

Tracking AI in climate innovation: Supplementary information

December 8, 2021

Supplementary Information

In this part we provide the supplementary results referred to from the main text and Methods.

Contents

1	Modeling details	2
2	Model diagnostics	3
3	Exploring the data and the covariates	32

1 Modeling details

- Time period for statistical analysis: 2010-2019
- Target variable: 3-year forward citation in CPC subclasses with the following covariates. This results in analysing forward citations for patents from 2010-2017.
- Pre-processing:
 - For regression modeling, filter out top 1% of most highly cited patents (outliers skewed the model fit quite a bit, removing them clearly improved the fit in diagnostics plots, as seen below).
 - Remove subgroups that lack patents in each of the year 2010-2017 (around 1% subgroups with little or no patenting activity).
- List of covariates/controls:
 - ai (factor, binary: AI classification according to the WIPO)
 - green_ycodes (factor, binary: if Y02 label on patent)
 - grantyear (factor, 8 levels: [2010,...,2017])
 - cpc_subarea (factor: for CPC subclass, groups; for groups, subgroups)
 - tct_type (factor, 4 levels: [0-5,5-10,10-15,15-]) (based on the data exploration below)
 - n_claims (log plus one transformed): number of claims in patent
 - n_inventors_individual (log plus one transformed): number of inventors
 - n_backward_patent_citations (log plus one transformed): number of patent citations
 - n_backward_research_citations (log plus one transformed): number of research citations
 - n_backward_other_citations (log plus one transformed): number of other citations
 - organizational (factor, binary: if firm on patent)
- GAMLSS method: General Additive Model for Location, Shape, and Scaleⁱ. This approach was chosen because the data turns out to be both zero-inflated and heavily skewed – variance that differs across subclasses and groups. Including model parameters for subclass- and group-level skewness led to models with good fit, compared to Poisson and Negative Binomial models.
- The specific model for the citation count data is: Generalized inverse Gaussian distribution (Sichel distribution), as this can model highly skewed and zero-inflated count data.^{ii,iii}
- Model specification on the level of technologies (subclasses/groups): We use repeated modeling for the different technology areas to get estimates of how AI predicts a difference in citation between the units in the data. Below, we also present residuals for the subclass- and group-level regressions.
- Model diagnostics: for the models, the normalized randomized quantile residuals^{iv} (for details, see section 6.2 in reference given by footnote (iii)). Model diagnostics for these can be done on a Q-Q plot for approximate normality.

ⁱRigby, R.A. and Stasinopoulos, D.M. (2005), Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 54: 507-554.

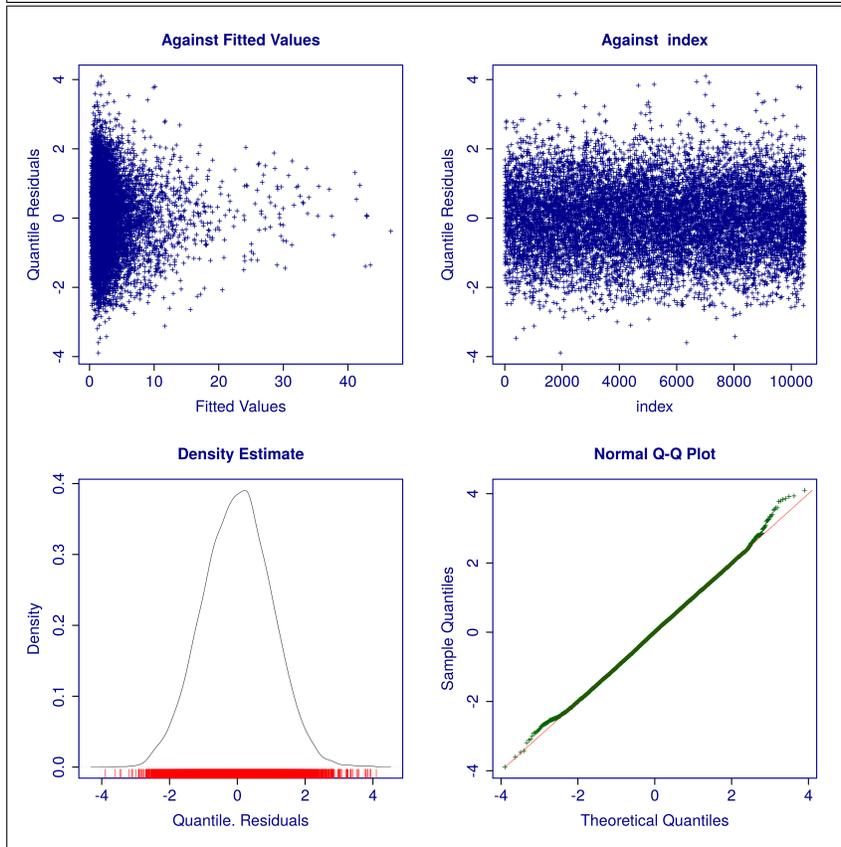
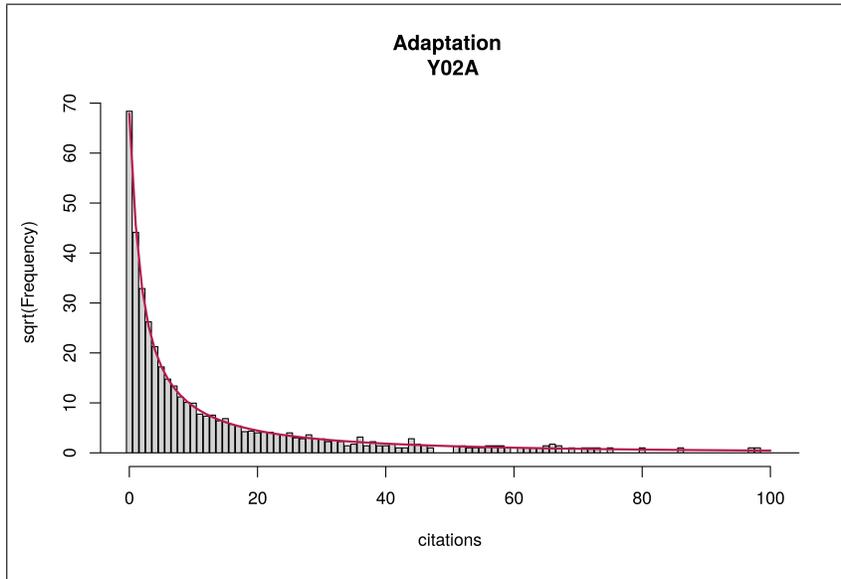
ⁱⁱGillian Z. Stein, Walter Zucchini and June M. Juritz. Parameter Estimation for the Sichel Distribution and its Multivariate Extension. (1987) *Journal of the American Statistical Association*. Vol. 82, No. 399, pp. 938-944

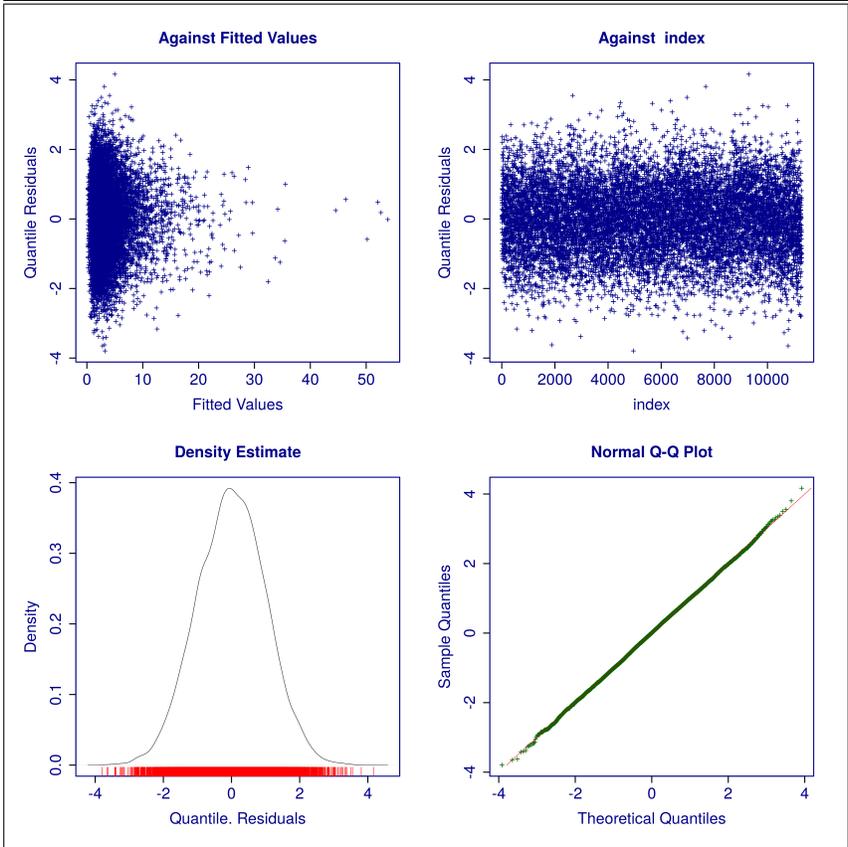
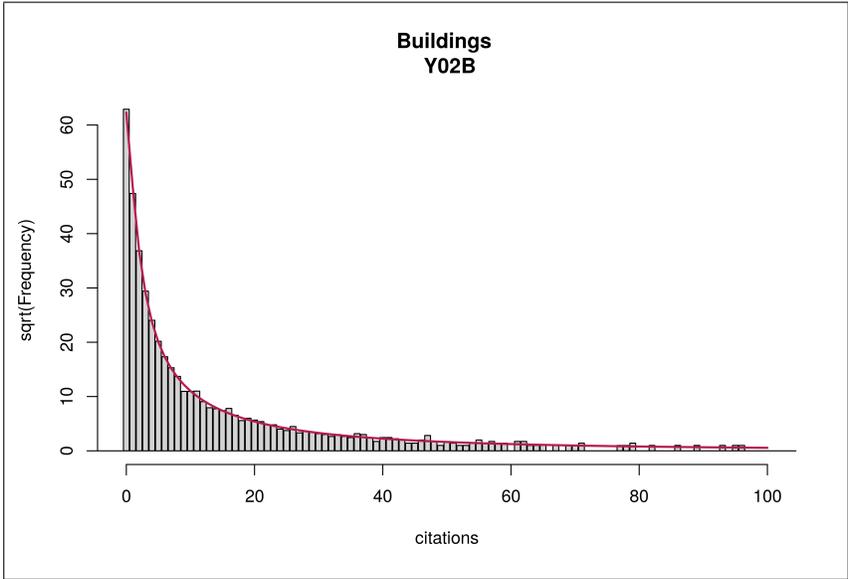
ⁱⁱⁱR. A. Rigby, D. M. Stasinopoulos, and C. Akantziliotou. (2008) A framework for modelling overdispersed count data, including the Poisson-shifted generalized inverse Gaussian distribution. *Computational Statistics and Data Analysis*, 53(2):381–393.

