# Meeting 15:
# Bayesian inference:
## Comparisons with forensic examples

# Comparisons

Consider the issue to test

$$H_0 : \theta_1 \geq \theta_2 \text{ (or } \theta_1 \leq \theta_2 \text{ ) against}$$
$$H_1: \theta_1 < \theta_2 \text{ (or } \theta_1 < \theta_2)$$

where $\theta_1$ and $\theta_2$ are of the same parameter type, but from two different populations.

By rewriting $H_0$ as $\theta_1 - \theta_2 \geq 0$ (and correspondingly $H_1$ as $\theta_1 - \theta_2 < 0$ we have transferred the problem back to a test of the value of the compound parameter $\theta = \theta_1 - \theta_2$

Statistical calculations for the "merging" of two populations will then be needed to sort out the problem – independent populations and independent samples gives product priors and product likelihoods.

# Two-sided testing

So far the hypotheses specified through parameters have been of the kind

$$H_0 : \theta \in \Theta_0$$
$$H_1 : \theta \in \Theta_1 \ (= \Theta \setminus \Theta_0)$$

where $\Theta$ is the parameter space.

This may be viewed as a classification problem – the issue is to choose between two hypotheses (two models) the one that best describes the data given prior probabilities for both hypotheses.

The specification through a parameter induces the need for integrating the likelihood function as soon as any of the hypotheses is composite.

Very often, though, the issue is to analyse two sets of data material for making inference about whether they originate from the same model.

$H_0$ : Model for data set 1 = Model for data set 2
$H_1$ : Model for data set 1 ≠ Model for data set 2

Not so much about the models themselves, but essentially about whether they are common or not – *n.b.* completely specified models, not just up to the value of a parameter.

*Examples*

• Are the effects of the two treatments equal or not?

• Do the two seizures of amphetamine come from the same manufacturing batch?

• Was the recovered shoeprint made by the suspects' shoe ?

In many situations a parametric model is common for the two data sets.

The difference is then expressed in terms of one or several parameters.

$$H_0 : \theta_1 = \theta_2$$
$$H_1 : \theta_1 \neq \theta_2$$

Assuming independent data sets the general likelihood function is

$$L(\theta_1, \theta_2 | Data) = L(\theta_1 | Data\ set\ 1) \cdot L(\theta_2 | Data\ set\ 2)$$

However, the likelihood function for $H_0$ is a function of one parameter value only:

$$L(H_0 | Data) = L(\theta | Data\ set\ 1) \cdot L(\theta | Data\ set\ 2)$$

since $H_0$ states $\theta_1 = \theta_2 = \theta$.

With prior density $p(\theta)$ for $\theta$ (common to both hypothesis) the Bayes factor becomes

$$B = \frac{\int_{\Theta} L(\theta | Data\ set\ 1) \cdot L(\theta | Data\ set\ 2) \cdot p(\theta) d\theta}{\int_{\Theta} L(\theta_1 | Data\ set\ 1) \cdot L(\theta_2 | Data\ set\ 2) \cdot p(\theta_1) \cdot p(\theta_2) d\theta_1\ d\theta_2} =$$

$$= \frac{\int_{\Theta} L(\theta | Data\ set\ 1) \cdot L(\theta | Data\ set\ 2) \cdot p(\theta) d\theta}{\int_{\Theta} L(\theta | Data\ set\ 1) \cdot p(\theta) d\theta \cdot \int_{\Theta} L(\theta | Data\ set\ 2) \cdot p(\theta) d\theta}$$

Lindley,
*Biometrika,* 1977

*Example*

Assume we are comparing two seizures of cannabis with respect to their mean concentration of THC (*Tetrahydrocannabinol, the active narcotic substance*).
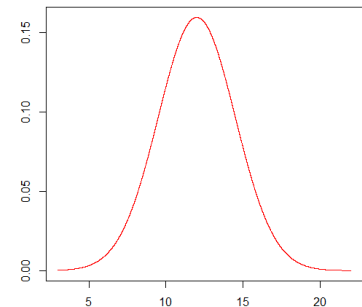
Denote the two means $\mu_1$ and $\mu_2$ respectively.

The concentration for the types of material in question (typically buds of cannabis plants) varies normally between 5 and 20 % with a peak around 12 %. We therefore assign a prior density for the mean concentration of *any* material of this type as

$$\mu \sim N\left(\text{mean} = 12, \text{standard deviation} = \frac{20-5}{6}\right) = N(12, 2.5)\,[\%]$$

For symmetric close-to-normal distributions a reasonable estimate of the standard deviation is Range/6.
The range is almost covered by the mean $\pm\,3$ standard deviations for a normal distribution.

To simply formulas denote the prior $N(\phi, \tau)$ where $\phi = 12$ and $\tau$, the standard deviation, is 2.5.

$$\Rightarrow p(\mu) = \frac{1}{\tau\sqrt{2\pi}} \cdot \exp\left\{-\frac{(\mu-\phi)^2}{2\tau^2}\right\}$$

Now, let's say we have a method of measurement that gives a value $x$, which is normally distributed $N(\mu, \sigma)$, where $\mu$ is the true mean concentration and $\sigma$ is the standard deviation.

Let's further assume that the method has been validated to provide a standard deviation of about 0.1 percentage points ➔ $x \sim N(\mu, 0.1)$

For $n_1$ measurements $x_{11}, \ldots, x_{1n_1}$ on the first seizure we obtain the likelihood function

$$L(\mu_1 | x_{11}, \ldots, x_{1n_1}) = \prod_{i=1}^{n_1} \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp\left\{-\frac{(x_{1i}-\mu_1)^2}{2\sigma^2}\right\} =$$

$$= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^2 \cdot \exp\left\{-\frac{1}{2\sigma^2}\sum_{i=1}^{n_1}(x_{1i}-\mu_1)^2\right\}$$

This can be shown to be

$$L\left(\mu_1 \big| x_{11}, \ldots, x_{1n_1}\right) = g_1\left(x_{11}, \ldots, x_{1n_1}, \sigma\right) \cdot \exp\left\{-\frac{n_1}{2\sigma^2}\left(\overline{x_1} - \mu_1\right)^2\right\} =$$

$$= g_1 \cdot \exp\left\{-\frac{n_1}{2\sigma^2}\left(\overline{x_1} - \mu_1\right)^2\right\}$$

where the function $g_1$ does not depend on $\mu_1$.

