**MATH60604A *Statistical Modelling***

**Midterm examination**

Exam booklet

Practice exam

Léo Belzile

**Instructions**: The time allotted for the examination is 180 minutes. A non-programmable calculator may be used. You are allowed a single-sided crib sheet.

There are a total of 40 marks available in the exam paper, the distribution of which can be found in the right margin.

You must hand back the **exam booklet** at the end of the examination.

Last name:

First name:

STUDENT ID:

| Question: | 1 | 2 | 3 | Total |
|-----------|----|----|----|-------|
| Points: | 12 | 12 | 16 | 40 |
| Score: | | | | |

**Question 1.** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .    $\boxed{12}$

In extreme value analysis, large-sample theory dictates that exceedances of $Y$ above a high threshold $u$ are well approximated by a **generalized Pareto** distribution $Z = Y - u \dot{\sim} \mathrm{GP}(\sigma, \xi)$ with scale $\sigma > 0$ and shape $\xi \in \mathbb{R}$, with distribution and density functions

$$F(z) = 1 - (1 + \xi z/\sigma)_+^{-1/\xi}, \qquad f(z) = \sigma^{-1}(1 + \xi z/\sigma)_+^{-1/\xi - 1},$$

with $(x)_+ = \max\{x, 0\}$; the case $\xi = 0$ is defined by continuity (exponential submodel).

We consider the largest fire insurance claims (in millions Danish krones) filed at the Copenhagen Re company in Denmark from January 1980 until the end of December 1990 ($n_Y = 11$ years worth of data). We model the $n = 109$ exceedances above 10 millions krones, corresponding to $\zeta = 0.0503$ proportion of the data. Our objective is to provide 100-years return levels for risk analysis.

The $T$-year return level $r(T)$, is a high quantile exceeded with tail probability $p$, where $p = \zeta n_Y / T$, with $\zeta$ the proportion of observations above the threshold, $n_Y$ the number of years worth of data and $T = 100$ the number of year. If we invert the distribution function, we find

$$r(T) = \frac{\sigma}{\xi}\left\{(\zeta n_Y / T)^{-\xi} - 1\right\}$$

1.1 Write the log likelihood function for an $n$-sample of independent threshold exceedances $z_i$,    [2]
$(i = 1, \ldots, n)$ for $\xi \neq 0$.

1.2 If we reparametrize the model in terms of $\xi$ and $\theta = \xi/\sigma$, show that we can derive an explicit formula for the profile log likelihood $\ell_{\mathrm{p}}(\theta)$, thereby reducing the optimization to a one-dimensional problem. [2]

```
1 Log-likelihood: -374.893
2 Sample size: 109
3 Proportion above threshold: 0.0503
4
5 Estimates                Standard errors
6 scale   shape            scale    shape
7 6.975   0.497            1.1135   0.1363
```
**Listing 1:** Maximum likelihood estimate for the generalized Pareto

1.3  An optimization routine has yielded the following estimates for the MLE in Listing 1. Ex-     [2]
plain how you could test whether $\xi = 0$ (exponential model) using (a) a Wald test using the
output from Listing 1 and (b) a likelihood ratio test if in addition

```
1 > sum(dexp(y, rate = 1/mean(y), log = TRUE))
2 -397.2921
```
**Listing 2:** Code for an exponential log likelihood

1.4  One can show that the Fisher information matrix is     [2]

$$\iota(\sigma,\xi) = n\begin{pmatrix} \sigma^{-2}(1+2\xi)^{-1} & \sigma^{-1}(1+\xi)^{-1}(1+2\xi)^{-1} \\ \sigma^{-1}(1+\xi)^{-1}(1+2\xi)^{-1} & 2(1+\xi)^{-1}(1+2\xi)^{-1} \end{pmatrix}$$