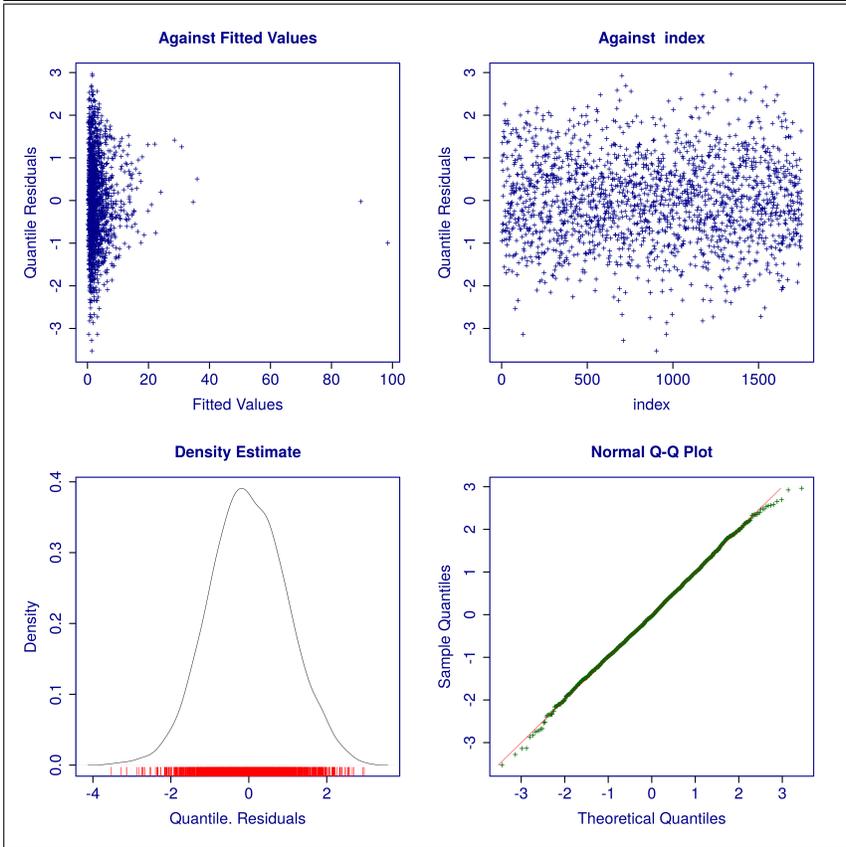
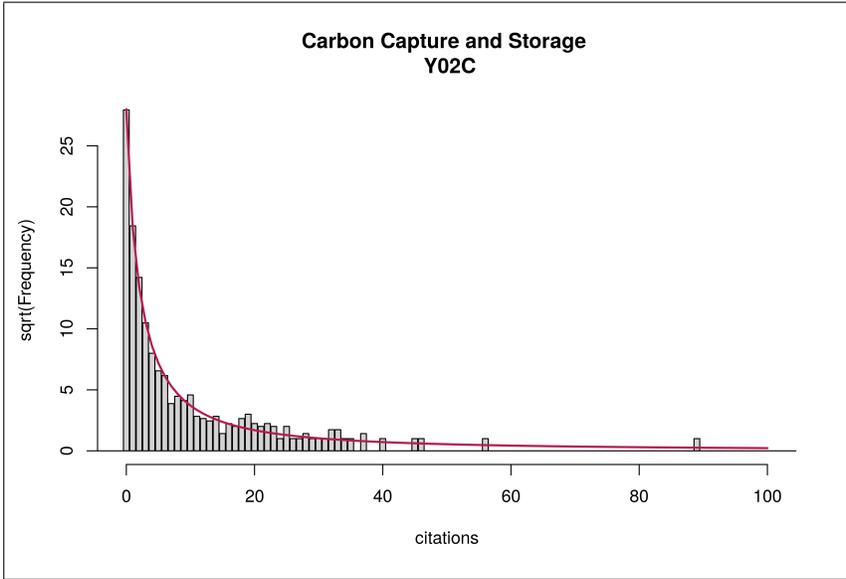
^{iv}Randomized Quantile Residuals. Peter K. Dunn & Gordon K. Smyth. (1996) *Journal of Computational and Graphical Statistics* Vol. 5, No. 3. Pages 236-244.

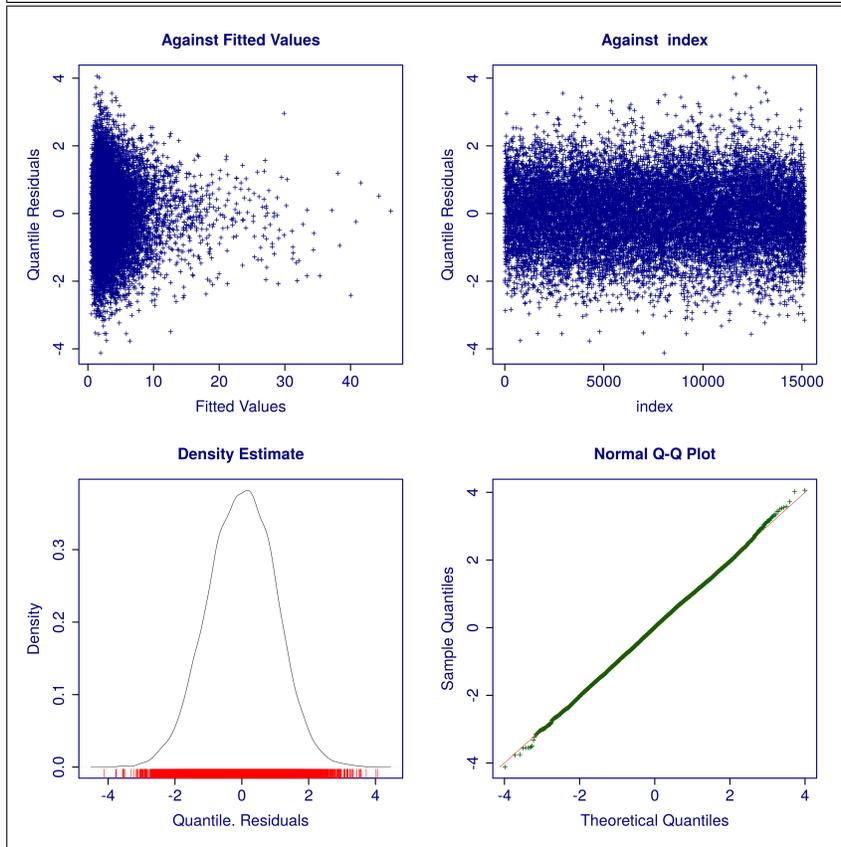
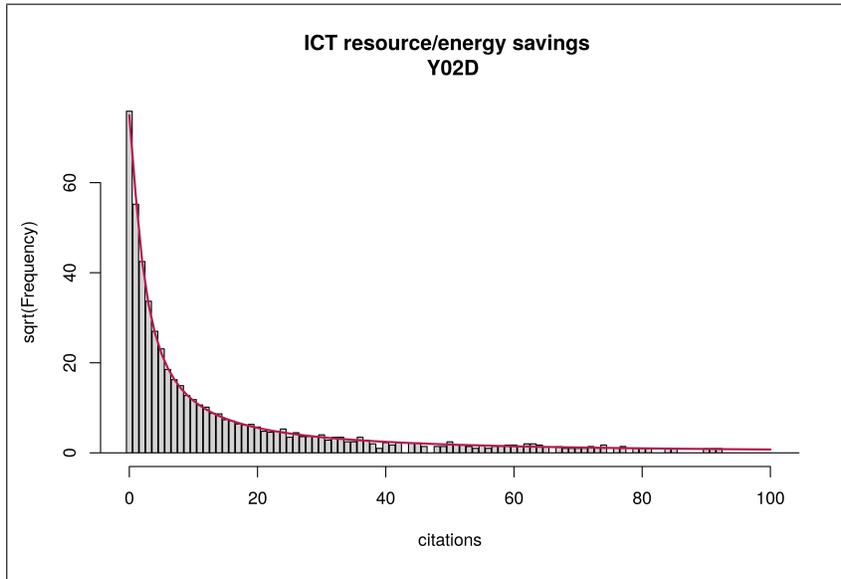
2 Model diagnostics

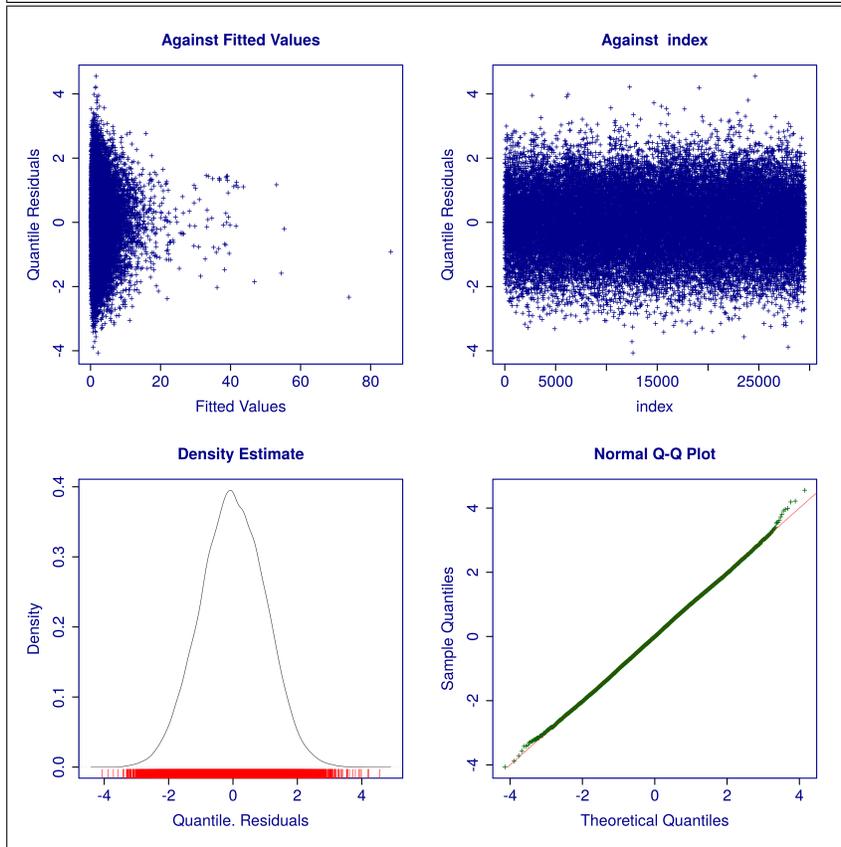
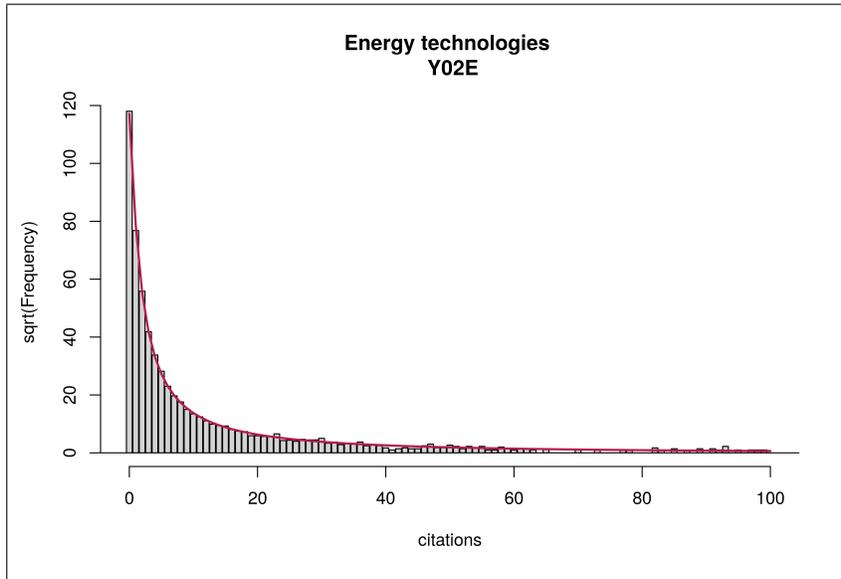
This section displays model diagnostics, both with rootograms – indicating whether or not the model (red curves) produces similar results as the data (the grey bins) – as well as residual plots. Quantile residuals should be approximately aligned with the diagonal of the Q-Q plots, which we find is the case for the final model.

