

Minimum Dispair & Maximum Despair: Minimum Disparity Statistics



Nils Lid Hjort

Department of Mathematics, University of Oslo

Building Bridges at Bislett
May 2017, Teknologihuset

Minimum disparity, i.i.d. case

Suppose y_1, \dots, y_n i.i.d. from g , and we use model $f_\theta(\cdot) = f(\cdot, \theta)$. With $B(\cdot)$ convex and smooth on the half-line, and with $B(1) = 0$,

$$d(g, f_\theta) = \int B\left(\frac{g}{f_\theta}\right) f_\theta \, dy$$

is a **divergence** (Johan Jensen inequality, 1906). Note that **dim ≥ 2 is ok.**

Hence it's a **good estimation idea** to use

$$\hat{\theta} = \operatorname{argmin} d(\tilde{g}, f_\theta),$$

where \tilde{g} is a 'wider' estimate of g . Can use

- ▶ \tilde{g} **nonparametric**, like the kernel density estimator **[99% of all MinDisp papers use this]**, or e.g. Hjort–Glad $f(y, \tilde{\theta})\tilde{r}(y)$;
- ▶ \tilde{g} via **Bayesian nonparametrics** (e.g. along with prior for θ);
- ▶ \tilde{g} via $g = (1 - \varepsilon)f_\theta + \varepsilon H$ modelling;
- ▶ $\tilde{g}(y) = g(y, \tilde{\theta}, \tilde{\gamma})$ from a **bigger parametric** family.

Minimising $\int B(g/f_\theta) f_\theta \, dy$ corresponds to solving

$$\int A\left(\frac{g}{f_\theta}\right) f_\theta u_\theta \, dy = 0,$$

where $u_\theta = \partial \log f_\theta / \partial \theta$ is score function and

$$A(\rho) = \rho B'(\rho) - B(\rho)$$

is the **Ratio Adjustment Function** (RAF). So, $\hat{\theta}$ solves

$$H_n(\theta) = \int A\left(\frac{\tilde{g}}{f_\theta}\right) f_\theta u_\theta \, dy = 0.$$

We have $A'(\rho) = B''(\rho)$, so $A(\rho)$ is increasing.

Properties of $A(\rho)$ drive robustness properties of $\hat{\theta}$.

A few Min Dispair schemes

1. $B(\rho) = \rho \log \rho$: then $d = \text{KL}(g, f_\theta)$, and $A(\rho) = \rho$. Estimator: solution to $\int \tilde{g}(y) u(y, \theta) dy = 0$... famous divergence, efficient at model, close to ML (inside and outside model), and not robust.

2. $B(\rho) = -\log \rho$: then $d = \text{KL}(f_\theta, g)$, the converse of what we learn in school, and $A(\rho) = \log \rho - 1$. Estimator: solution to

$$H_n(\theta) = \int \log \frac{\tilde{g}(y)}{f(y, \theta)} f(y, \theta) u(y, \theta) dy = 0.$$

3. $B(\rho) = 1 - \sqrt{\rho}$: then $d = 1 - \int \sqrt{g f_\theta} dy$, the Hellinger distance; $A(\rho) = \frac{1}{2} \sqrt{\rho} - 1$, and the estimator solves

$$H_n(\theta) = \int (\tilde{g} f_\theta)^{1/2} u_\theta dy = 0.$$

4. $B(\rho) = 1 - \rho^{1/p}$: then $d = 1 - \int g^{1/p} f_\theta^{1/q} dy$, with $1/p + 1/q = 1$. Estimator maximises $\int \tilde{g}^{1/p} f_\theta^{1/q} dy$.

Issues to pursue

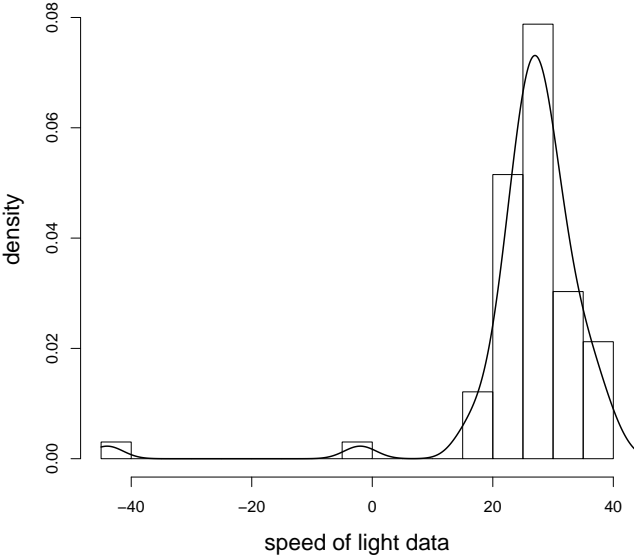
So, $\hat{\theta}$ is the minimiser of $d(\tilde{g}, f_{\theta}) = \int B(\tilde{g}/f_{\theta})f_{\theta} \, d\gamma$.