Explain how you could use the above result to derive standard errors for the parameters $\sigma$
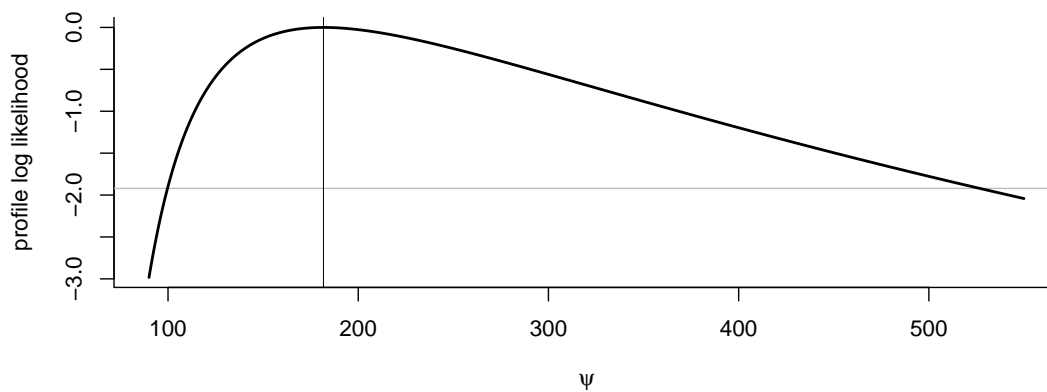and $\xi$.

**Figure 1:** Profile log likelihood for the return level. The profile has been shifted to equal zero at the MLE, and the gray horizontal line indicates the cutoff for 95% profile confidence intervals based on the asymptotic $\chi_1^2$ distribution.

1.5 Give the maximum likelihood estimate of the 100 years return level for the Danish insurance data. [2]

1.6 Figure 1 shows the profile log likelihood for the 100-year return level, or value-at-risk at level $p = 1/100$, with the gray horizontal line indicating the cutoff for 95% profile confidence intervals based on the asymptotic $\chi_1^2$ distribution. Would a Wald-based confidence intervals be similar looking? [2]

**Question 2.** ............................................          12

Grossmann and Kross (2014) study the question "Are people wiser when reflecting on other people's problems compared with their own?". They "randomly assigned participants to

1. reason about their own problem from an immersed perspective (self-immersed condition),
2. reason about their friend's problem from an immersed perspective (other-immersed condition),
3. reason about their own problem from a distanced perspective (self-distanced condition), or
4. reason about their friend's problem from a distanced perspective (other-distanced condition)".

The study below also considered age as a separate factor.

The variables include

- limits: response variable, the mean-centered score for the question on the "recognition of limits of knowledge".
- target: factor, is the target self or other.
- perspective: factor, either immersed or distanced.
- age: age group, either young (20–40 years old) or old, (60–80 years old).

**Table 1:** Coefficients and standard errors for the three-way full factorial model.

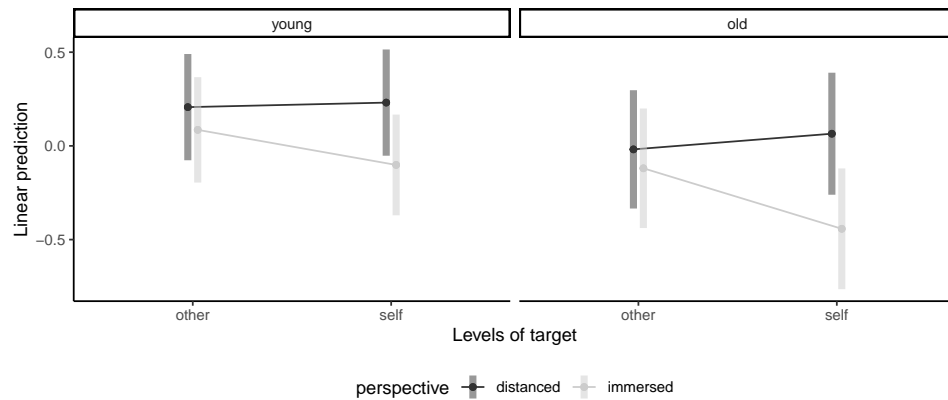|  | coef. | std. error |
|---|---|---|
| (Intercept) | 0.207 | 0.144 |
| target [self] | 0.024 | 0.204 |
| perspective [immersed] | −0.121 | 0.203 |
| age [old] | −0.225 | 0.216 |
| target [self]:perspective [immersed] | −0.211 | 0.284 |
| target [self]:age [old] | 0.059 | 0.308 |
| perspective [immersed]:age [old] | 0.020 | 0.306 |
| target [self]:perspective [immersed]:age [old] | −0.195 | 0.433 |

