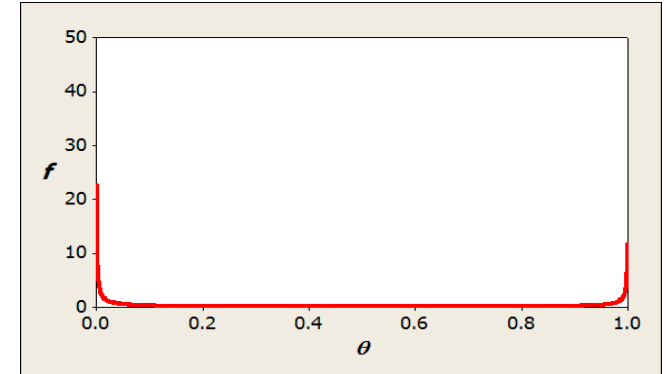
$$\hat{\nu}_1 = 0.075 \quad \text{and} \quad \hat{\nu}_2 = 0.224$$

Bias adjusted estimates are

$$\hat{\nu}_1^* = 0.038 \quad \text{and} \quad \hat{\nu}_2^* = 0.133$$

and upper 90% confidence limits for the true values of  $\nu_1$  and  $\nu_2$  are

$$\nu_1 \leq 0.062 \quad \text{and} \quad \nu_2 \leq 0.262$$



Should confidence limits be used in empirical Bayes?

Now, assume the forthcoming sample of  $n$  units will consist entirely of Ecstasy pills. (*Otherwise the case will be considered “non-standard”*)

The sample size is determined so that the *posterior* probability of  $\theta$  being higher than a certain proportion, say 50 %, is at least say 99% (referred to as 99% *credibility*)

For large seizures the posterior distribution of  $\theta$  given all  $n$  sample units consist of Ecstasy is also *beta*:

$$f(\theta | n, \nu_1, \nu_2) = \frac{\theta^{\nu_1+n-1} \cdot (1-\theta)^{\nu_2-1}}{B(\nu_1 + n, \nu_2)} ; 0 \leq \theta \leq 1$$

Thus we solve for  $n$

$$\int_{0.50}^1 f(\theta | n, \nu_1, \nu_2) d\theta \geq 0.99 \quad \Leftrightarrow \quad \frac{\int_{0.50}^1 \theta^{\nu_1+n-1} \cdot (1-\theta)^{\nu_2-1} d\theta}{B(\nu_1 + n, \nu_2)} \geq 0.99$$

where  $\nu_1$  and  $\nu_2$  are replaced by their (adjusted) point estimates (or upper confidence limits).

For the above case we find that with the bias-adjusted point estimates

$$\hat{v}_1^* = 0.038 \text{ and } \hat{v}_2^* = 0.133$$

the required sample size is at least **3** and with the upper confidence limits used instead (i.e with 0.062 and 0.262) the required sample size is at least **4**

There are in general no large differences between different choices of estimated parameters, nor between different colours of Ecstasy pills.

A general sampling rule of  $n=5$  can therefore be used to state with 99% credibility that at least 50% of the seizure consists of Ecstasy pills. For a higher proportion, a sample size around 12 appears to be satisfactory.

For smaller seizures it is more wise to rephrase the requirement in terms of the number of Ecstasy units in the non-sampled part of the seizure.

The posterior beta distribution is then replaced with a *beta-binomial* distribution.

## 2. Homogeneity stated upon inspection only

Consider now a case with a (large) seizure of drug pills of which the forensic scientist cannot directly suspect the contents.

Visual inspection  $\Rightarrow$  All pills seem to be identical

Can we substitute the “experience” from the Ecstasy case?

*UV-lightning*

Pills can be inspected under UV light.

The fluorescence differs between pills with different chemical composition and looking at a number of pills under UV light would thus reveal (to greatest extent) heterogeneity.

*does not work for capsules and ampoules*

Uncertainty of this procedure lies mainly with the person who does the inspection  
 $\Rightarrow$  Experiment required!

Assume a prior  $g(\theta)$  for the proportion of pills in the seizure that contains a certain (but possibly unknown) illicit drug.

For sake of simplicity, assume that pills may be of two kinds (the illicit drug or another substance).

Let  $Y$  be a random variable associated with the inspection such that

$$Y = \begin{cases} 0 & \text{if inspection gives "all pills are identical"} \\ 1 & \text{if inspection gives "differences among pills"} \end{cases}$$

Relevant case is  $Y = 0$  (*Otherwise the result of the UV-inspection has rejected the assumption of homogeneity.*)

Now,  $P(Y = 0 \mid \theta)$  for  $0 < \theta < 1$

is the *false positive probability* as a function of  $\theta$  (if a positive result means that no heterogeneity is detected)

while  $P(Y = 0 \mid \theta = 0) + P(Y = 0 \mid \theta = 1)$

is the true positive probability.