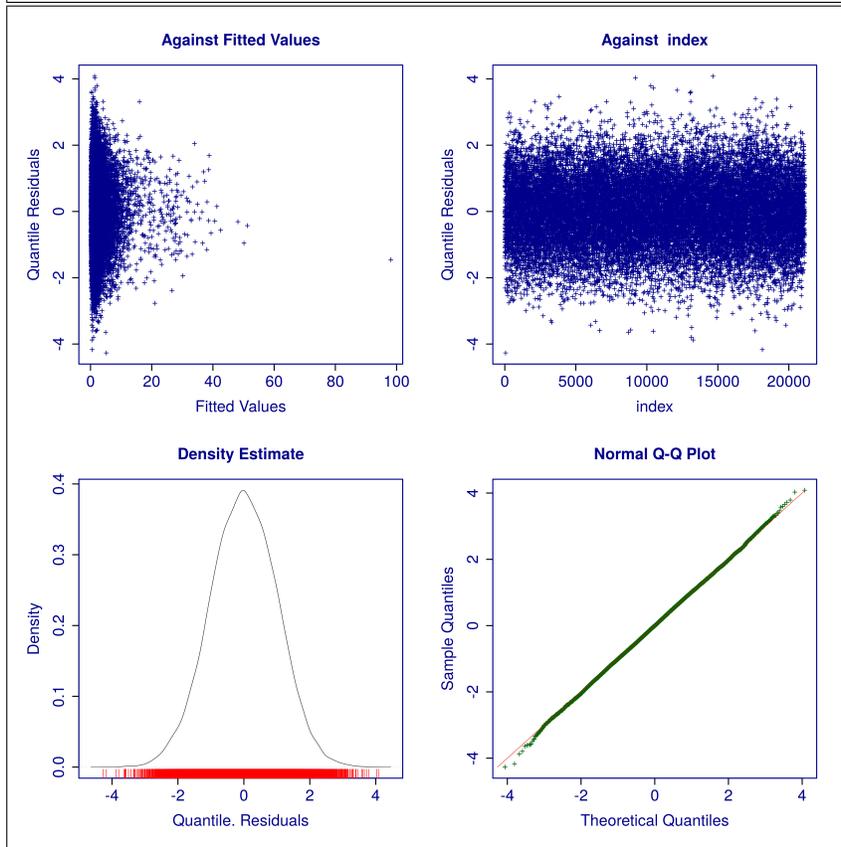
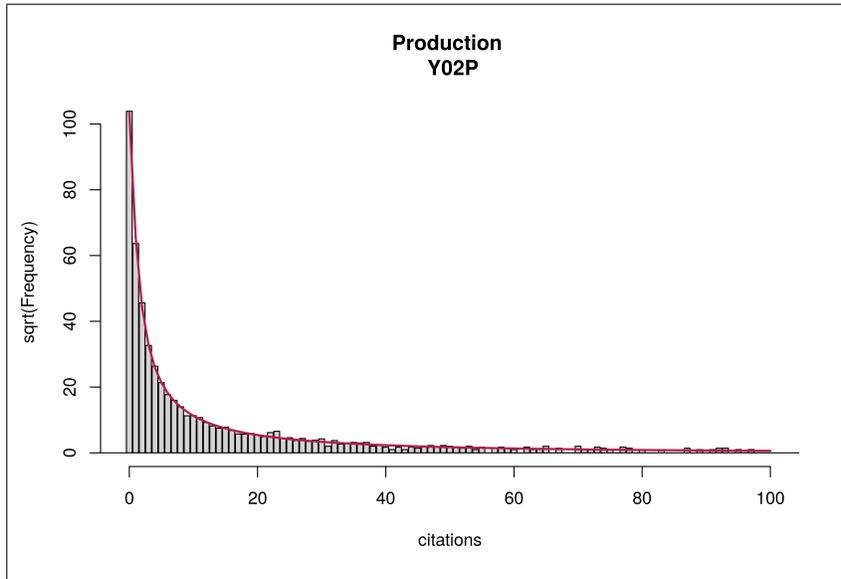


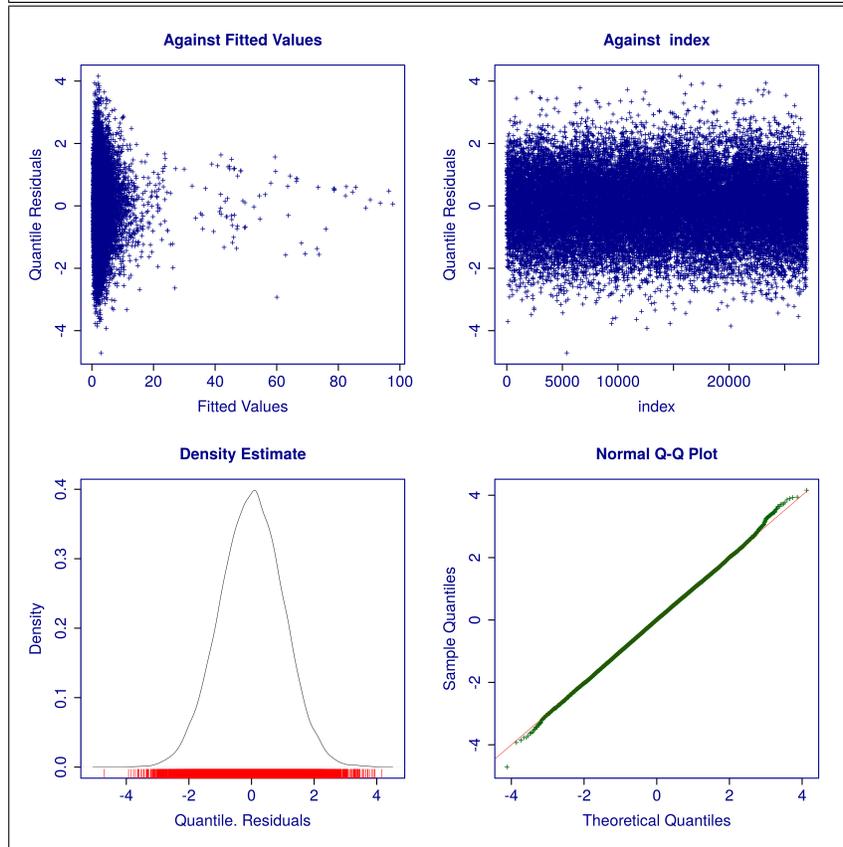
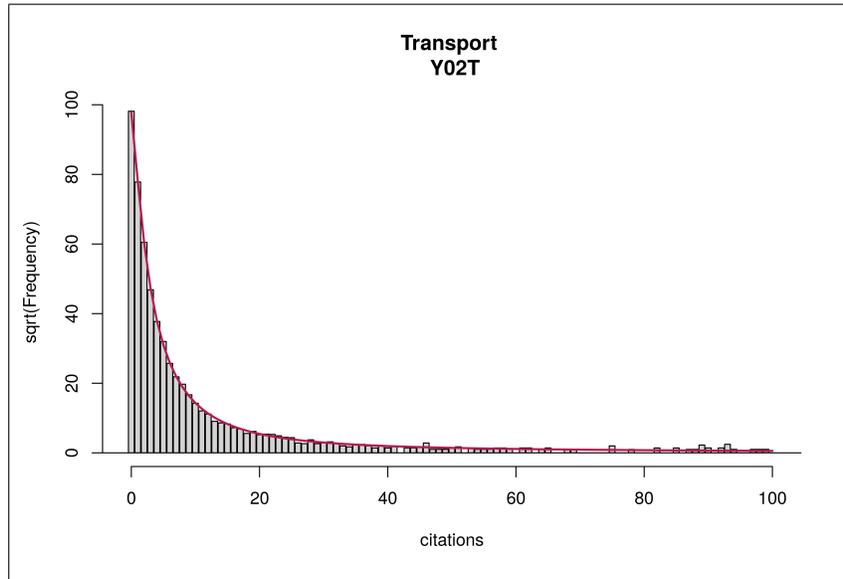


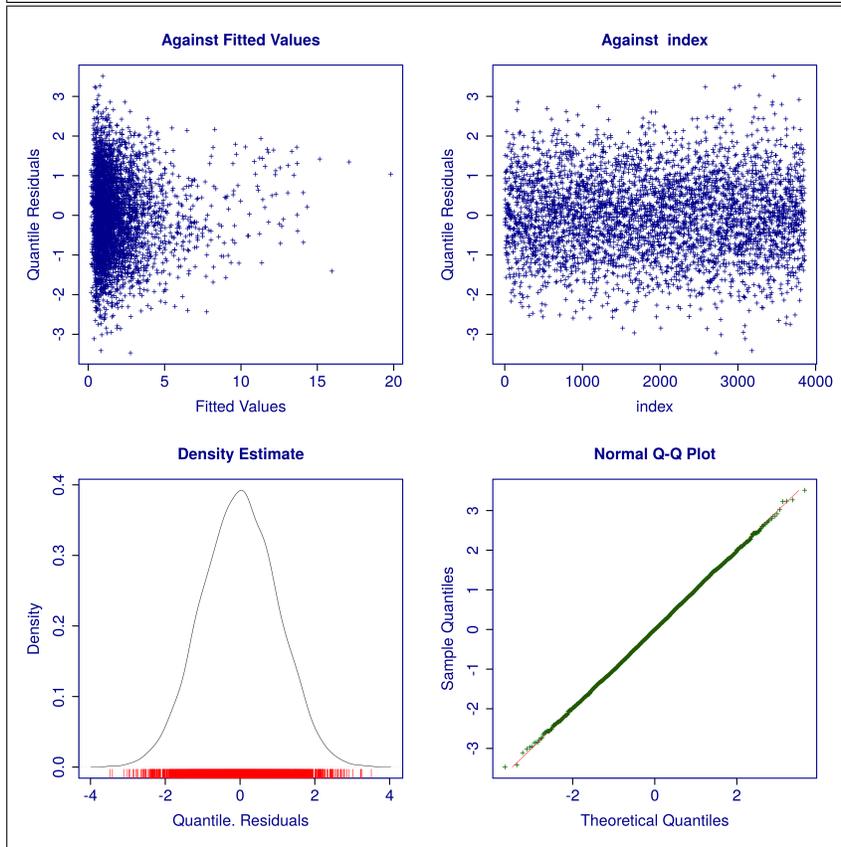
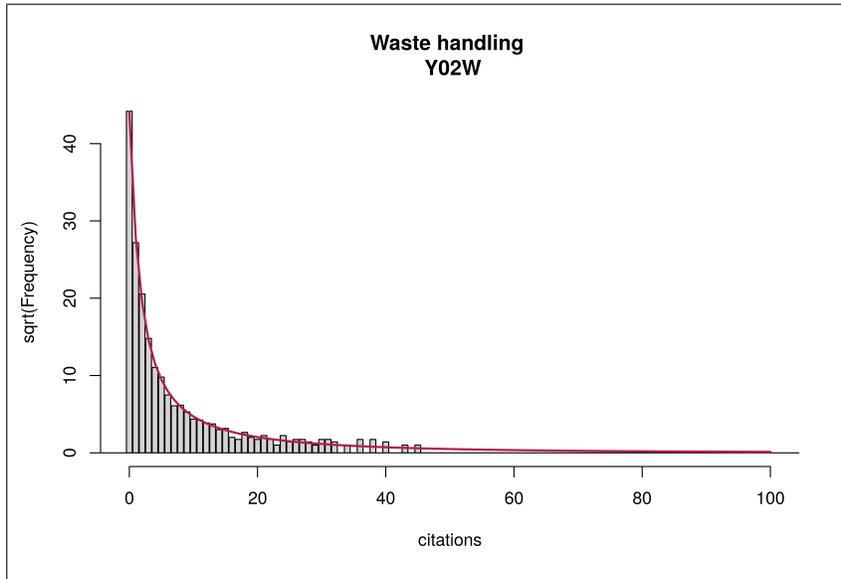


