

Notes on Robust Student t Regression

Max Leung

January 24, 2024

1 Arbitrary Elliptically Symmetric Family of Densities

Suppress notation i for observation i . $\dim(\mathbf{y}) = k$

1.1 Probability Density Function

$$p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Omega}, \nu) = |\boldsymbol{\Omega}|^{-1/2} g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)$$

1.1.1 Special Case: k-variate generalized t distribution

$$g(s, \nu) = \frac{\Gamma((\nu + k)/2)}{\Gamma(1/2)^k \Gamma(\nu/2) \nu^{k/2}} \left(1 + \frac{s}{\nu}\right)^{-(\nu+k)/2}$$
$$\mathbf{y} \sim t_k(\boldsymbol{\mu}(\boldsymbol{\theta}), \boldsymbol{\Omega}, \nu)$$

Property: $\mathbb{E}(\mathbf{y}) = \boldsymbol{\mu}(\boldsymbol{\theta})$

1.1.2 Special Case: k-variate generalized power-exponential family

$$g(s, \nu) = c(\nu) e^{-s^\nu/2}$$

1.2 Gradient Vector

Assume $\boldsymbol{\mu}$ is a function of $\boldsymbol{\theta}$ i.e., $\boldsymbol{\mu}(\boldsymbol{\theta})$

$$\begin{aligned} \frac{\partial \ln[p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Omega}, \nu)]}{\partial \boldsymbol{\theta}} &= \frac{\partial \ln[|\boldsymbol{\Omega}|^{-1/2} g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)]}{\partial \boldsymbol{\theta}} \\ &= \frac{\partial (\ln[|\boldsymbol{\Omega}|^{-1/2}] + \ln[g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)])}{\partial \boldsymbol{\theta}} \\ &= \frac{\partial \ln[g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)]}{\partial \boldsymbol{\theta}} \\ &= \frac{\partial \ln[g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)]}{\partial g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)} \frac{\partial g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)}{\partial \boldsymbol{\theta}} \\ &= \frac{1}{g} \left(\frac{\partial g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)}{\partial (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu})} \frac{\partial (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\partial \boldsymbol{\theta}} + \frac{\partial g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)}{\partial \nu} \frac{\partial \nu}{\partial \boldsymbol{\theta}} \right) \\ &= \frac{1}{g} (g_1 \cdot \frac{\partial (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\partial \boldsymbol{\theta}} + \frac{\partial g((\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}), \nu)}{\partial \nu} \cdot 0) \\ &= \frac{g_1}{g} \frac{\partial (\mathbf{y} - \boldsymbol{\mu})'}{\partial \boldsymbol{\theta}} (\boldsymbol{\Omega}^{-1} + \boldsymbol{\Omega}^{-1'}) (\mathbf{y} - \boldsymbol{\mu}) \\ &= \frac{g_1}{g} \left(-\frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})'}{\partial \boldsymbol{\theta}} \right) (2 \cdot \boldsymbol{\Omega}^{-1}) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\theta})) \\ &= -2 \frac{g_1}{g} \frac{\partial \boldsymbol{\mu}(\boldsymbol{\theta})'}{\partial \boldsymbol{\theta}} \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\theta})) \end{aligned}$$

1.2.1 Special Case: k-variate generalized t distribution

$$\begin{aligned}
g_1 &= \frac{d \frac{\Gamma((\nu+k)/2)}{\Gamma(1/2)^k \Gamma(\nu/2) \nu^{k/2}} (1 + \frac{s}{\nu})^{-(\nu+k)/2}}{ds} \\
&= -\frac{\Gamma((\nu+k)/2)}{\Gamma(1/2)^k \Gamma(\nu/2) \nu^{k/2}} ((\nu+k)/2) (1 + \frac{s}{\nu})^{-(\nu+k)/2-1} \frac{1}{\nu} \\
\frac{g_1}{g} &= \frac{-\frac{\Gamma((\nu+k)/2)}{\Gamma(1/2)^k \Gamma(\nu/2) \nu^{k/2}} ((\nu+k)/2) (1 + \frac{s}{\nu})^{-(\nu+k)/2-1} \frac{1}{\nu}}{\frac{\Gamma((\nu+k)/2)}{\Gamma(1/2)^k \Gamma(\nu/2) \nu^{k/2}} (1 + \frac{s}{\nu})^{-(\nu+k)/2}} \\
&= -((\nu+k)/2) (1 + \frac{s}{\nu})^{-1} \frac{1}{\nu} \\
&= -\frac{1}{2} (\nu+k) \frac{\nu}{\nu+s} \frac{1}{\nu} \\
&= -\frac{1}{2} \frac{\nu+k}{\nu+s}
\end{aligned}$$

where $s = (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{y} - \boldsymbol{\mu})$