**Figure 2:** Estimated marginal means with 95% confidence intervals for each subgroup

**Table 2:** Analysis of variance table (type 2 decomposition) for the full factorial three-way ANOVA model.

|  | sum of squares | df | $F$-stat | $p$-value |
|---|---|---|---|---|
| target | 1.12 | 1 | 0.87 | 0.35 |
| perspective | 7.59 | 1 | 5.88 | 0.02 |
| age | 6.12 | 1 | 4.74 | 0.03 |
| target:perspective | 2.45 | 1 | 1.90 | 0.17 |
| target:age | 0.04 | 1 | 0.03 | 0.86 |
| perspective:age | 0.16 | 1 | 0.13 | 0.72 |
| target:perspective:age | 0.26 | 1 | 0.20 | 0.65 |
| residuals | 570.17 | 442 |  |  |

**Table 3:** Estimated marginal means for the four experimental conditions.

| perspective | target | emmean | std. error | lower CL | upper CL |
|---|---|---|---|---|---|
| distanced | other | 0.09 | 0.11 | −0.12 | 0.31 |
| immersed | other | −0.02 | 0.11 | −0.23 | 0.20 |
| distanced | self | 0.15 | 0.11 | −0.07 | 0.36 |
| immersed | self | −0.27 | 0.11 | −0.48 | −0.06 |

**Table 4:** Marginal contrasts for the two-way marginal model.

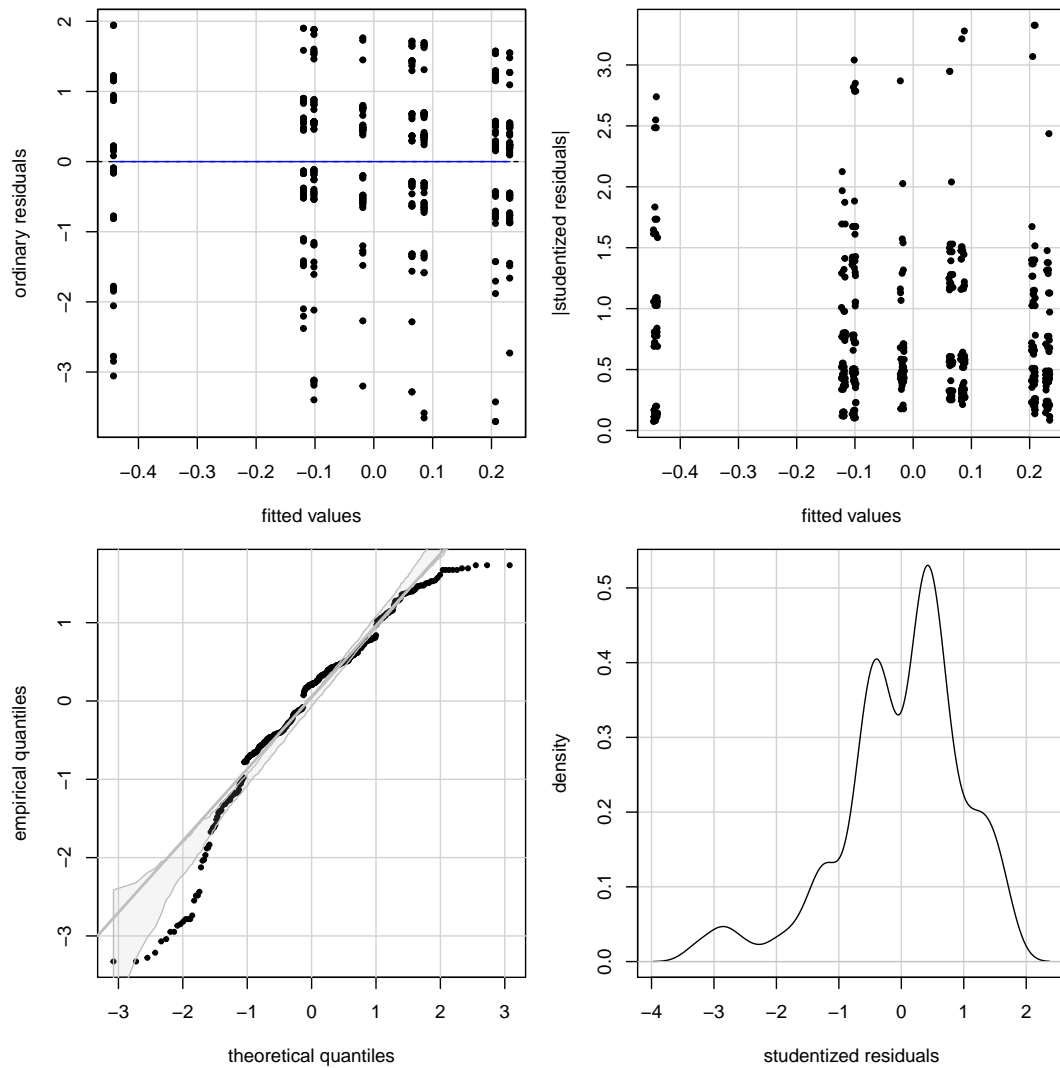| contrast | estimate | std. error | $t$-stat | $p$-value |
|---|---|---|---|---|
| $C_1$ | 0.311 | 0.131 | 2.366 | 0.018 |
| $C_2$ | −0.109 | 0.134 | −0.818 | 0.414 |
| $C_3$ | 0.420 | 0.153 | 2.742 | 0.006 |
| $C_4$ | 0.111 | 0.153 | 0.727 | 0.468 |