The prior  $g$  can be updated using this information (when available)

$$h(\theta | Y = 0) = \frac{\Pr(Y = 0 | \theta) \cdot g(\theta)}{\int_0^1 \Pr(Y = 0 | \lambda) \cdot g(\lambda) d\lambda}$$

Note that an *non-informative prior* (i.e.  $g(\theta) \equiv 1$  ;  $0 \leq \theta \leq 1$ ) can be used.

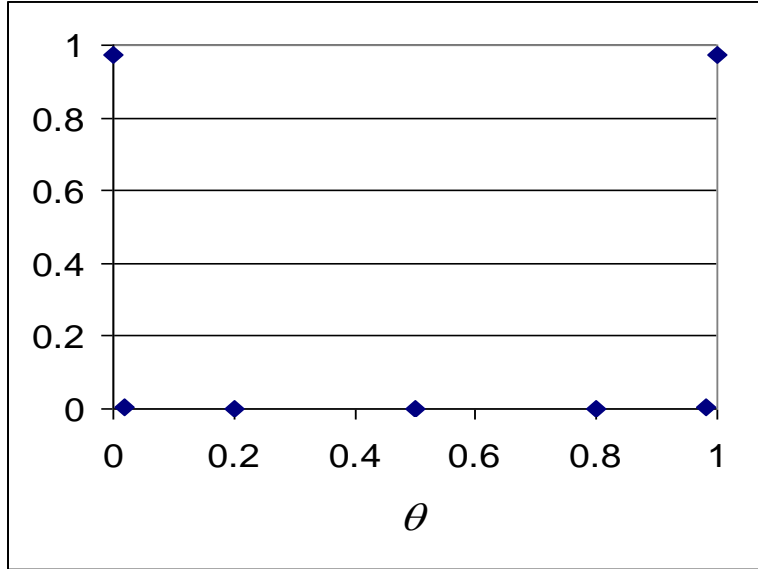
The updated prior (i.e. the posterior upon UV-inspection) can then be used analogously to the previous case (Ecstasy).



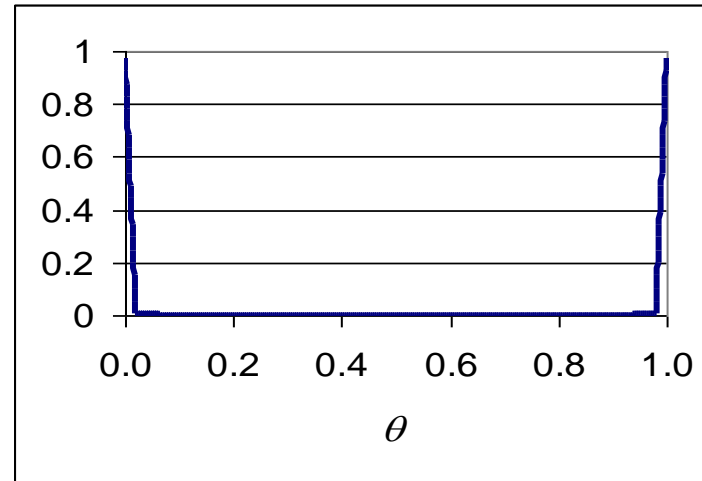
*Example* Experiment (conducted at SKL (*now* NFC))

- 8 types of pills with different substances were used to form 9 different mixtures (i.e. of two proportions) of 2 types of pills
- Each mixture was prepared by randomly shuffling 100 pills with the current proportions on a tray that was put under UV-light
- 10 case-workers made inspections in random order such that a total of 114-117 inspections were made for each mixture

Data can be illustrated by plotting estimated probabilities for  $Y = 0$  vs.  $\theta$