- ▶ 1 \tilde{g} nonparametric: properties $\hat{\theta}$?
- ▶ 2 And what is really required of \tilde{g} ? Fine-tuning?
- ▶ 3 $\hat{\theta}$ is often efficient at the model, with $\sqrt{n}(\hat{\theta} - \hat{\theta}_{\text{ML}}) \rightarrow_{\text{pr}} 0$; how different is it (then) from ML?
- ▶ 4 \tilde{g} bigger-parametric: properties $\hat{\theta}$?
- ▶ 5 Which $B(\rho)$ make the procedures robust?
- ▶ 6 Choice of $B(\rho)$ function (in a class of good ones)?
- ▶ 7 Generalising from i.i.d. to regression?
- ▶ 8 Goodness of fit, based on

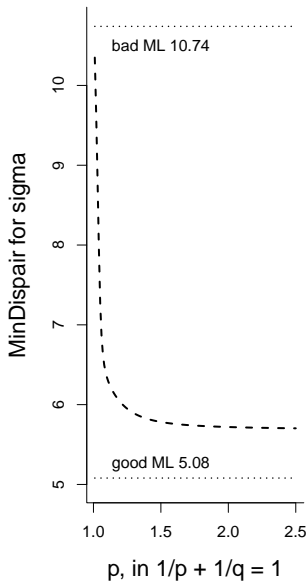
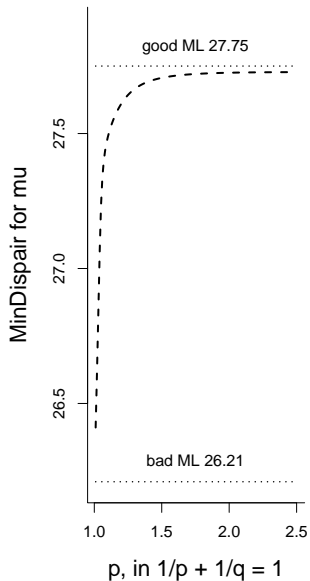
$$D_n = \min d(\tilde{g}, f_{\theta}) = d(\tilde{g}, f(\cdot, \hat{\theta}))?$$

- ▶ 9 Model selection using Min Dispair? Is there a MDIC, Minimum Divergence Information Criterion?

Example: Speed of light data, $n = 66$, two outliers to the left.



Minimising $1 - \int \tilde{g}^{1/p} f_{\theta}^{1/q} dy$, for different p , where $1/p + 1/q = 1$:



Minimising d (nonparametric, parametric)

Starting with a $B(\rho)$, use $A(\rho) = \rho B'(\rho) - B(\rho)$, to define the least false θ_0 as solution to

$$\int A\left(\frac{g}{f_\theta}\right) f_\theta u_\theta \, dy = 0$$

and the MinDisp estimator $\hat{\theta}$ as solution to

$$U_n(\theta) = \int A\left(\frac{\tilde{g}}{f_\theta}\right) f_\theta u_\theta \, dy = 0.$$

Under some conditions,

$$\sqrt{n}U_n(\theta_0) = \sqrt{n} \int \left\{ A\left(\frac{\tilde{g}}{f_{\theta_0}}\right) - A\left(\frac{g}{f_{\theta_0}}\right) \right\} f_{\theta_0} u_{\theta_0} \, dy \rightarrow_d U \sim N_p(0, K).$$

Along with 2nd order derivative matrix to a J , and other technical details, we're led to

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d J^{-1}U \sim N_p(0, J^{-1}KJ^{-1}).$$

The conditions for convergence of

$$\sqrt{n}U_n(\theta_0) = \sqrt{n} \int \left\{ A\left(\frac{\tilde{g}}{f_{\theta_0}}\right) - A\left(\frac{g}{f_{\theta_0}}\right) \right\} f_{\theta_0} u_{\theta_0} dy$$

basically amount to

$$\sqrt{n} \int A'\left(\frac{g}{f_{\theta_0}}\right) (\tilde{g} - g) u_{\theta_0} dy \rightarrow_d U \quad [\text{think } \sqrt{nh^2} \rightarrow 0]$$

and

$$\sqrt{n} \int A''\left(\frac{g}{f_{\theta_0}}\right) \frac{(\tilde{g} - g)^2}{f_{\theta_0}} u_{\theta_0} dy \rightarrow_{pr} 0 \quad [\text{think } \sqrt{nh} \rightarrow \infty].$$

1. The literature is **far too dominated** by 'always' using the **kernel density estimator** (granted, even there it's 'technical').
2. **At the model**, with $g = f_{\theta_0}$, matters simply, both J and K are proportional to J_{fish} , and **sandwich becomes J_{fish}^{-1}** : the MinDisp is as efficient as ML.
3. **Sam-Erik** needs within-reach generalisations of these matters to **$O(1/\sqrt{n})$ neighbourhoods** around given models.