**Figure 3:** Diagnostic plots for the three-way ANOVA model: scatterplots of fitted versus ordinary residuals (top left), fitted vs absolute value of externally studentized residuals (top right), Student quantile-quantile plot of externally studentized residuals (bottom left) and density plot of the externally studentized residuals (bottom right).

2.1  Interpret the intercept parameter of the full factorial model fitted in Table 1.          [2]

2.2  Based on Table 2, how many participants were recruited in the study?          [2]

2.3  Based on the analysis of variance reported in Table 2, is it correct to marginalize over age          [2]
and consider the main effects of (target, perspective) by pooling data for both young and
old? What is the benefit of doing so?

2.4  Write down the weight vectors for the following four contrasts, given the order of the cate-          [2]
gories in Table 3:
$C_1$  (_____, _____, _____, _____) other (immersed and distance) vs self-immersed
$C_2$  (_____, _____, _____, _____) other (immersed and distanced) vs self-distanced
$C_3$  (_____, _____, _____, _____) self-distanced vs self-immersed
$C_4$  (_____, _____, _____, _____) other-distanced vs other-immersed

2.5 Based on the output of Table 4, comment on the contrast analysis.    [2]

2.6 Based on Figure 3, what model assumption is violated? Justify your answer and discuss the impacts on inference.    [2]

- independence
- incorrect mean model specification
- additivity
- homoscedasticity
- normality
- absence of outliers

**Question 3.** . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . ⟨16⟩

    **Think-aloud and the impact of learning to read.** The data from Baumann *et al.* (1992) study the effect of different instruction methods for reading comprehension. The data we consider are the following

- `pretest2`: score (out of 15) on second pretest using a comprehension monitoring questionnaire
- `posttest2`: response, score (out of 18) on an expanded comprehension monitoring questionnaire
- `group`: experimental group, one of directed reading-thinking activity (DRTA), think-aloud (TA) and directed reading group (DR)

We fit the models for the `posttest2` score with the **sum-to-zero** constraint for the categorical variable experimental `group`. The model matrix for the dummies is given in Table 5

|      | (Intercept) | group1 | group2 |
|------|:-----------:|:------:|:------:|
| DR   | 1           | 1      | 0      |
| DRTA | 1           | 0      | 1      |
| TA   | 1           | −1     | −1     |

**Table 5:** Model matrix with the sum-to-zero constraint.

```
1 > model1 <- lm(posttest2 ~ group, data = BSJ92)
2 > model2 <- lm(posttest2 ~ offset(pretest2) + group, data = BSJ92)
3 > model3 <- lm(posttest2 ~ pretest2 + group, data = BSJ92)
4 > model4 <- lm(posttest2 ~ pretest2 * group, data = BSJ92)
```