Analogously, for $n_2$ measurements $x_{21}, \ldots, x_{2n_2}$ on the second seizure we obtain the likelihood function

$$L\left(\mu_2 \big| x_{21}, \ldots, x_{2n_2}\right) = g_2\left(x_{21}, \ldots, x_{2n_2}, \sigma\right) \cdot \exp\left\{-\frac{n_2}{2\sigma^2}\left(\overline{x_2} - \mu_2\right)^2\right\} =$$

$$= g_2 \cdot \exp\left\{-\frac{n_2}{2\sigma^2}\left(\overline{x_2} - \mu_1\right)^2\right\}$$

Hence, the Bayes factor is

$$B = \frac{\int_{-\infty}^{\infty} g_1 \cdot \exp\left\{-\frac{n_1}{2\sigma^2}(\overline{x_1} - \mu)^2\right\} \cdot g_2 \cdot \exp\left\{-\frac{n_2}{2\sigma^2}(\overline{x_2} - \mu)^2\right\} \cdot \frac{1}{\tau\sqrt{2\pi}} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu}{\int_{-\infty}^{\infty} g_1 \cdot \exp\left\{-\frac{n_1}{2\sigma^2}(\overline{x_1} - \mu)^2\right\} \cdot \frac{1}{\tau\sqrt{2\pi}} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu \cdot \int_{-\infty}^{\infty} g_2 \cdot \exp\left\{-\frac{n_2}{2\sigma^2}(\overline{x_2} - \mu)^2\right\} \cdot \frac{1}{\tau\sqrt{2\pi}} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu} =$$

$$= \frac{\int_{-\infty}^{\infty} \exp\left\{-\frac{n_1}{2\sigma^2}(\overline{x_1} - \mu)^2\right\} \cdot \exp\left\{-\frac{n_2}{2\sigma^2}(\overline{x_2} - \mu)^2\right\} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu}{\frac{1}{\tau\sqrt{2\pi}} \cdot \int_{-\infty}^{\infty} \exp\left\{-\frac{n_1}{2\sigma^2}(\overline{x_1} - \mu)^2\right\} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu \cdot \int_{-\infty}^{\infty} \exp\left\{-\frac{n_2}{2\sigma^2}(\overline{x_2} - \mu)^2\right\} \cdot \exp\left\{-\frac{(\mu - \phi)^2}{2\tau^2}\right\} d\mu}$$

*A little bit tedious to sort out*

# Simplification when comparing two normal means

Express the hypotheses as

$$H_0 : \mu_1 - \mu_2 = \delta = 0$$
$$H_1 : \mu_1 - \mu_2 = \delta \neq 0$$

Above, it was shown that $\quad L\left(\mu_1 \middle| x_{11}, \ldots, x_{1n_1}\right) = g_1 \cdot \exp\left\{-\frac{n_1}{2\sigma^2}\left(\overline{x_1} - \mu_1\right)^2\right\}$

If we "reduce" the sample $x_{11}, \ldots, x_{1n_1}$ to its sample mean, i.e. $\overline{x_1}$
the likelihood function for $\mu_1$ becomes

$$L\left(\mu_1 \middle| \overline{x_1}\right) = \frac{1}{\left(\sigma/\sqrt{n_1}\right)\sqrt{2\pi}} \cdot \exp\left\{-\frac{\left(\overline{x_1} - \mu_1\right)^2}{2 \cdot \left(\sigma^2/n_1\right)}\right\}$$

since $\quad \overline{x_1} \sim N\left(\text{mean} = \mu_1, \text{standard deviation} = \frac{\sigma}{\sqrt{n_1}}\right)$

Note that

$$L_1\left(\mu_1 \middle| x_{11}, \ldots, x_{1n_1}\right) \propto \exp\left\{-\frac{n_1}{2\sigma^2}\left(\overline{x_1} - \mu_1\right)^2\right\}$$

$$L\left(\mu_1 \middle| \overline{x_1}\right) \propto \exp\left\{-\frac{\left(\overline{x_1} - \mu_1\right)^2}{2 \cdot \left(\sigma^2/n_1\right)}\right\} = \exp\left\{-\frac{n_1}{2\sigma^2}\left(\overline{x_1} - \mu_1\right)^2\right\}$$

Both likelihood functions are proportional to a common essential part, i.e. the part that contains $\mu_1$.

This is due to that the sample mean is a *sufficient statistic* for the population mean.

Analogously

$$L\left(\mu_2 \middle| \overline{x_2}\right) \propto \exp\left\{-\frac{n_2}{2\sigma^2}\left(\overline{x_2} - \mu_1\right)^2\right\}$$

Now, for independent samples $\quad \overline{x_1} - \overline{x_2} \sim N\left( \mu_1 - \mu_2, \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}} \right)$

With $\mu_1 - \mu_2 = \delta$ and $\sigma_1 = \sigma_2 = \sigma$ we get

$$\overline{x_1} - \overline{x_2} \sim N\left( \delta, \sigma\sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}} \right)$$

$\overline{x_1} - \overline{x_2}$ is a sufficient statistic for $\mu_1 - \mu_2 = \delta$ and "reducing" the two data sets to $\overline{x_1} - \overline{x_2}$ gives the likelihood function for $\delta$ as

$$L\left(\delta \mid \overline{x_1} - \overline{x_2}\right) = \dfrac{1}{\sigma\sqrt{1/n_1 + 1/n_2}\,\sqrt{2\pi}} \cdot \exp\left\{ -\dfrac{\left(\overline{x_1} - \overline{x_2} - \delta\right)^2}{2 \cdot \sigma^2 \cdot \left(1/n_1 + 1/n_2\right)} \right\}$$