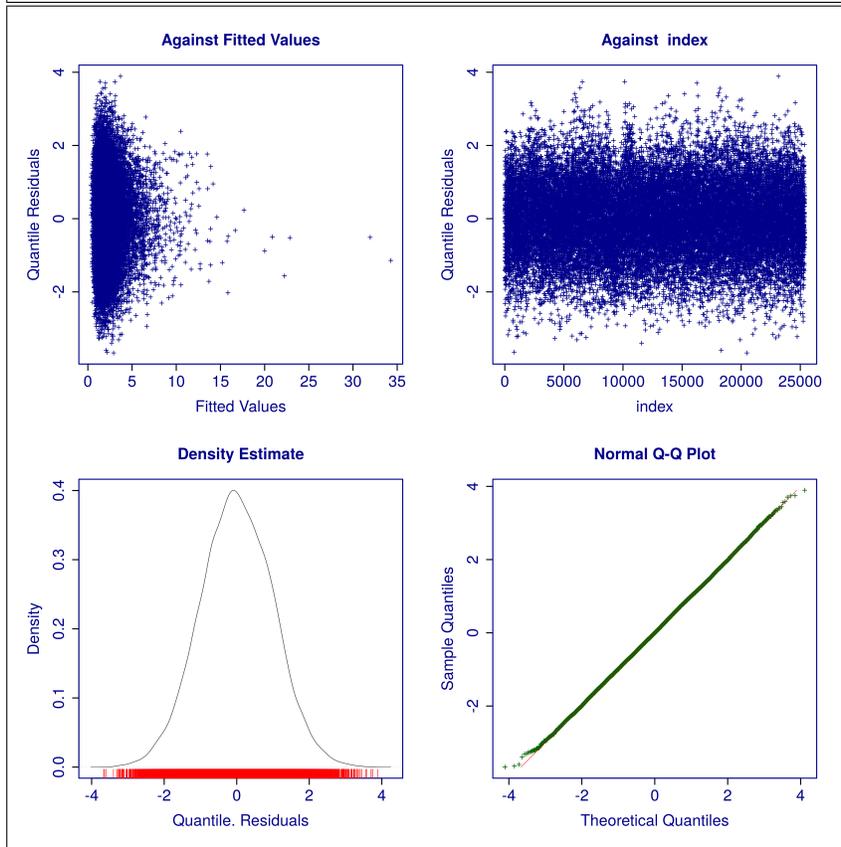
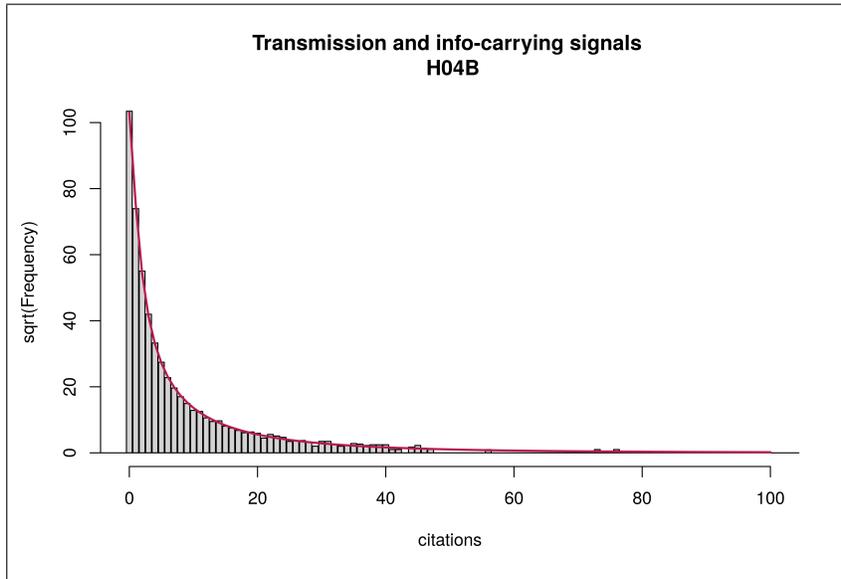


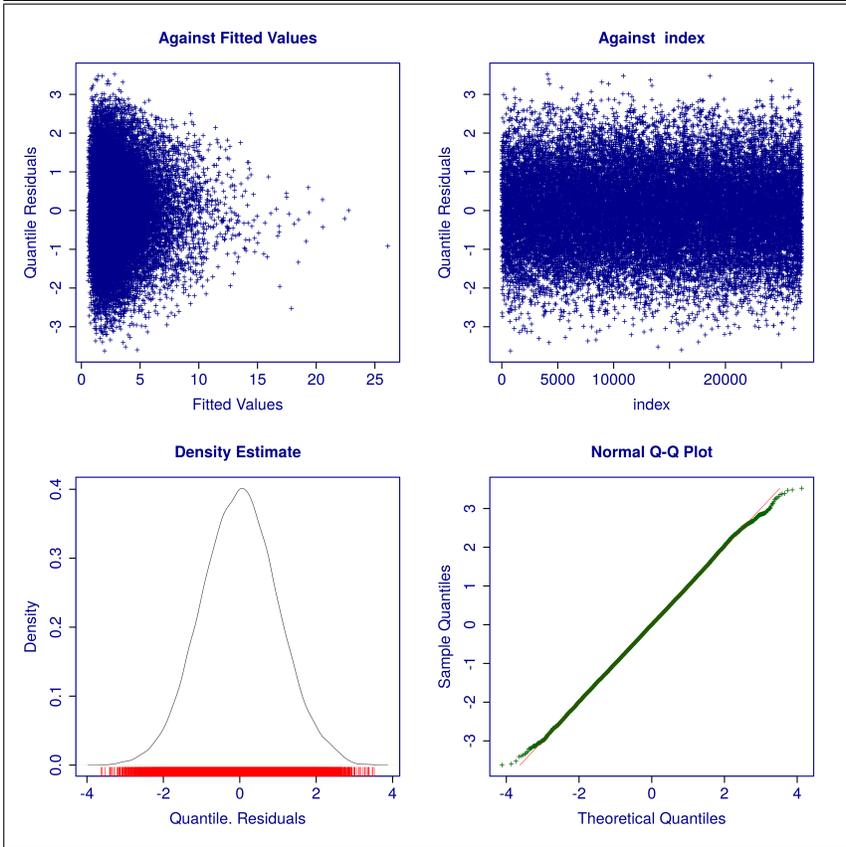
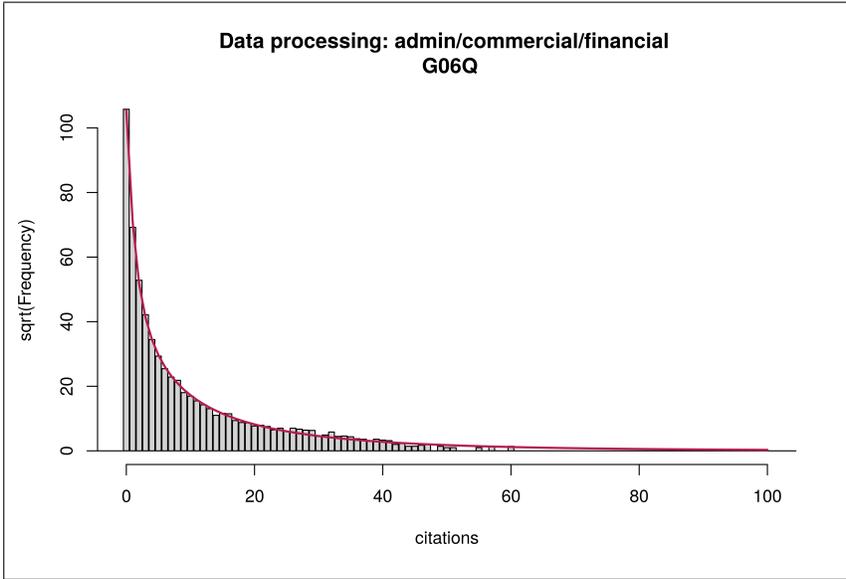


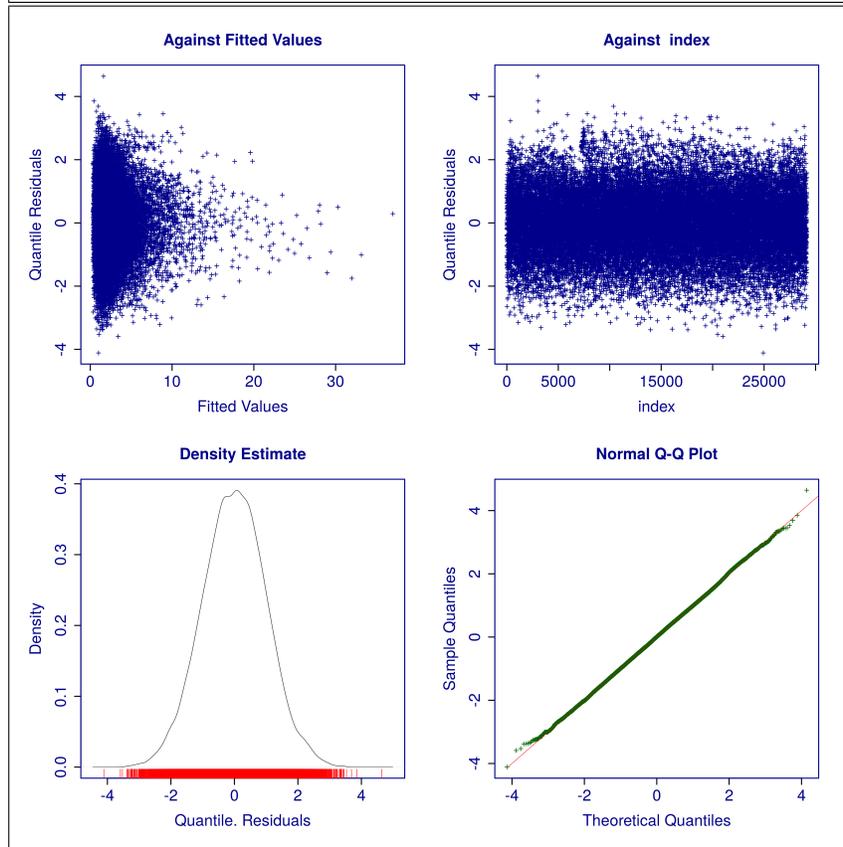
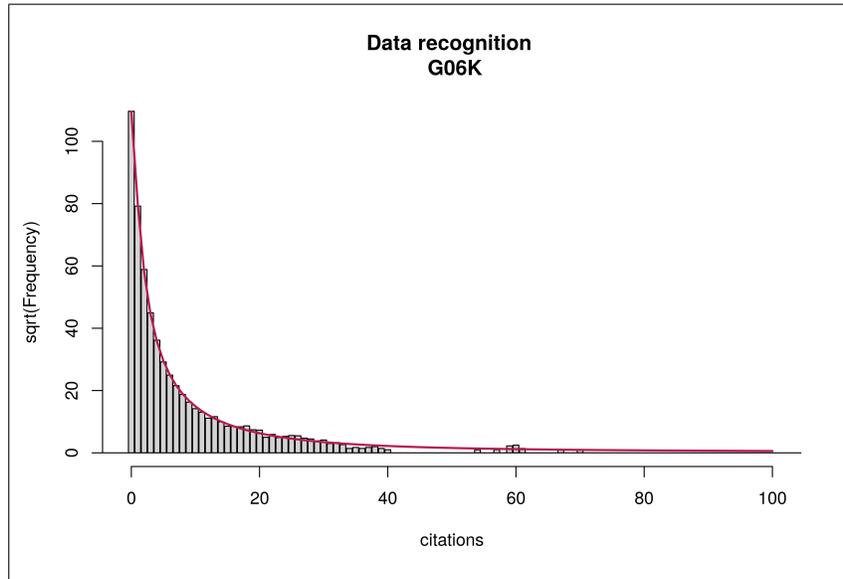