1.3 Maximum Likelihood Estimation

Assume independent \mathbf{y}_i with common density p

$$\begin{aligned}
\frac{\partial \ln L}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} &= \frac{\partial \ln \prod_{i=1}^N p(\mathbf{y}_i | \boldsymbol{\mu}_i, \boldsymbol{\Omega}_i, \nu)}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} = \mathbf{0} \\
&\sum_{i=1}^N \frac{\partial \ln [p(\mathbf{y}_i | \boldsymbol{\mu}_i, \boldsymbol{\Omega}_i, \nu)]}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} = \mathbf{0} \\
&\sum_{i=1}^N \frac{g_{1,i}}{g_i} \frac{\partial \boldsymbol{\mu}_i(\boldsymbol{\theta})'}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} \boldsymbol{\Omega}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}}_{mle})) = \mathbf{0}
\end{aligned}$$

1.3.1 Special Case: k-variate generalized t distribution

$$s_i = (\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}}_{mle}))' \boldsymbol{\Omega}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}}_{mle}))$$

$$\begin{aligned}
&\sum_{i=1}^N \left(-\frac{1}{2} \frac{\nu + k_i}{\nu + s_i} \right) \frac{\partial \boldsymbol{\mu}_i(\boldsymbol{\theta})'}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} \boldsymbol{\Omega}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}}_{mle})) = \mathbf{0} \\
&\sum_{i=1}^N \frac{\nu + k_i}{\nu + s_i} \frac{\partial \boldsymbol{\mu}_i(\boldsymbol{\theta})'}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} \boldsymbol{\Omega}_i^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_i(\hat{\boldsymbol{\theta}}_{mle})) = \mathbf{0}
\end{aligned}$$

The resulting ML estimator is robust in the sense that outlier observation with large s_i has smaller weight in the log likelihood function.

If $k_i = 1$ for $\forall i$, $y_i \sim t_1(\mu(\boldsymbol{\theta}, \mathbf{x}_i), \sigma_i^2, \nu)$ where $\mu(\boldsymbol{\theta}, \mathbf{x}_i)$ can be linear $\mathbf{x}_i' \boldsymbol{\theta}$

$$\sum_{i=1}^N \frac{\nu + 1}{\nu + s_i} \frac{\partial \mu(\boldsymbol{\theta}, \mathbf{x}_i)'}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}_{mle}} \sigma_i^{-2} (y_i - \mu(\hat{\boldsymbol{\theta}}_{mle}, \mathbf{x}_i)) = 0$$

where $s_i = (y_i - \mu(\hat{\boldsymbol{\theta}}_{mle}, \mathbf{x}_i))^2 / \sigma_i^2$

1.4 Student-t is a mix density

The density of k-variate generalized t random variable can also be written as a mixture of multivariate normal and scaled inverse chi-square densities.

$$\begin{aligned}
t(\mathbf{y} | \boldsymbol{\mu}, \boldsymbol{\Omega}, \nu) &= \int_0^\infty N(\mathbf{y} | \boldsymbol{\mu}, z \cdot \boldsymbol{\Omega}) \chi^{-2}(z | \nu, 1) dz \\
&= \int_0^\infty p(\mathbf{y}, z | \boldsymbol{\mu}, \boldsymbol{\Omega}, \nu) dz
\end{aligned}$$

Note that $IG(z|\frac{\nu}{2}, \frac{\nu}{2}) = \chi^{-2}(z|\nu, 1)$ where IG = Inverse Gamma. $IG(z|a, b) = \frac{b^a}{\Gamma(a)} z^{-(a+1)} \exp(-b/z)$

If $k = 1$,

$$\begin{aligned} t(y|\mu, \sigma^2, \nu) &= \int_0^\infty N(y|\mu, z \cdot \sigma^2) \chi^{-2}(z|\nu, 1) dz \\ &= \int_0^\infty p(y, z|\mu, \sigma^2, \nu) dz \end{aligned}$$

2 EM estimation of robust regression with univariate generalized t distribution

Regarding z_i as missing data, assume ν is known (or specify a grid of possible ν values).