The Bayes factor becomes

$$B = \frac{L\left(\delta = 0 \middle| \overline{x_1} - \overline{x_2}\right)}{\int_{-\infty}^{\infty} L\left(\delta \middle| \overline{x_1} - \overline{x_2}\right) \cdot p(\delta) d\delta} =$$

$$= \frac{\dfrac{1}{\sigma\sqrt{1/n_1 + 1/n_2}\sqrt{2\pi}} \cdot \exp\left\{\dfrac{\left(\overline{x_1} - \overline{x_2}\right)^2}{2 \cdot \sigma^2 \cdot \left(1/n_1 + 1/n_2\right)}\right\}}{\displaystyle\int_{-\infty}^{\infty} \dfrac{1}{\sigma\sqrt{1/n_1 + 1/n_2}\sqrt{2\pi}} \cdot \exp\left\{\dfrac{\left(\overline{x_1} - \overline{x_2} - \delta\right)^2}{2 \cdot \sigma^2 \cdot \left(1/n_1 + 1/n_2\right)}\right\} \cdot \dfrac{1}{\tau\sqrt{2}\sqrt{2\pi}} \cdot \exp\left(-\dfrac{\delta^2}{4\tau^2}\right) d\delta} =$$

$$= \frac{\exp\left\{\dfrac{\left(\overline{x_1} - \overline{x_2}\right)^2}{2 \cdot \sigma^2 \cdot \left(1/n_1 + 1/n_2\right)}\right\}}{\dfrac{1}{\tau\sqrt{2}\sqrt{2\pi}} \cdot \displaystyle\int_{-\infty}^{\infty} \exp\left\{\dfrac{\left(\overline{x_1} - \overline{x_2} - \delta\right)^2}{2 \cdot \sigma^2 \cdot \left(1/n_1 + 1/n_2\right)}\right\} \cdot \exp\left(-\dfrac{\delta^2}{4\tau^2}\right) d\delta}$$

A little bit easier to sort out. Can be shown to be

$$B = \sqrt{1 + \frac{2 \cdot n_1 \cdot n_2 \cdot \tau^2}{\sigma^2 \cdot (n_1 + n_2)}} \cdot \exp\left\{-\frac{1}{2} \cdot \left(\frac{1}{\sigma^2 \cdot (1/n_1 + 1/n_2)} - \frac{1}{\sigma^2 \cdot (1/n_1 + 1/n_2) + 2\tau^2}\right) \cdot \left(\overline{x_1} - \overline{x_2}\right)^2\right\}$$

Inserting the values of $\sigma$ $(= 0.1)$ and $\tau$ $(= n2.5)$ we obtain

$$B = \sqrt{1 + \frac{12.5 \cdot n_1 \cdot n_2}{0.01 \cdot (n_1 + n_2)}} \cdot \exp\left\{-\frac{1}{2} \cdot \left(\frac{1}{0.01 \cdot (1/n_1 + 1/n_2)} - \frac{1}{0.01 \cdot (1/n_1 + 1/n_2) + 12.5}\right) \cdot (\overline{x}_1 - \overline{x}_2)^2\right\}$$

Let's say we have obtained the sample means 18.1 % and 18.2 % and that our prior odds for $H_0$ are 1 (non-informative prior).

| $n_1$ | $n_2$ | $B$ | Posterior odds | Classical $P$-value |
|-------|-------|-------|----------------|---------------------|
| 2 | 2 | 21.46 | 21.46 | 0.3173 |
| 2 | 3 | 21.27 | 21.27 | 0.2733 |
| 5 | 5 | 16.03 | 16.03 | 0.1138 |
| 5 | 7 | 14.05 | 14.05 | 0.0877 |
| 10 | 10 | 6.49 | 6.49 | 0.0253 |
| 30 | 30 | 0.08 | 0.08 | 0.0001 |
| 100 | 100 | 0.00 | 0.00 | 0.0000 |

# Comparison of multivariate means

Assume there are two multivariate populations with means $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ respectively. The statistical model is that each population can be described in terms of a $p$-dimensional random vector

$$X = (X_1, \ldots, X_p)^{\mathrm{T}}$$

with a probability density function (or probability mass function) depending on the mean $\boldsymbol{\mu}$ : $f(\boldsymbol{x} \mid \boldsymbol{\mu}) = f(x_1, \ldots, x_p \mid \boldsymbol{\mu})$

The pair of hypotheses for comparing the two (vector) means is:

$$H_0: \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$$
$$H_1: \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$$

The random variation can be decomposed into (at least) two *levels*:

1. Within-variation: Accounts for the variability in $X$ given a specific value of $\mu$. This is modelled by $f(x \mid \mu)$.
2. Between-variation: Accounts for the uncertainty about the true value of $\mu$. This can be modelled in terms of a prior density for $\mu$, i.e. $p(\mu)$, $\mu \in \mathbf{M}$

This is referred to as a *two-level model* for the variation in which the variation at different levels is stochastic (cf. analysis of variance with random effects).

The Bayes factor for assessing the pair of hypothesis ($H_0$, $H_1$) from two sets of data, one from each population (*Data set* 1, *Data set* 2) is

$$B = \frac{\int_{\mathbf{M}} L(\mu|Data\ set\ 1) \cdot L(\mu|Data\ set\ 2) \cdot p(\mu) d\mu}{\int_{\mathbf{M}} L(\mu|Data\ set\ 1) \cdot p(\mu) d\mu \cdot \int_{\mathbf{M}} L(\mu|Data\ set\ 2) \cdot p(\mu) d\mu}$$

i.e. the same structure as for univariate data

When (vector) means are compared, sufficient statistics from the data sets are the two sample (vector) means:

*Data set* 1:

$$\boldsymbol{x}_1 = \left(\boldsymbol{x}_{11}, \ldots, \boldsymbol{x}_{1n_1}\right) \quad ; \quad \boldsymbol{X}_{1i} \sim f\left(\boldsymbol{x}_{1j} \middle| \boldsymbol{\mu}_1\right)$$

Sufficient statistic : $\overline{\boldsymbol{X}}_{1,n_1} \sim f\left(\overline{\boldsymbol{x}}_{1\cdot} \middle| \boldsymbol{\mu}_1\right)$

*Data set* 2:

$$\boldsymbol{x}_2 = \left(\boldsymbol{x}_{21}, \ldots, \boldsymbol{x}_{2n_2}\right) \quad ; \quad \boldsymbol{X}_{2i} \sim f\left(\boldsymbol{x}_{2j} \middle| \boldsymbol{\mu}_2\right)$$