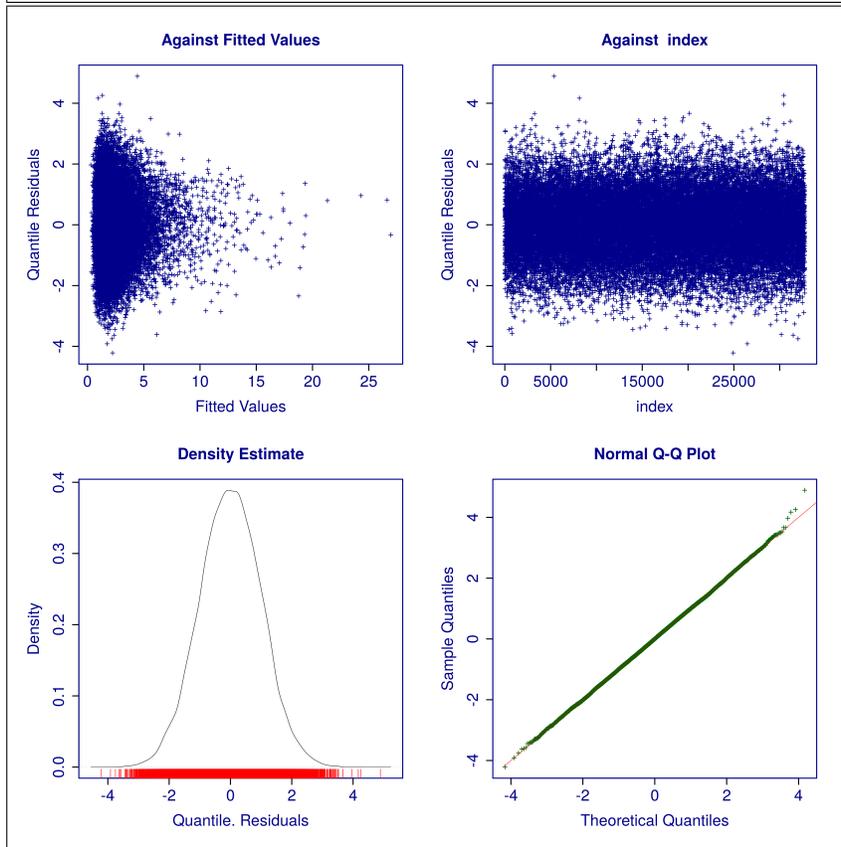
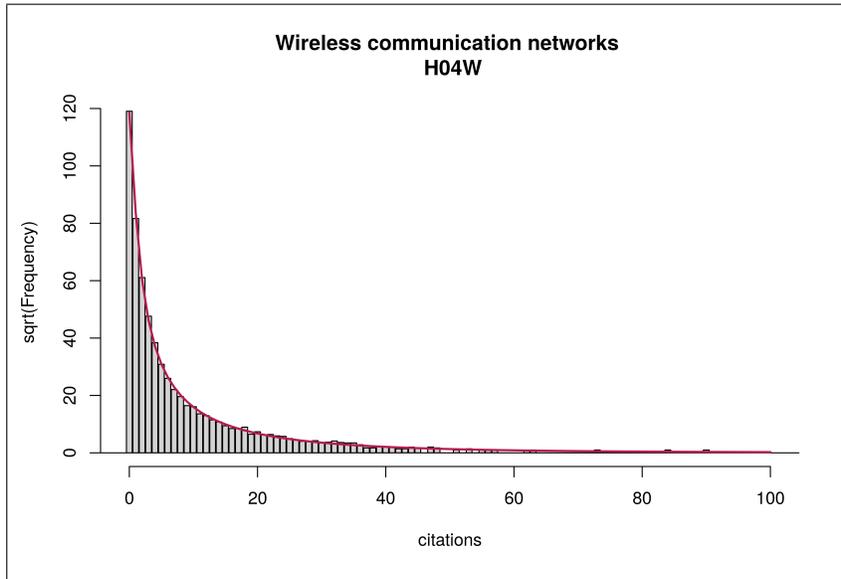


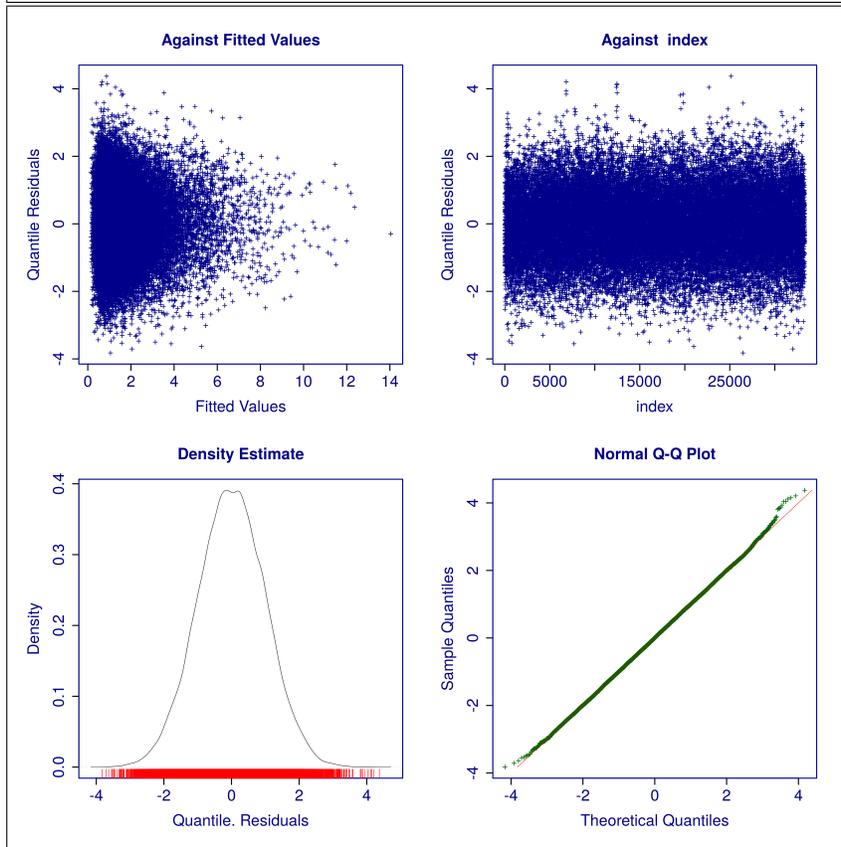
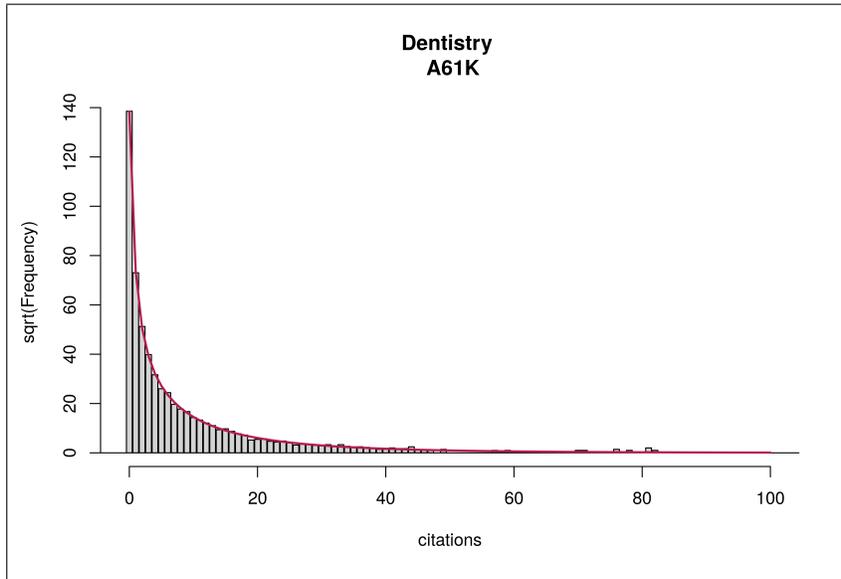


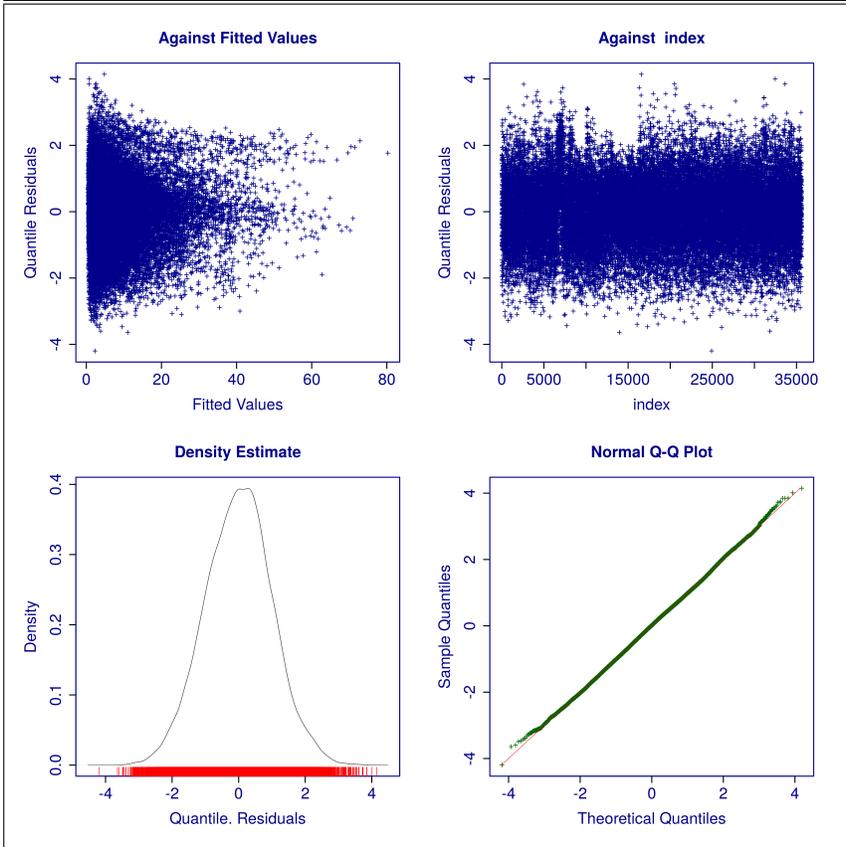
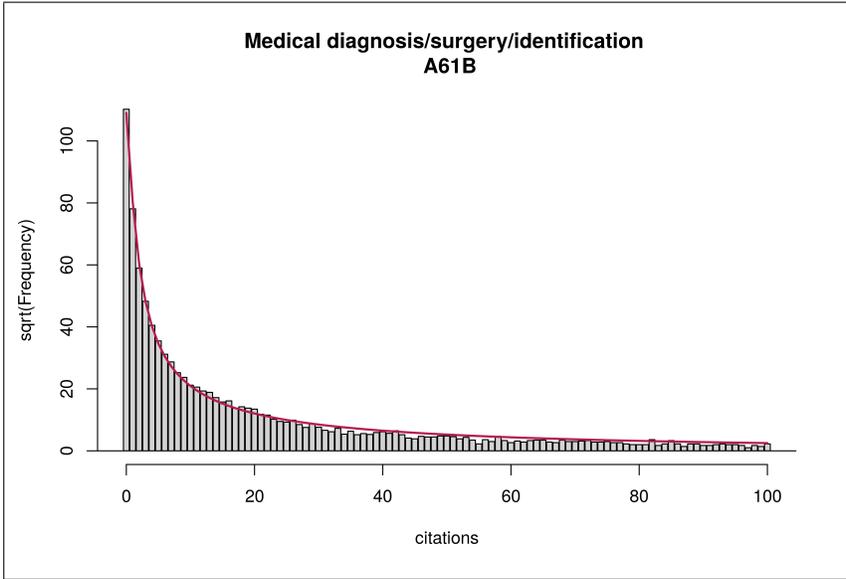