Minimum Dispair with Bigger-Parametric Start

Suppose y_1, \dots, y_n are i.i.d. from some g . We fit data to $g(y, \theta, \gamma)$, which contains $f(y, \theta)$ as a special case:

$$f(y, \theta) = g(y, \theta, \gamma_0) \quad \text{for a known } \gamma_0.$$

With $(\tilde{\theta}, \tilde{\gamma})$ the ML in bigger model, the minimum dispair estimator for θ is $\hat{\theta}$, the minimiser of

$$d(\tilde{g}, f_\theta) = \int B\left(\frac{g(y, \tilde{\theta}, \tilde{\gamma})}{f(y, \theta)}\right) f(y, \theta) dy.$$

We have $\hat{\theta} \xrightarrow{\text{pr}} \theta_0$, the least false value minimising

$$\int B\left(\frac{g(y, \theta_1, \gamma_1)}{f(y, \theta)}\right) f(y, \theta) dy,$$

where (θ_1, γ_1) are KL minimisers in the bigger model. **If small model is correct**, $\hat{\theta}$ is consistent for the right θ_0 .

Assume smaller model correct, $g(y, \theta_0, \gamma_0) = f(y, \theta_0)$.

Theorem A: MinDisp is efficient at the model:

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow_d N_p(0, J(\theta_0)^{-1}).$$

Theorem B: Can use

$$D_n = \min_{\theta} d(\tilde{g}, f(\cdot, \theta)) = \int B\left(\frac{g(y, \tilde{\theta}, \tilde{\gamma})}{f(y, \hat{\theta})}\right) f(y, \hat{\theta}) dy$$

as a goodness-of-fit test statistic:

$$nD_n \rightarrow_d \frac{1}{2} B''(1) \chi_q^2$$

under model, where $q = \dim(\gamma)$.

Theorem C: Can work out $\sqrt{n}(\hat{\theta} - \theta_1)$ outside model conditions.

Remarks & questions re subset of Sam-Erik's PhD

A. Why isn't everyone using **Minimum Dispair**? Why don't we teach **minimum disparity estimators** in our courses?

We can choose B and \tilde{g} , and compute

$$\hat{\theta} = \operatorname{argmin} \int B\left(\frac{\tilde{g}}{f_{\theta}}\right) f_{\theta} dy.$$

It is **efficient at the model** (close enough to the ML then) and **robust**.

- ▶ **slim following** (Basu et collegae; students of Lindsay; Hooker; a few other mild epicentres – though we've invented **MWL** and **WIC**).
- ▶ technicalities, fine-tuning choices, theorems are still lacking, and they're **nitty-gritty demanding**?
- ▶ statisticians are **happy with ML** and closer cousins?
- ▶ robustness **isn't judged an important issue** (despite propaganda in the 70ies and 80ies)?
- ▶ doesn't easily generalise to regression?

B. We need clearer conditions on \tilde{g} (and I'm bored by the default kernel method).

C. We should work out a suitable MDIC, a Minimum Disparity Information Criterion (non-trivial task).

The distance from truth to fitted model is $\int B(g/f_{\hat{\theta}})f_{\hat{\theta}} dy$. What we directly observe is $D_n = \int B(\tilde{g}/f_{\hat{\theta}})f_{\hat{\theta}} dy$. Need analysis of

$$\begin{aligned}\Delta_n &= \int B\left(\frac{\tilde{g}}{f_{\hat{\theta}}}\right)f_{\hat{\theta}} dy - \int B\left(\frac{g}{f_{\hat{\theta}}}\right)f_{\hat{\theta}} dy \\ &\doteq \int B'\left(\frac{g}{f_{\hat{\theta}}}\right)(\tilde{g} - g) dy + \frac{1}{2} \int B''\left(\frac{g}{f_{\hat{\theta}}}\right)\frac{(\tilde{g} - g)^2}{f_{\hat{\theta}}} dy.\end{aligned}$$

For some \tilde{g} , may show

$$(nh)^{1/2}\Delta_n \rightarrow_d \int B'\left(\frac{g}{f_{\theta_0}}\right)Z(y) dy,$$

which leads to a MDIC = $D_n + (nh)^{-1/2}\hat{q}$ (it's messy, though).

D. Some of the schemes working well in $\dim = 1$ have trouble in $\dim \geq 2$ (other biases kick in). Can these matters be ameliorated?

E. Can MinDisp estimators take the place of ML estimators in FIC (the Gerda-Nils Focused Information Criterion)? I believe yes – and Sam-Erik is working on this. The root of the matter for ML is this: If data stem from $f_n(y) = f(y, \theta_0)\{1 + r(y)/\sqrt{n}\}$, then

$$\sqrt{n}(\hat{\theta}_{\text{ML}} - \theta_0) \rightarrow_d N_p(J^{-1}b, J^{-1}), \quad b = \int f(y, \theta_0)r(y)u(y, \theta_0) dy.$$

Task: Demonstrate that this holds (along with a certain list of other technical things) for classes of MinDisp.

F. Regression setting, model $f(y_i | x_i, \theta)$: choose $\hat{\theta}$ to minimise

$$n^{-1} \sum_{i=1}^n \int B\left(\frac{\tilde{g}(y | x_i)}{f(y | x_i, \theta)}\right) f(y | x_i, \theta) dy.$$

G. MinDisp needs to be shown at work in clear well-motivated application stories.