Sufficient statistic : $\overline{\boldsymbol{X}}_{2,n_2} \sim f\left(\overline{\boldsymbol{x}}_{2\cdot} \middle| \boldsymbol{\mu}_2\right)$

$$\Rightarrow L\left(\boldsymbol{\mu}_i \middle| Data\ set\ i\right) = k\left(\boldsymbol{x}_{i1}, \ldots, \boldsymbol{x}_{in_i}\right) \cdot f\left(\overline{\boldsymbol{x}}_{i\cdot} \middle| \boldsymbol{\mu}_i\right)$$

This in turn gives

$$B = \frac{\int_{\boldsymbol{M}} f\left(\overline{\boldsymbol{x}}_1 \middle| \boldsymbol{\mu}\right) \cdot f\left(\overline{\boldsymbol{x}}_2 \middle| \boldsymbol{\mu}\right) \cdot p(\boldsymbol{\mu}) d\boldsymbol{\mu}}{\int_{\boldsymbol{M}} f\left(\overline{\boldsymbol{x}}_1 \middle| \boldsymbol{\mu}\right) \cdot p(\boldsymbol{\mu}) d\boldsymbol{\mu} \cdot \int_{\boldsymbol{M}} f\left(\overline{\boldsymbol{x}}_2 \middle| \boldsymbol{\mu}\right) \cdot p(\boldsymbol{\mu}) d\boldsymbol{\mu}}$$

Now, assume that both within- and between-variation is normal:

$$X_{ij}\big|\boldsymbol{\mu_i}, \boldsymbol{\Sigma}_i \sim N_p\left(\boldsymbol{\mu_i}; \boldsymbol{\Sigma}_i\right)$$

where $\boldsymbol{\Sigma}_i$ is the (variance-)covariance matrix.

This gives

$$\overline{X}_i \sim N_p\left(\boldsymbol{\mu_i}; n_i^{-1}\boldsymbol{\Sigma}_i\right) \implies$$

$$f\left(\overline{x}_i\big|\boldsymbol{\mu_i}\left[, \boldsymbol{\Sigma}_i\right]\right) = \left(2\pi \cdot \det\left(n_i^{-1}\boldsymbol{\Sigma}_i\right)\right)^{-p/2} \cdot \exp\left\{-\left(\overline{x}_i - \boldsymbol{\mu_i}\right)^{\mathrm{T}} \cdot n_i\boldsymbol{\Sigma}_i^{-1} \cdot \left(\overline{x}_i - \boldsymbol{\mu_i}\right)\right\}$$

*Full Bayesian approach:*

$$\boldsymbol{\mu}_i \sim N_p\left(\boldsymbol{\theta}_i; \boldsymbol{\Omega}_i\right)$$

Prior distribution is multivariate normal with hyper mean $\boldsymbol{\theta}_i$ and hyper covariance matrix $\boldsymbol{\Omega}_i$

$$\boldsymbol{\Sigma}_i \sim W_p^{-1}\left(d_i; \mathbf{T}^{-1}\right)$$

Prior distribution is p-dimensional Inverse Wishart with $d_i$ degrees of freedom and precision matrix $\mathbf{T}^{-1}$

*Simplified approaches*

Assume that the covariance matrix $\Sigma$ is common to both populations. This is a most reasonable assumption when $\Sigma$ represents measurement uncertainty.

It is also a reasonable assumption when the variation in a population does not depend substantially on the mean (*read normal distributions*) .

The comparability of the populations is self-understood – why would we otherwise compare the means?

Assume further that the prior distribution of $\mu$ is common to both populations – why could that be reasonable?

Instead of applying the full Bayesian approach with these simplifications we may use estimates from sufficient background data of

- $\Sigma$ (to avoid the use of an Inverse Wishart distribution)
- $\boldsymbol{\theta}$ and $\boldsymbol{\Omega}$ (*empirical Bayes* approach)

Assume we have a background data set of $N$ observed vectors from each of $M$ different populations

$$\{\boldsymbol{y}_{l,r}\}_{l=1,\dots,M;r=1,\dots,N} \quad \text{where}$$

$$\boldsymbol{y}_{l,r} = \left(y_{l,r,1},\dots,y_{l,r,p}\right)^{\mathrm{T}} \quad \text{is assumed to be } N_p\left(\boldsymbol{\mu};\Sigma\right) \text{ and}$$

$$\boldsymbol{\mu} \quad \text{is assumed to be } N_p\left(\boldsymbol{\theta};\boldsymbol{\Omega}\right)$$

An estimate of $\Sigma$ is then obtained as

$$\hat{\Sigma} = \frac{1}{(N-1) \cdot M} \sum_{l=1}^{M} \sum_{r=1}^{N} \left( y_{l,r} - \bar{y}_{l,\cdot} \right) \cdot \left( y_{l,r} - \bar{y}_{l,\cdot} \right)^{\mathrm{T}} = \frac{S_w}{(N-1) \cdot M}$$

The analogue of *SSE* (residual sum of squares) for multivariate data.

where $\bar{y}_{l,\cdot} = N^{-1} \sum_{r=1}^{N} y_{l,r}$

An estimate of $\Omega$ is

$$\hat{\Omega} = \frac{1}{M-1} \sum_{l=1}^{M} \left( \bar{y}_{l,\cdot} - \bar{y}_{\cdot,\cdot} \right) \cdot \left( \bar{y}_{l,\cdot} - \bar{y}_{\cdot,\cdot} \right)^{\mathrm{T}} - \frac{\hat{\Sigma}}{N} = \frac{S_*}{M-1} - \frac{S_w}{N^2(M-1)}$$

where $\bar{y}_{\cdot,\cdot} = (N \cdot M)^{-1} \sum_{l=1}^{M} \sum_{r=1}^{N} y_{l,r}$

…and an estimate of $\theta$ is

$$\hat{\theta} = \bar{y}_{\cdot,\cdot} = (N \cdot M)^{-1} \sum_{l=1}^{M} \sum_{r=1}^{N} y_{l,r}$$