**Listing 3: R** syntax for the four fitted models

```
1 > anova(model3, model4)
2 Analysis of Variance Table
3
4 Model 3: posttest2 ~ pretest2 + group
5 Model 4: posttest2 ~ pretest2 * group
6   Res.Df    RSS Df Sum of Sq      F Pr(>F)
7 1     62 332.19
8 2     60 331.07  2    1.1264 0.1021 0.9031
```

**Listing 4:** Comparison of two models

**Table 6:** Coefficients and standard errors for different linear regression models.

|  | estimate | std. error |
|---|---|---|
| (Intercept) | 6.712 | 0.293 |
| group1 | −1.167 | 0.414 |
| group2 | −0.485 | 0.414 |

**(a)** Coefficients for Model 1 for the post-test score 2 as a function of the experimental group (with sum-to-zero parametrization).

|  | estimate | std. error |
|---|---|---|
| (Intercept) | 5.301 | 0.722 |
| pretest2 | 0.276 | 0.130 |
| group1 | −1.213 | 0.404 |
| group2 | −0.481 | 0.403 |

**(b)** Coefficients for Model 2 for the post-test score 2 as a function of the experimental group (with sum-to-zero parametrization) and `pretest2`.

|  | estimate | std. error |
|---|---|---|
| (Intercept) | 5.398 | 0.765 |
| pretest2 | 0.256 | 0.140 |
| group1 | −1.619 | 0.994 |
| group2 | −0.236 | 1.113 |
| pretest2:group1 | 0.079 | 0.176 |
| pretest2:group2 | −0.047 | 0.204 |

**(c)** Coefficients for Model 3 for the post-test score 2 as a function of the experimental group (with sum-to-zero parametrization), `pretest2` and their interaction.

**Table 7:** Estimated marginal means for Model 3.

| group | emmean | std. error | df | lower CL | upper CL |
|---|---|---|---|---|---|
| DR | 5.499 | 0.494 | 62 | 4.512 | 6.487 |
| DRTA | 6.231 | 0.494 | 62 | 5.245 | 7.218 |
| TA | 8.406 | 0.494 | 62 | 7.418 | 9.393 |

**Table 8:** Pairwise contrasts based on marginal means of Model 3.

| contrast | estimate | std. error | df | $t$-stat. | $p$-value |
|---|---|---|---|---|---|
| DR − DRTA | -0.732 | 0.698 | 62 | -1.048 | 0.550 |
| DR − TA | -2.906 | 0.699 | 62 | -4.157 | $< 10^{-3}$ |
| DRTA − TA | -2.174 | 0.698 | 62 | -3.114 | 0.008 |

3.1  Compute the sample mean of each group based on Table 6. [2]

3.2  It is common (but often incorrect) to fit an analysis of variance model for the difference [2]
post/pre, i.e., posttest2−pretest2 (Model 2) rather than fitting the linear regression as in
Model 3. This is equivalent to using an offset, i.e., a covariate with a known coefficient (or
$\beta = 1$). Does the data support this hypothesis?

3.3  The authors compared different models. Which of the above four models from Listing 3 are [2]
nested?

3.4 Write down the mean for each group of Model 4 and show that Model 3 is a simplification [2] of the latter. Write the null and alternative hypotheses in terms of model parameter and conclude based on the output of Listing 4.

3.5 The authors computed the marginal mean from Model 3 (Table 7), and the pairwise con- [2] trasts, reported in Table 8. Based on these, can you conclude about a ranking as to which treatment is the most effective (if higher scores are better)?

3.6 If the slopes for `pretest2` for each `group` are not parallel, explain why the comparison   [2]
    using marginal means per experimental groups are hazardous.

3.7 The authors report Levene's test of homogeneity of variance,   [2]

```
1 > car::leveneTest(rstudent(model3) ~ group,
2 +                 data = BSJ92,
3 +                 center = "mean")
```
**Listing 5: R** call for Levene's test of homogeneity of variance

returns a table with the statistic $F(2,63) = 1.51$, and a $p$-value of $p = 0.23$. What is the pur-
pose of this test, conclude and explain how this impacts your conclusion, if at all.

3.8 Given that the `pretest2` and `posttest2` are correlated, does this indicate a violation of the   [2]
    independence assumption? Discuss.