Linear interpolation gives

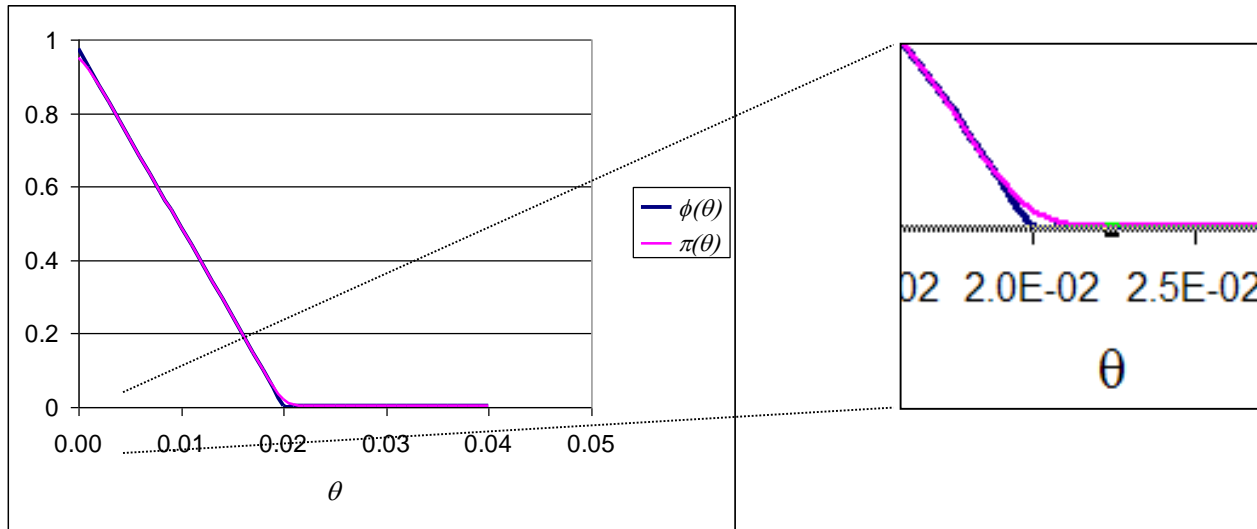


$$\hat{P}(Y = 0 | \theta) = \underline{\underline{\phi(\theta)}} = \begin{cases} 0.97 - 48.5 \cdot \theta & 0 \leq \theta \leq 0.02 \\ 0.005 - 0.024 \cdot \theta & 0.02 \leq \theta \leq 0.20 \\ 0 & 0.20 < \theta < 0.80 \\ -0.019 + 0.024 \cdot \theta & 0.80 \leq \theta < 0.98 \\ -47.5 + 48.5 \cdot \theta & 0.98 \leq \theta \leq 1 \end{cases}$$

To avoid the vertices at  $\theta = 0.02, 0.20, 0.80$  and  $0.98$ , the linearly interpolated values are smoothed using a Kernel function:

$$\pi(\theta) = \int_0^1 K(\theta - \lambda) \cdot \phi(\lambda) d\lambda$$

where  $K(x)$  is a symmetric function integrating to one over its support.



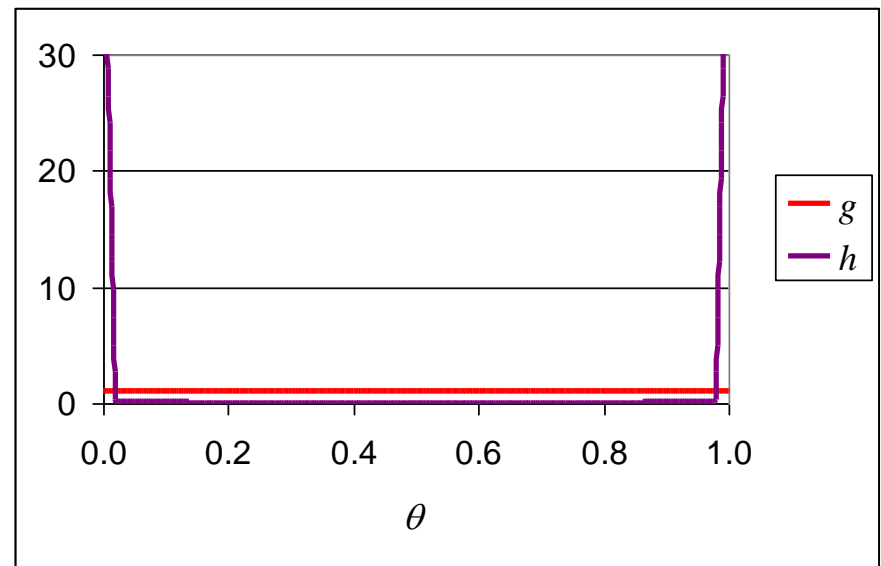
Now, the prior can be updated using this smoothed function as an estimate of  $\Pr(Y = 0 | \theta)$ , i.e.

$$h(\theta | Y = 0) = \frac{\pi(\theta) \cdot g(\theta)}{\int_0^1 \pi(\lambda) \cdot g(\lambda) d\lambda}$$

(With a non-informative prior  $g$ , this simplifies into

$$h(\theta | Y = 0) = \frac{\pi(\theta)}{\int_0^1 \pi(\lambda) d\lambda} \quad )$$

Comparison of the non-informative prior and the updated prior



Now, let  $x$  be the number of illicit drug pills found in a sample of  $n$  pills.