For obtaining the Bayes factor

$$B = \frac{\int_{\mathbf{M}} f(\bar{x}_1|\mu) \cdot f(\bar{x}_2|\mu) \cdot p(\mu) d\mu}{\int_{\mathbf{M}} f(\bar{x}_1|\mu) \cdot p(\mu) d\mu \cdot \int_{\mathbf{M}} f(\bar{x}_2|\mu) \cdot p(\mu) d\mu}$$

this implies

$$\bar{X}_i \text{ is assumed to be } N_p\left(\mu; n_i^{-1}\hat{\Sigma}\right) \Rightarrow$$

$$f(\bar{x}_i|\mu) = \left(2\pi \cdot \det\left(n_i^{-1}\hat{\Sigma}\right)\right)^{-p/2} \cdot \exp\left\{-(\bar{x}_i - \mu)^{\mathrm{T}} \cdot n_i\hat{\Sigma}^{-1} \cdot (\bar{x}_i - \mu)\right\}$$

and

$$\mu \text{ is assumed to be } N_p\left(\hat{\theta}; \hat{\Omega}\right) \Rightarrow$$

$$p(\mu) = \left(2\pi \cdot \det\left(\hat{\Omega}\right)\right)^{-p/2} \cdot \exp\left\{-(\mu - \hat{\theta})^{\mathrm{T}} \cdot \hat{\Omega}^{-1} \cdot (\mu - \hat{\theta})\right\}$$

Aitken CGG & Lucy D (2004). Evaluation of trace evidence in the form of multivariate data. *Applied Statistics* **53**(1): 109-122

– use the full likelihoods (i.e. not the densities of the sufficient means) in

$$B = \frac{\int_{\mathbf{M}} L(\boldsymbol{\mu}|Data\ set\ 1)\cdot L(\boldsymbol{\mu}|Data\ set\ 2)\cdot p(\boldsymbol{\mu})d\boldsymbol{\mu}}{\int_{\mathbf{M}} L(\boldsymbol{\mu}|Data\ set\ 1)\cdot p(\boldsymbol{\mu})d\boldsymbol{\mu} \cdot \int_{\mathbf{M}} L(\boldsymbol{\mu}|Data\ set\ 2)\cdot p(\boldsymbol{\mu})d\boldsymbol{\mu}}$$

i.e.

$$L(\boldsymbol{\mu}|Data\ set\ i) = \prod_{j=1}^{n_i} f(\boldsymbol{x}_{ij}|\boldsymbol{\mu}) = \prod_{j=1}^{n_i} (2\pi\cdot\det(\hat{\boldsymbol{\Sigma}}))^{-p/2}\cdot\exp\left\{-(\boldsymbol{x}_{ij}-\boldsymbol{\mu})^{\mathrm{T}}\cdot\hat{\boldsymbol{\Sigma}}^{-1}\cdot(\boldsymbol{x}_{ij}-\boldsymbol{\mu})\right\}$$

– derive explicit formulas for the Bayes factor when
  • prior distribution of $\boldsymbol{\mu}$ is normal
  • prior distribution of $\boldsymbol{\mu}$ is estimated by a kernel density

*From the "real" world*

# Comparison of amphetamine seizures

- Issue:
  - To evaluate findings from two seizures of amphetamine against the propositions

  > $H_m$ : The two seizures originate from the same precipitation batch
  > $H_a$ : The two seizures originate from different precipitation batches

- Framework constraint:
  - The findings should be based on impurity profiling for a set of 30 impurities agreed on in a European cooperation program

- Problem:
  - Today's "standard of evaluation" is manual inspection of profiles – time consuming, especially since in a case there are usually pairwise comparisons of several seizures

# The production of amphetamine

- Choose a recipe

- Produce amphetamine oil

- Precipitate the amphetamine powder ⇒ ***precipitation batch***

# *Impurities (by-products) come with production:*

Diluents: Sugar

Adulterants:
Caffeine, Phenazone,
1-phenylethylamine

Amphetamine
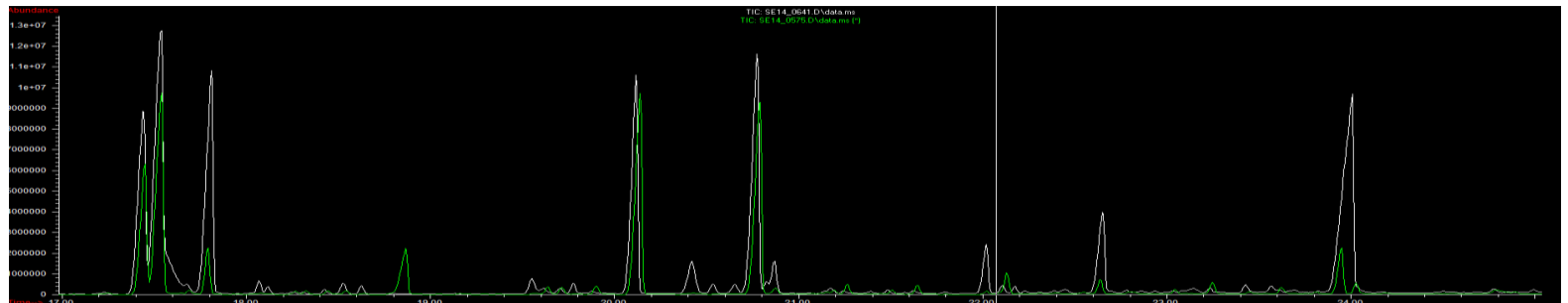
By-products

# Manual inspection/comparison of impurity profiles

Inspecting overlaid *chromatograms*

**Green** and **violet** peaks representing equal or different profiles?

I.



II.

# Evaluative statement

The National Forensic Centre (NFC) uses a common (verbal) scale of conclusions based on intervals of likelihood ratios (LR) [or Bayes factors].