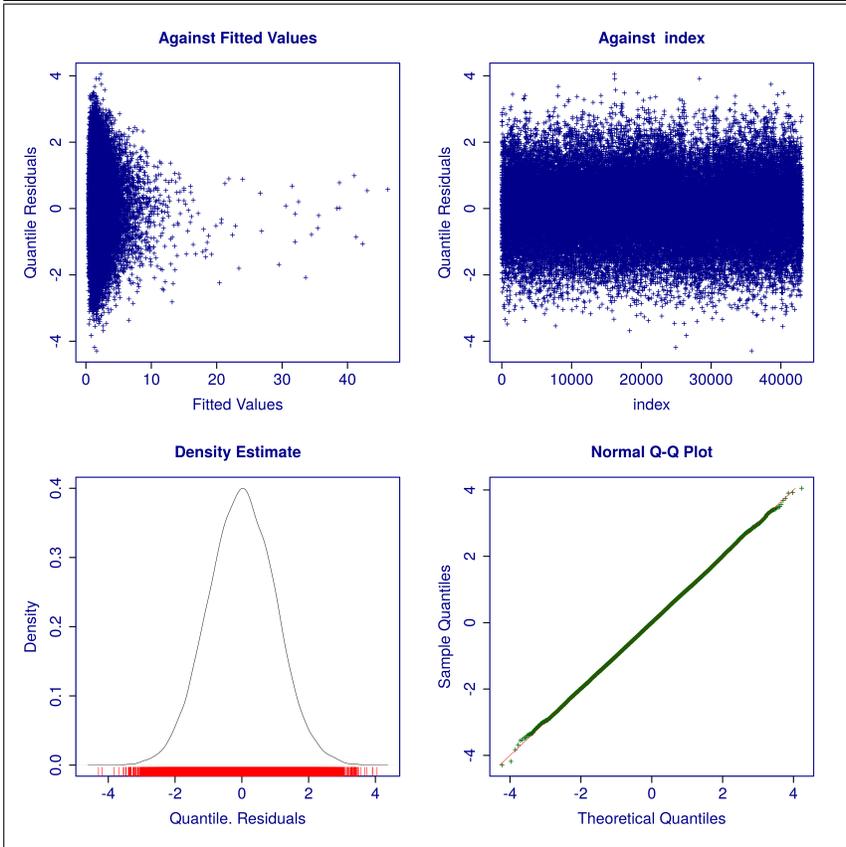
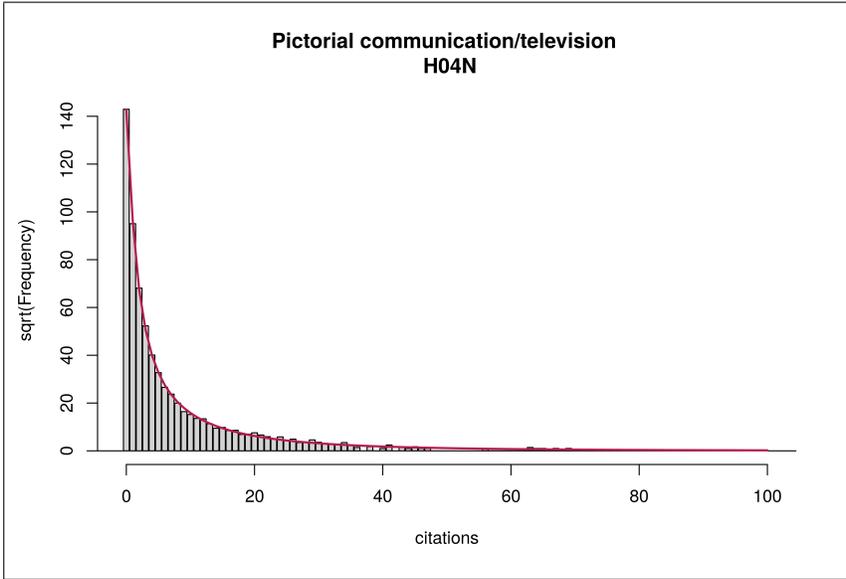


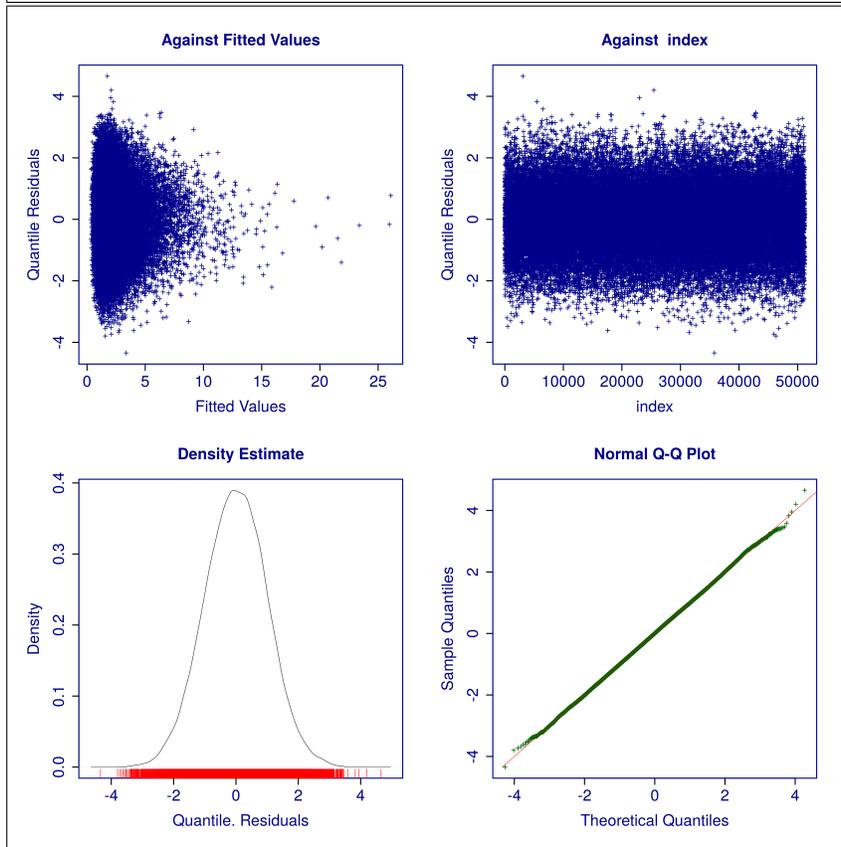
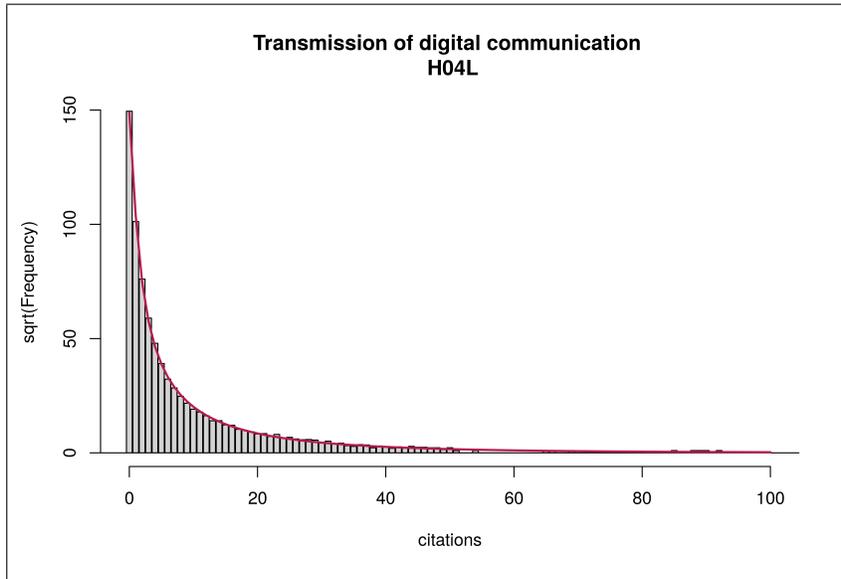


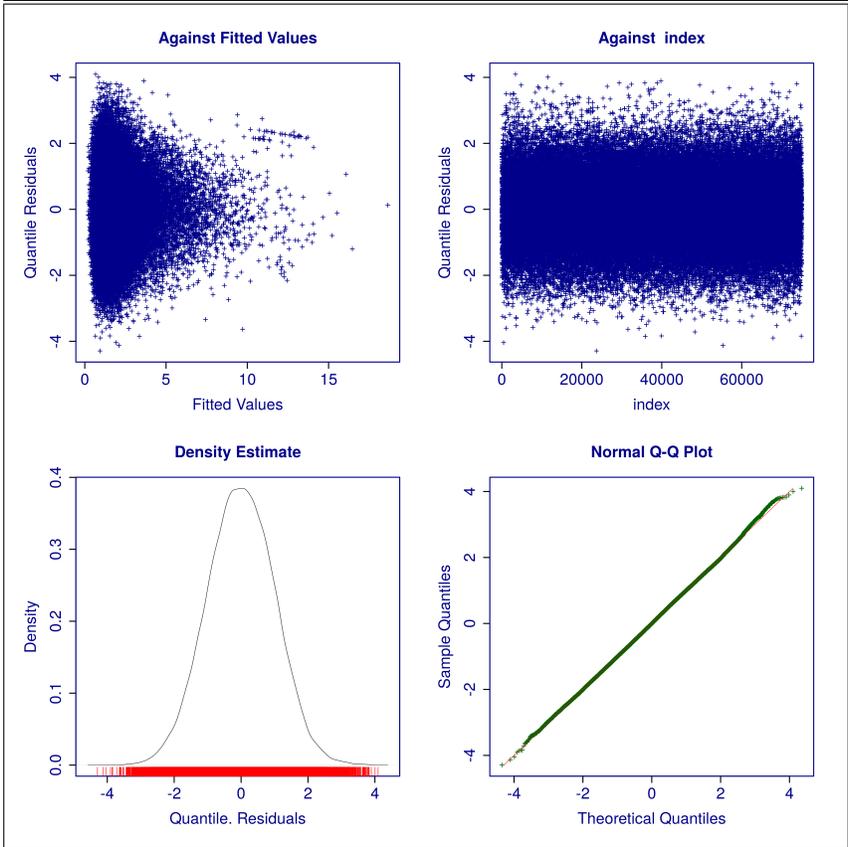
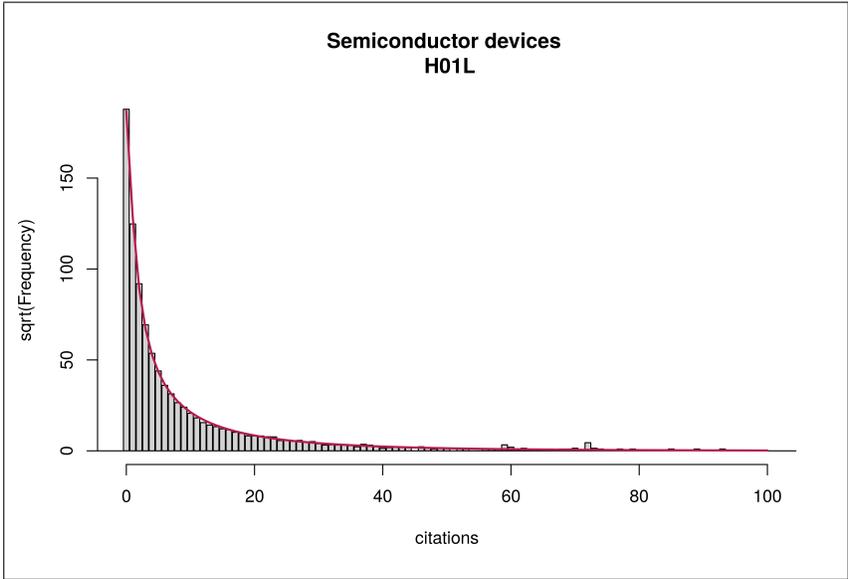