2.1 E step

$$\begin{aligned} \ln[p(y_i, z_i | \overbrace{\mu(\boldsymbol{\theta}, \mathbf{x}_i)}^{\boldsymbol{\theta}, \mathbf{x}_i}, \sigma_i^2, \nu)] &= \ln[N(y_i | \boldsymbol{\theta}, \mathbf{x}_i, z_i \cdot \sigma_i^2) \chi^{-2}(z_i | \nu, 1)] \\ &= \ln[N(y_i | \boldsymbol{\theta}, \mathbf{x}_i, z_i \cdot \sigma_i^2)] + \ln[\chi^{-2}(z_i | \nu, 1)] \\ &= \ln\left[\frac{1}{\sqrt{z_i \cdot \sigma_i^2 \cdot 2\pi}} \exp\left(-\frac{1}{2} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2}\right)\right] + \ln\left[\frac{1}{\Gamma(\nu/2)} \left(\frac{\nu}{2}\right)^{\frac{\nu}{2}} z_i^{-\nu/2-1} \exp\left(-\frac{\nu}{2z_i}\right)\right] \\ &= -\frac{1}{2} \ln(z_i \cdot \sigma_i^2 \cdot 2\pi) - \frac{1}{2} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2} + \left(-\frac{\nu}{2} - 1\right) \ln[z_i]^{-\nu/2-2} - \frac{\nu}{2z_i} + \text{constant} \end{aligned}$$

Thus, given independence

$$\begin{aligned} \mathbb{E}[\ln[p(\mathbf{y}, \mathbf{z}_i | \boldsymbol{\theta}, \mathbf{x}_i, \sigma_i^2, \nu)] | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu] &= \mathbb{E}[\ln\left[\prod_{i=1}^N p(y_i, z_i | \boldsymbol{\theta}, \mathbf{x}_i, \sigma_i^2, \nu)\right] | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu] \\ &= \mathbb{E}\left[\sum_{i=1}^N \ln[p(y_i, z_i | \boldsymbol{\theta}, \mathbf{x}_i, \sigma_i^2, \nu)] | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu\right] \\ &= \mathbb{E}\left[\sum_{i=1}^N \left\{-\frac{1}{2} \ln(z_i \cdot \sigma_i^2 \cdot 2\pi) - \frac{1}{2} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2} + \left(-\frac{\nu}{2} - 1\right) \ln[z_i]^{-\nu/2-2} - \frac{\nu}{2z_i} + \text{constant}\right\} | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu\right] \end{aligned}$$

Omitting terms without $\boldsymbol{\theta}$, as those will be differentiated away in M step

$$\begin{aligned} &\approx \mathbb{E}\left[-\frac{1}{2} \frac{\sum_{i=1}^N (y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2} | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu\right] \\ &= -\sum_{i=1}^N \underbrace{\mathbb{E}\left[\frac{1}{z_i} | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu\right]}_{w_i} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{2\sigma_i^2} \end{aligned}$$

In order to evaluate w_i , we have to find out the posterior distribution $z_i | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu$

$$\begin{aligned} p(z_i | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu) &\propto p(y_i | z_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2) p(z_i | \nu) \\ &= N(y_i | \boldsymbol{\theta}, \mathbf{x}_i, z_i \cdot \sigma_i^2) \chi^{-2}(z_i | \nu, 1) \\ &= \frac{1}{\sqrt{z_i \cdot \sigma_i^2 \cdot 2\pi}} \exp\left(-\frac{1}{2} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2}\right) \cdot \frac{1}{\Gamma(\nu/2)} \left(\frac{\nu}{2}\right)^{\frac{\nu}{2}} z_i^{-\nu/2-1} \exp\left(-\frac{\nu}{2z_i}\right) \\ &= z_i^{-1/2} z_i^{-\nu/2-1} \exp\left(-\frac{1}{2} \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{z_i \sigma_i^2}\right) \exp\left(-\frac{\nu}{2z_i}\right) \cdot C \\ &= z_i^{-(v+1)/2-1} \exp\left(-\frac{1}{2z_i} \left(\frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{\sigma_i^2} + \nu\right)\right) \cdot C \\ &= z_i^{-[(v+1)/2+1]} \exp\left(-\left[\frac{1}{2} \left(\frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{\sigma_i^2} + \nu\right)\right] / z_i\right) \cdot C \end{aligned}$$