"The findings from the comparison of the two seizures...

| | |
|---|---|
| *...support extremely strongly...* | $(10^6 \leq LR\,)$ |
| *...support strongly...* | $(6000 \leq LR < 10^6\,)$ |
| *...support...* | $(100 \leq LR < 6000\,)$ |
| *...support to some extent...* | $(6 \leq LR < 100\,)$ |
| ... that the two seizures originate from the same precipitation batch | |
| *...support neither* that the two seizures originate from the same precipitation batch *nor* that they originate from different precipitation batches | $(1/6 < LR < 6\,)$ |
| *...support extremely strongly/strongly//to some extent* that the two seizures originate from different precipitation batches | $LR \leq 1/6$ , $1/100$, $1/6000$, $1/10^6$ |

Nordgaard et al, *Law, probability and Risk* , 2012

# How do we know that the correct magnitude of the likelihood ratio has been assigned?

Always address the likelihood ratio when evaluating:

$$LR = \frac{\Pr(E|H_m)}{\Pr(E|H_a)}$$

with $E$ representing the findings from the inspection/comparison

The probabilities $\Pr(E|H_m)$ and $\Pr(E|H_a)$ are assigned from experience and knowledge…

…but are by necessity objected to a substantial amount of subjectivity.

Manual comparisons should always be made, but it would be good to have the evaluation assisted by an objective probabilistic model.

# Applying a probabilistic model to the findings

Monitoring 30 peaks in a chromatogram: *Multivariate inference*

*In theory:* Model the multivariate distributions of the ***peak areas***
under each of the propositions and compute the Bayes factor for
the comparison:

With $x$ = observed vector of peak areas for seizure 1

$y$ = observed vector of peak areas for seizure 2

$\theta$ = true mean vector of peak areas

$p(\theta)$ = prior density for $\theta$

$$B = \frac{\int f(x|\theta) \cdot f(y|\theta) \cdot p(\theta)d\theta}{\int f(x|\theta)p(\theta)d\theta \times \int f(y|\theta)p(\theta)d\theta}$$

# Simplifying…

*Assume…*

- the peak area variation between samples from the same precipitation batch – the *within-variation* – is normal

- the peak area variation between samples from different precipitation batches – the *between-variation* – is not (necessarily) normal

*Graphical description of experimental data (see below) do not contradict these assumptions.*

With these assumptions Aitken & Lucy (*Applied Statistics*, 2004) give explicit functions for the numerator and denominator of the Bayes factor.

$$B = \frac{f_m(x, y | \mu, U, C)}{f_a(x, y | \mu, U, C)}$$

where $f_m$ and $f_a$ are density functions obtained by integrating multivariate normal distributions of $x$ and $y$ with a kernel estimate of the prior density of $\theta$.

### *The Multivariate Kernel Likelihood Ratio (MVK)*

$B$ will then depend on
* the mean vector of the prior distribution of peak areas: $\mu$
* the variance-covariance matrix for the within-variation: $U$
* the variance-covariance matrix for the between-variation: $C$

*Taking these as hyperparameters we need stable estimates!*

*Could also be assigned prior distributions, but that is not done here.*

# Experimental study

- Five different recipes were chosen

- Five different oils were produced with each recipe

- Three precipitation batches per oil were taken

- *One precipitation was not successful*

$\Rightarrow$ In total 74 different precipitation batches

From each precipitation batch 4-9 replicate samples were taken for profiling.

4-5 replicates were taken when the batch was fresh. For a number of batches another 4 replicates were taken upon storage (in freezer) and airing.

Ignoring the full hierarchical structure of the data…

Two levels of variation:

- within precipitation (between replicates from the same batch)
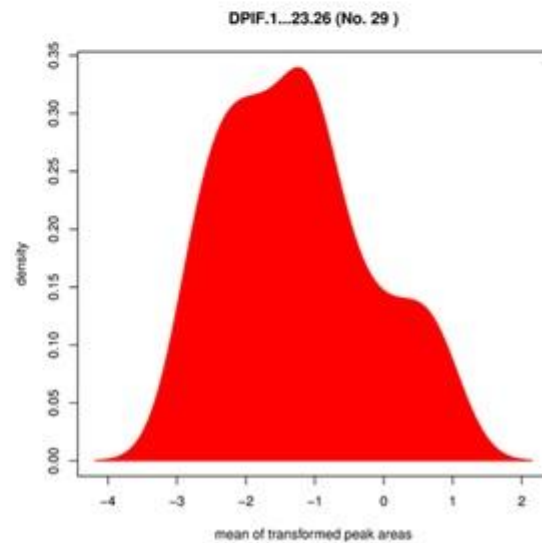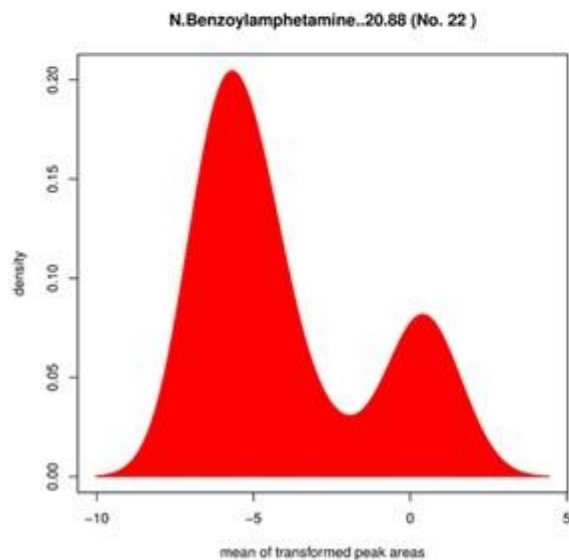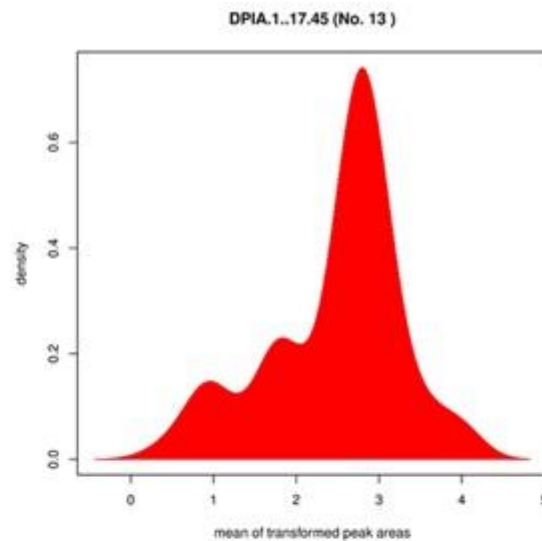- between precipitations (sum of variance components from recipe, oil and precipitation)
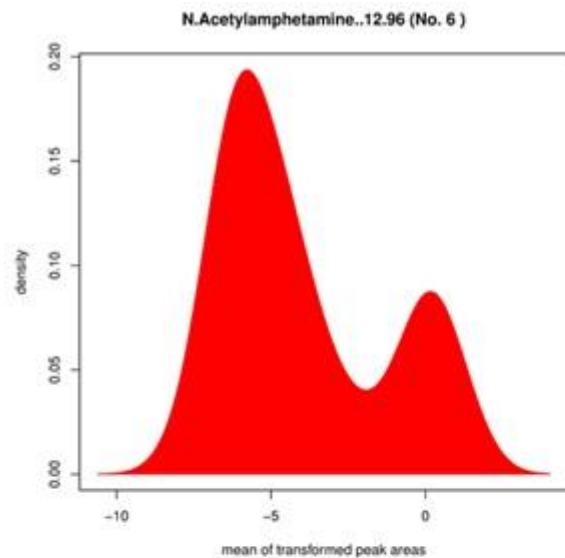
# Harmonising and variance stabilizing transformations:

- When profiling using GC/MS an *internal standard* is run in each sequence

- Each sample is diluted before analysis so that the linear range of the instrument is not exceeded. How much depends on the content of by-products in each sample. The degree of dilution is represented as a *sample multiplier* (factor)

- Moreover, the dry concentrations of the precipitation batches will vary substantially

To account for these observable sources of variation and to allay the effect of outliers (w.r.t. assumed normal within-variation) the original peak area (*PA*) is transformed as

$$pa = \log\left( \frac{PA}{PA_{\text{internal standard}}} \times \frac{\text{sample multiplier}}{\text{dry concentration}} \right)$$

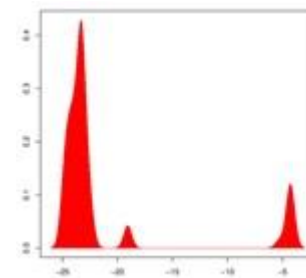# Between-variation for some of the impurities:

# Data reduction

With 74 different batches and sometimes only 4 replicate measurements on each it is not numerically feasible to estimate variance-covariance matrices for a 30-dimensional vector.

Several of the impurities could however be discarded by "expert reasons":

- Some impurities are known to be volatile
- In experimental (and casework) data some impurities show very low entropy (are seldom detected)

zero peak area

Upon such considerations 18 impurities remain – still too many!

# Solution used here:

Expert choice of 12 impurities, from experience known to vary between seizures indisputably from different sources.

N.Acetylamphetamine

N.Formylamphetamine

Benzylamphetamine

DPIA.1

alfa.Methyldiphenetyletylamine

N.Benzoylamphetamine

Unknown.B2

X2.6.Dimethyl.3.5.diphenylpyridine

X2.4.Dimetyl.3.5.diphenylpyridine

Pyridine.7.and.14

X2.6.Diphenyl.3.4.dimethylpyridine

DPIF.1

*Of interest primarily for chemists and people involved in the profiling cooperation projects*

Works numerically, but stable estimates may still be hard to obtain.

# Investigating the stability by resampling

**Resample** the experimental data by drawing batch samples with replacement within batches.

Compute estimates (MLE) of $U$ and $C$ (from original data and resampled data).

$$\Rightarrow \hat{U}, \hat{C} \text{ and } \{U^*\}, \{C^*\}$$

Assess the traces and determinants of $\{U^*\}$ and $\{C^*\}$ with respect to relative standard deviation over the resamples for three size of set of impurities:

     (i) 3 impurities   (ii) 6 impurities  (iii) 12 impurities (as was used)

Results for 500 resamples:

|  | 3 impurities | 6 impurities | 12 impurities |
|---|---|---|---|
| **tr(U), %σ** | 4.9 | 4.8 | 6.0 |
| **det(U), %σ** | 18.4 | 19.4 | 21.9 |
| **tr(C), %σ** | 1.3 | 1.7 | 1.5 |
| **det(C), %σ** | 5.1 | 5.5 | 9.2 |

# Does it work for the evaluation?

How can we assess whether the suggested method works or not?

*In general:*
- The false positive rate should be very low (*Do we accept non-zero rates?*)
- The false negative rate should be very low (*Needs to be of the same magnitude as the false positive rate?*)

"It is better that ten guilty escape than one innocent suffer."

*Lord William Blackstone*

- Evidence strongly supporting a proposition should generate a large Bayes factor and vice versa

Recall the NFC scale of conclusions:

| Level | Interval of B | Level | Interval of B |
|---|---|---|---|
| +4 | $10^6 \leq B$ | −1 | $1/100 < B \leq 1/6$ |
| +3 | $6000 \leq B < 10^6$ | −2 | $1/6000 < B \leq 1/100$ |
| +2 | $100 \leq B < 6000$ | −3 | $1/10^6 < B \leq 1/6000$ |
| +1 | $6 \leq B < 100$ | −4 | $B \leq 1/10^6$ |
| 0 | $1/6 < B < 6$ | | |

Hence, what matters here is
- whether a level at the "correct" side of the scale is attained **<u>and</u>**
- whether findings on seizures with *a common source* generally lead to high levels (+3, +4); and whether findings on seizures with *different sources* generally lead to low levels ( −3, −4)

# Cross-validatory study:

$$B = \frac{f_m(\boldsymbol{x}, \boldsymbol{y} | \boldsymbol{\mu}, U, C)}{f_a(\boldsymbol{x}, \boldsymbol{y} | \boldsymbol{\mu}, U, C)}$$

1. (a) Leave out one precipitation batch at a time, estimate $\boldsymbol{\mu}$, $U$ and $C$ from the rest
   (b) Split the replicates of the left-out batch into two "seizures with common origin" (use different combinations for the split)
   (c) Apply the estimated "model" to each combination and compute the median Bayes factor

$\Rightarrow$ 74 values of $B$ for assessing the false negative rate and the "strength consistency" for a common origin.