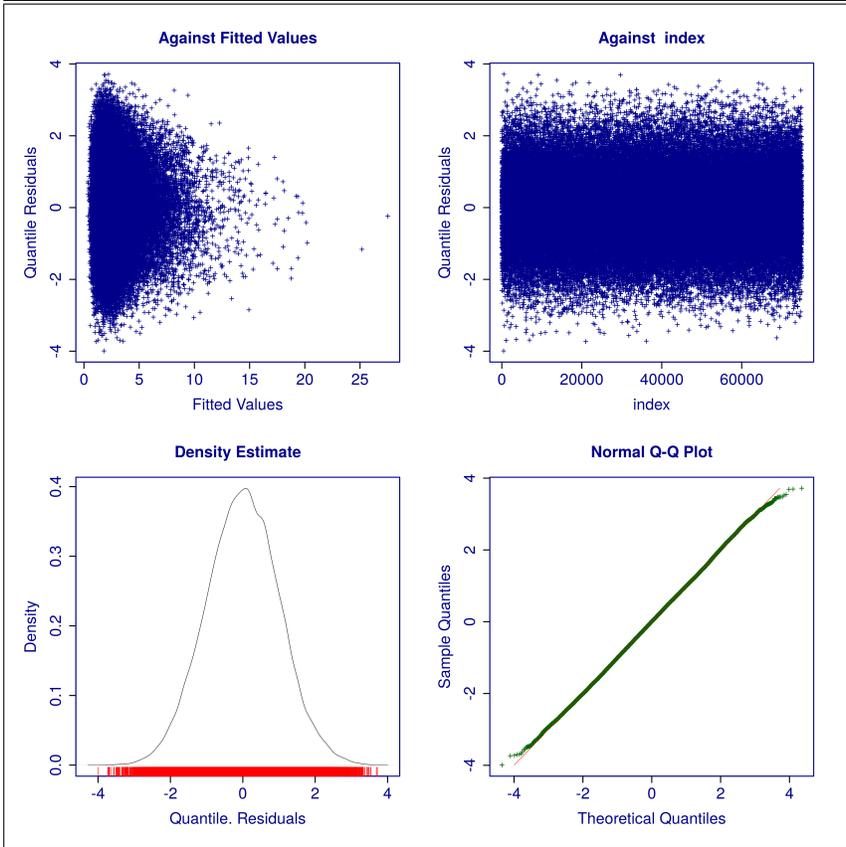
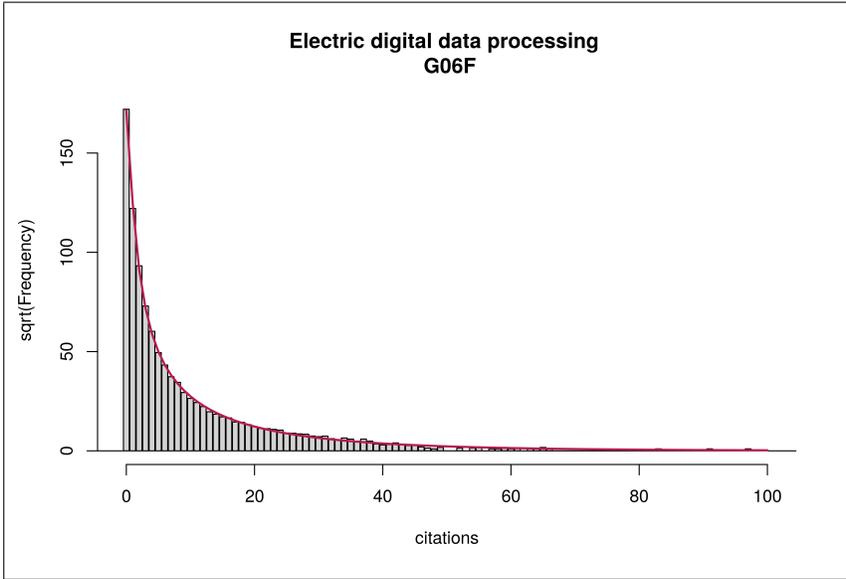




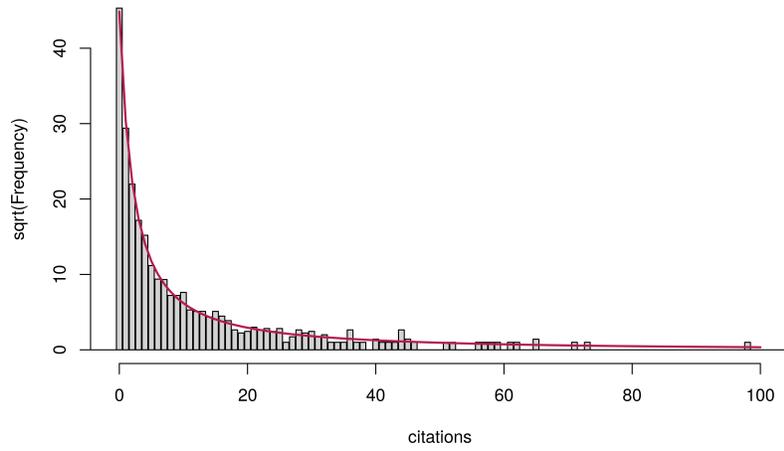




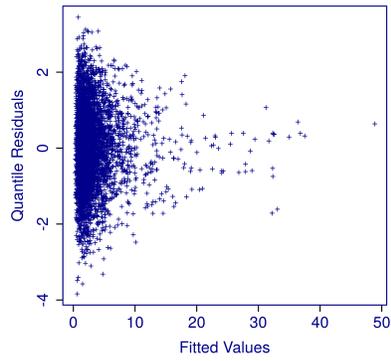




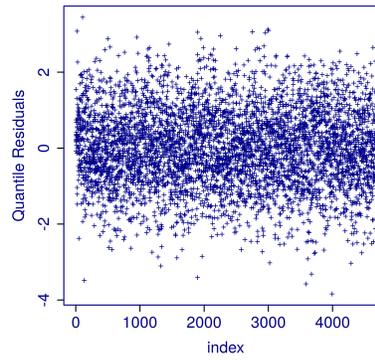
**Adaptation: Human health protection
Y02A.50**



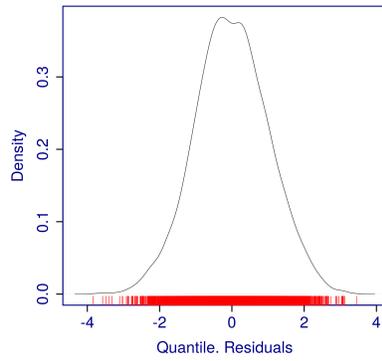
Against Fitted Values



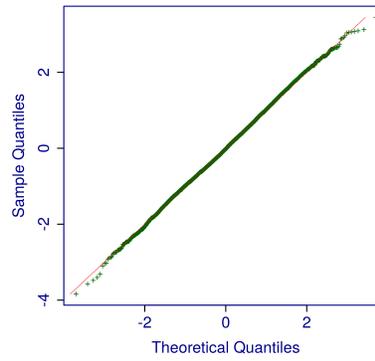
Against index



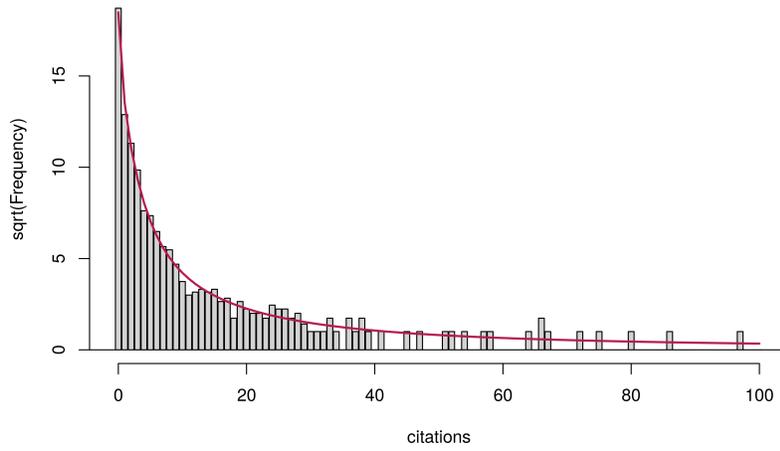
Density Estimate



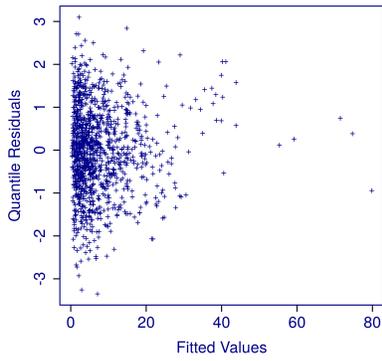
Normal Q-Q Plot



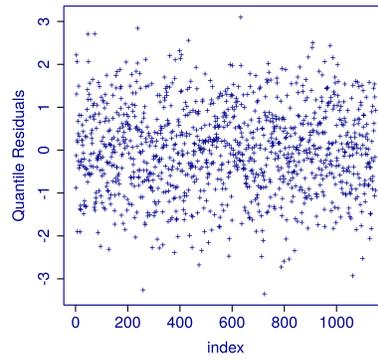
**Indirect contribution to adaptation: ICT/meteorology/healthcare/resource assessme
Y02A.90**