Given the fact that inverse gamma is the conjugate prior of normal-inverse-gamma mix, $z_i|y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu \sim IG(\frac{\nu+1}{2}, \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2 / \sigma_i^2 + \nu}{2})$. Therefore, $\frac{1}{z_i} \sim G(\frac{\nu+1}{2}, \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2 / \sigma_i^2 + \nu}{2})$. By the property of Gamma random variable, $w_i = \mathbb{E}[\frac{1}{z_i} | y_i, \mathbf{x}_i, \boldsymbol{\theta}, \sigma_i^2, \nu] = (\frac{\nu+1}{2}) / (\frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2 / \sigma_i^2 + \nu}{2}) = \frac{\nu+1}{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2 / \sigma_i^2 + \nu}$

2.2 M Step

$$\begin{aligned}\boldsymbol{\theta}^{t+1} &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}[\ln[p(\mathbf{y}, z_i | \boldsymbol{\theta}, \mathbf{x}_i, \sigma_i^2, \nu)] | y_i, \mathbf{x}_i, \boldsymbol{\theta}^t, \sigma_i^{2,t}, \nu] \\ &= \arg \max_{\boldsymbol{\theta}} \left\{ - \sum_{i=1}^N w_i^t \frac{(y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{2\sigma_i^{2,t}} \right\} \\ &= \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{i=1}^N \frac{w_i^t}{\sigma_i^{2,t}} (y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2 \right\}\end{aligned}$$

where $w_i^t = \frac{\nu+1}{(y_i - \mu(\boldsymbol{\theta}^t, \mathbf{x}_i))^2 / \sigma_i^{2,t} + \nu}$ and $\sigma_i^{2,t+1} = \sum_{i=1}^N \frac{w_i^t}{\sigma_i^{2,t}} (y_i - \mu(\boldsymbol{\theta}^{t+1}, \mathbf{x}_i))^2 / N$

It is a weighted non-linear least square (NLS) problem with weight $\frac{w_i}{\sigma_i^2}$. If $\mu(\boldsymbol{\theta}, \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\theta}$, EM algorithm is non-iterative. If $\mu(\boldsymbol{\theta}, \mathbf{x}_i)$ is non-linear, EM algorithm is iterative. Thus, some says that EM algorithm is Iteratively Re-weighted NLS. Its First Order Condition (FOC) is the same as MLE's FOC.

$$\begin{aligned}\frac{\partial \sum_{i=1}^N \frac{w_i^t}{\sigma_i^{2,t}} (y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}^{t+1}} &= \mathbf{0} \\ \sum_{i=1}^N \frac{w_i^t}{\sigma_i^{2,t}} \frac{\partial (y_i - \mu(\boldsymbol{\theta}, \mathbf{x}_i))^2}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}^{t+1}} &= \mathbf{0} \\ \sum_{i=1}^N \frac{w_i^t}{\sigma_i^{2,t}} \cdot (-2) \cdot (y_i - \mu(\boldsymbol{\theta}^{t+1}, \mathbf{x}_i)) \frac{\partial \mu(\boldsymbol{\theta}, \mathbf{x}_i)'}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}^{t+1}} &= \mathbf{0} \\ \sum_{i=1}^N w_i^t \frac{\partial \mu(\boldsymbol{\theta}, \mathbf{x}_i)'}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}^{t+1}} \sigma_i^{-2,t} (y_i - \mu(\boldsymbol{\theta}^{t+1}, \mathbf{x}_i)) &= \mathbf{0}\end{aligned}$$

3 Reference

Lange, K. L., Roderick J. A. Little, & Jeremy M. G. Taylor. (1989). Robust Statistical Modeling Using the t Distribution. Journal of the American Statistical Association, 84(408), 881–896. <https://doi.org/10.2307/2290063>

Kevin P. Murphy. (2023). Probabilistic Machine Learning: Advanced Topics, MIT Press