Analogously with the Ecstasy case  $n$  should be determined so that if  $x = n$  a 99% credible lower limit for  $\theta$  is 50% (or even higher).

With the updated prior derived the following table of posterior probabilities is obtained.

$n$	$\Pr(\theta > 0.5 \mid x = n, Y = 0)$
3	0.99996032237
4	0.99999475894
5	0.99999924614
6	0.99999988597
7	0.99999998211
8	0.99999999711
9	0.99999999952
10	0.99999999992

Thus, a sample size of  $n = 3$  units is satisfactory.

Slightly higher values may be recommended due to the limits of the experiment

# The decision-theoretic approach

As was previously taken up, the decision about sampling (and how much to sample) builds on the expected value of sample information,  $EVSI(n)$ , and the optimal sample size is the value of  $n$  for which the *expected net gain of sampling*

$$ENG S(n) = EVSI(n) - CS(n)$$

is maximised.

*Example:* Return to the examples with illicit pills

Assume we should make a decision on whether the proportion,  $\theta$ , of Ecstasy pills in a seizure of 1000 pills is less than or at least 50 %.

The possible actions are  $a_1 = “\theta < 50 \text{ \%}”$  and  $a_2 = “\theta \geq 50 \text{ \%}”$

Assume a “0– $k_i$ ” loss function as

	$\theta < 50 \text{ \%}$	$\theta \geq 50 \text{ \%}$
$a_1$	0	1
$a_2$	10	0

Assume a prior distribution of  $\theta$  as  $Beta(\nu_1, \nu_2)$  with  $\nu_1 = 0.038$  and  $\nu_2 = 0.133$  (the point estimates from the empirical Bayes procedure)

$$\Rightarrow P(\theta < 0.50 | \nu_1 = 0.038, \nu_2 = 0.133) = \int_0^{0.5} \frac{\theta^{0.038-1} \cdot (1-\theta)^{0.133-1}}{B(0.038, 0.133)} d\theta \approx 0.780$$

$\Rightarrow$

$$EL(a_1) = 0 \cdot 0.780 + 1 \cdot 0.220 = 0.22$$

$$EL(a_2) = 10 \cdot 0.780 + 0 \cdot 0.220 = 7.8$$

$$\Rightarrow a'_{opt} = a_1$$

The number of Ecstasy pills in a sample of  $n$  pills is the  $Bin(n, \theta)$ .

Pre-assuming the sample to be completely homogeneous, i.e. either all are Ecstasy pills or all are non-Ecstasy pills gives the posterior distribution to be any of

$Beta(\nu_1 + n, \nu_2)$  [all are Ecstasy] and  $Beta(\nu_1, \nu_2 + n)$  [all are non-Ecstasy]



With  $Beta(\nu_1 + n, \nu_2)$  as posterior the expected posterior losses are

$$EL^{(i)}(a_1|n) = 0 \cdot \Pr(\theta < 0.5|n, n) + 1 \cdot \Pr(\theta \geq 0.5|n, n) = \int_{0.5}^1 \frac{\theta^{0.038+n-1} (1-\theta)^{0.133-1}}{B(0.038+n, 0.133)} d\theta$$

$$EL^{(i)}(a_1|n) = 10 \cdot \Pr(\theta < 0.5|n, n) + 0 \cdot \Pr(\theta \geq 0.5|n, n) = 10 \int_0^{0.5} \frac{\theta^{0.038+n-1} (1-\theta)^{0.133-1}}{B(0.038+n, 0.133)} d\theta$$

With  $Beta(\nu_1, \nu_2 + n)$  as posterior the expected posterior losses are

$$EL^{(ii)}(a_1|n) = 0 \cdot \Pr(\theta < 0.5|0, n) + 1 \cdot \Pr(\theta \geq 0.5|0, n) = \int_{0.5}^1 \frac{\theta^{0.038-1} (1-\theta)^{0.133+n-1}}{B(0.038, 0.133+n)} d\theta$$

$$EL^{(ii)}(a_2|n) = 10 \cdot \Pr(\theta < 0.5|0, n) + 0 \cdot \Pr(\theta \geq 0.5|0, n) = 10 \int_0^{0.5} \frac{\theta^{0.038-1} (1-\theta)^{0.133+n-1}}{B(0.038, 0.133+n)} d\theta$$

How would we obtain the optimal sample size?