2. (a) Leave out two precipitation batches at a time (use all combinations) to be two "seizures" with different origins, estimate $\boldsymbol{\mu}$, $U$ and $C$ from the rest
   (b) Apply the "estimated model" to the two seizures, compute the Bayes factor

$\Rightarrow$ 2701 values of $B$ for assessing the false positive rate and the "strength consistency" for different origins.

# Results:

## *Assessment of the false negative rate*

For the 74 (median) values of the Bayes factor, *B* for the proposition of a common origin :

| Scale level | −4 | −3 | −2 | −1 | 0 | +1 | +2 | +3 | +4 |
|---|---|---|---|---|---|---|---|---|---|
| *Frequency* | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 28 | 44 |

Hence, the false negative rate is 1/74 ≈ 1.4 %

Using the average Bayes factor instead of the median renders level +4 only

# Assessment of the false positive rate

For the 2701 values of the Bayes factor, $B$ for the proposition of a common origin :

| Scale level | −4 | −3 | −2 | −1 | 0 | +1 | +2 | +3 | +4 |
|---|---|---|---|---|---|---|---|---|---|
| Frequency | 2633 | 4 | 3 | 0 | 2 | 4 | 4 | 12 | 39 |

For the two occurrences of level 0 the Bayes factor is 0.39 and 0.34 respectively.

Hence, the false positive rate is 59/2701 ≈ 2.2 %

*Disturbing:* The false positive results are dominated by high scale levels (+3 and +4)

Further inspection of the comparisons rendering levels +3 and +4 shows that

- all but 3 are comparisons of batches from the same oil (the 3 concerns two oils from the same recipe)

- the Bayes factor for the comparison is mostly of a lower magnitude when comparing samples that were stored and aired compared to the corresponding comparison of fresh samples – in a few cases the level shifted from +3 to –2

- manual comparisons of the chromatograms
  - also resulted mostly in high scale levels (+3, +4) for the fresh samples
  - resulted in lower levels (from –2 to +3) for stored and aired samples

***Assessment of the strength consistency***

## Empirical cross-entropy, *ECE*

In the cross-validatory study, let

$S_m$ = The set of comparisons of samples with common origin
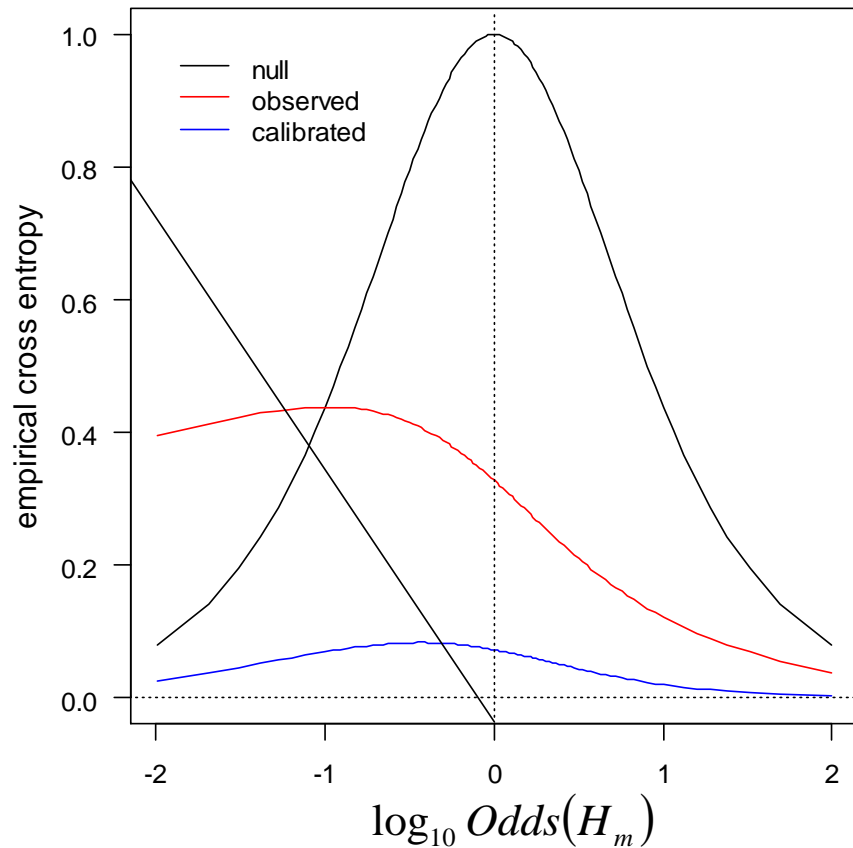$N_m$ = Number of comparisons of samples with common origin (i.e. 74)
$S_a$ = The set of comparisons of samples with different origins
$N_a$ = Number of comparisons of samples with different origins (i.e. 2701)

$$ECE = \frac{\Pr(H_m)}{N_m} \sum_{i \in S_m} \log_2 \left( 1 + \frac{1}{B_i \cdot Odds(H_m)} \right) +$$

$$+ \frac{\Pr(H_a)}{N_a} \sum_{j \in S_a} \log_2 \left( 1 + B_j \cdot Odds(H_a) \right)$$

# *ECE* plot    (provided by R-package *comparison*)



The closer the **red** curve (observed *ECE*) is to the **blue** curve (calibrated *ECE*) the more robust are the Bayes factors that are calculated this way.

For prior odds where the **red** curve exceeds the **black** curve (representing *ECE* for non-informative evidence – *B* is always = 1), the Bayes factors can be misleading.

Hence, the Bayes factors produced with the 12 selected impurities can be misleading if the prior odds are below 0.1.

# References

Aitken C.G.G, Lucy D. (2004). Evaluation of trace evidence in the form of multivariate data. *Appl.Statist.* 53(1): 109-122.

Lindley D.V. (1977). A problem in forensic science. *Biometrika* 64(2): 207-213.

Nordgaard A, Ansell R., Drotz W., Jaeger L. (2012). Scale of conclusions for the value of evidence. *Law, probability and Risk* 11(1): 1-24.

Ramos D., Gonzales-Rodriguez J., Zadora G., Aitken C. (2013). Information-Theoretical Assessment of Likelihood Ratio Computation Methods. *J.Forensic Sci* 58(6): 1503-1518.