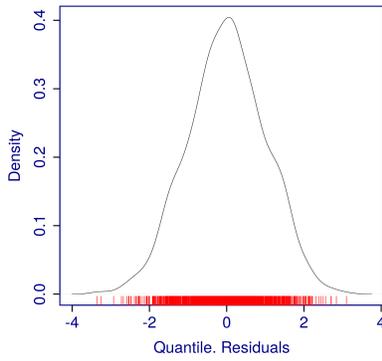
Against Fitted Values



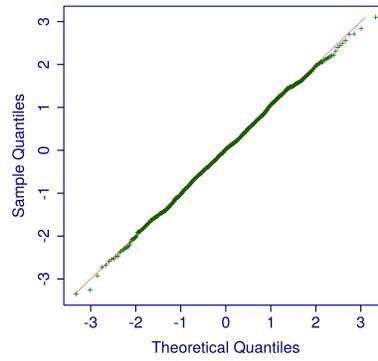
Against index



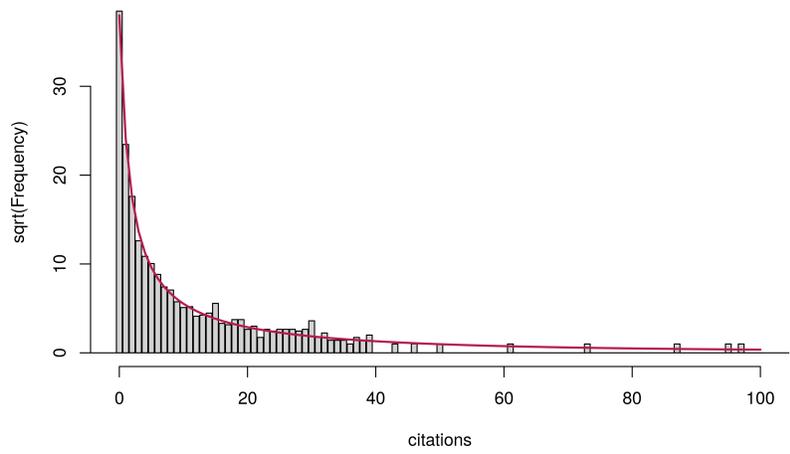
Density Estimate



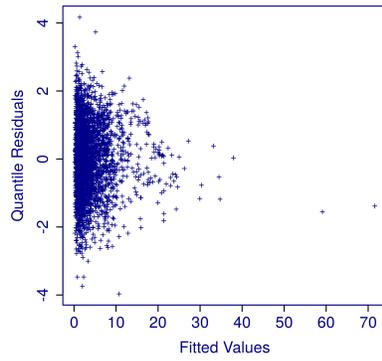
Normal Q-Q Plot



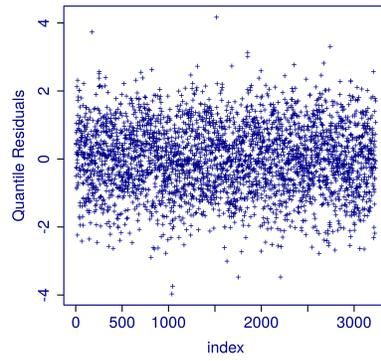
**Technologies for the production of fuel of non-fossil origin
Y02E.50**



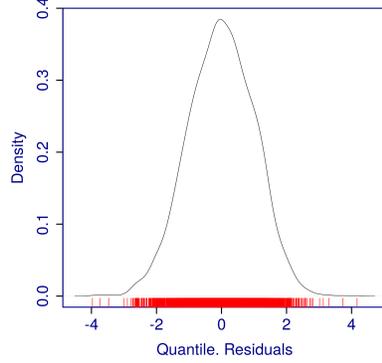
Against Fitted Values



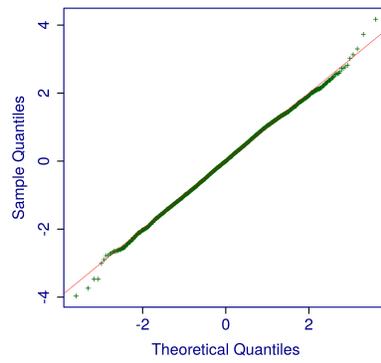
Against index



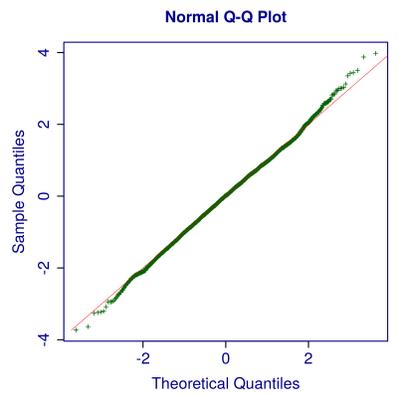
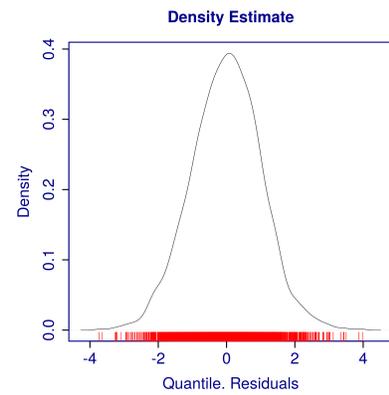
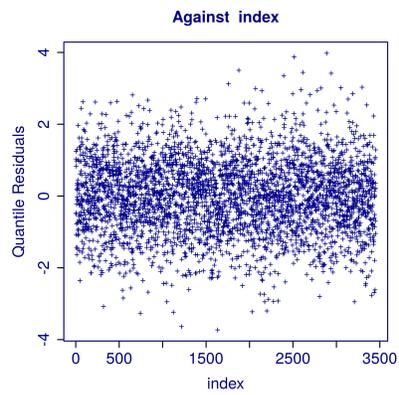
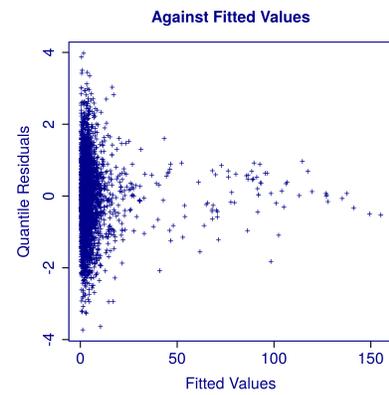
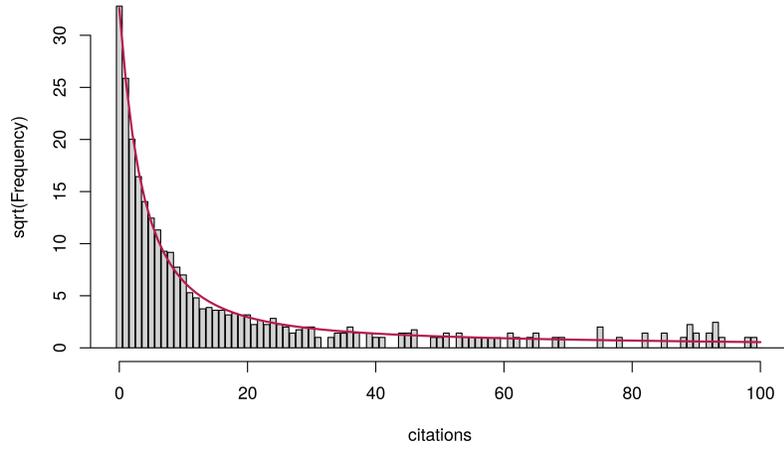
Density Estimate

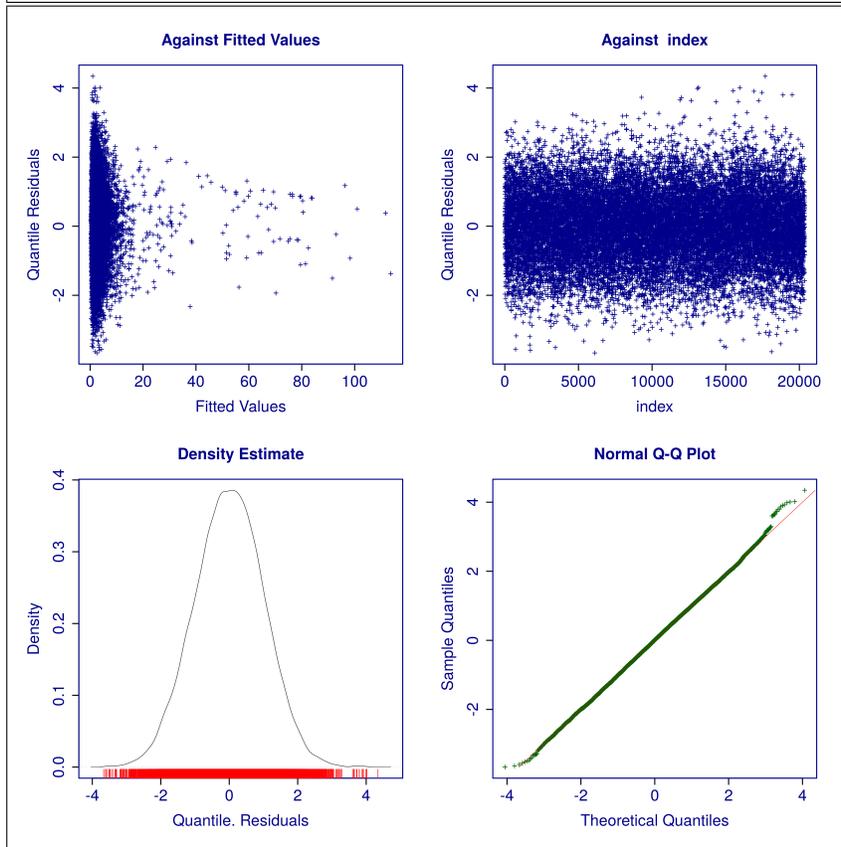
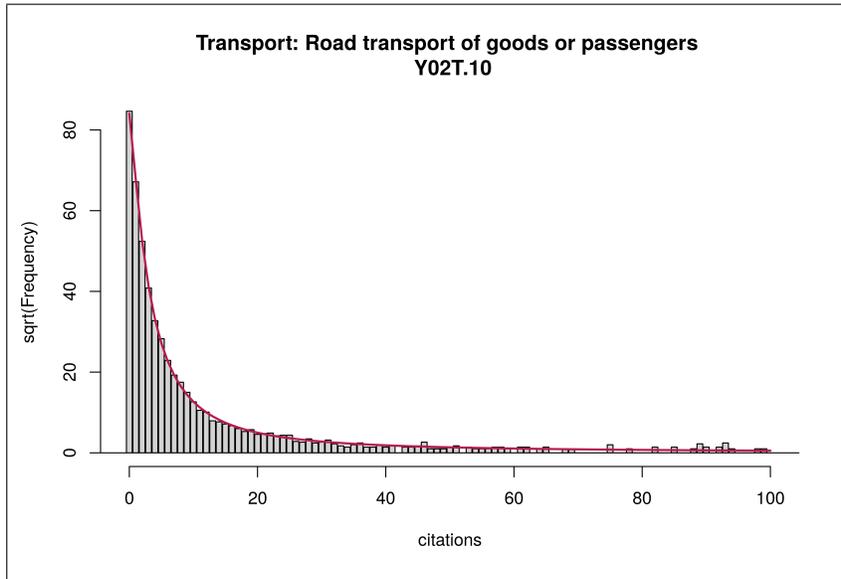


Normal Q-Q Plot

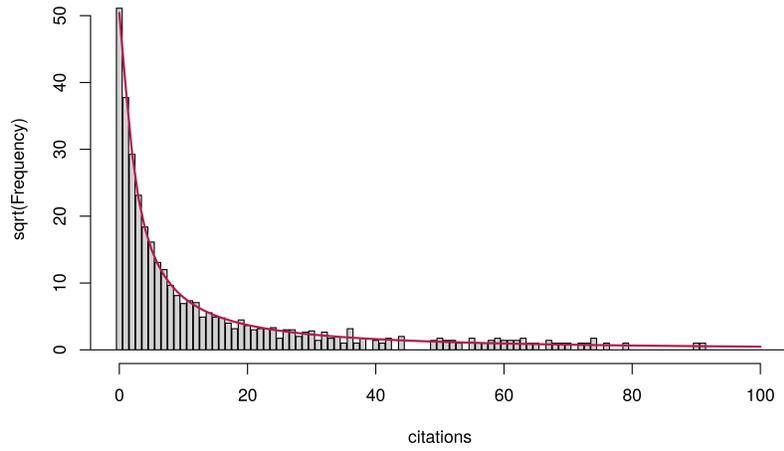


**transport: Enabling technologies (EV charging, plug-in, vehicle ICT, fuel cells, hydro
Y02T.90**

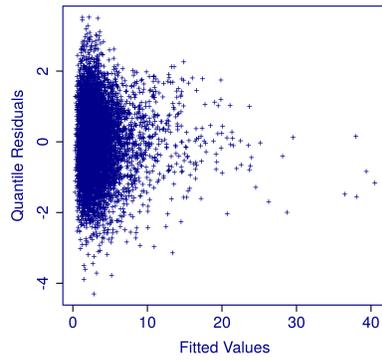




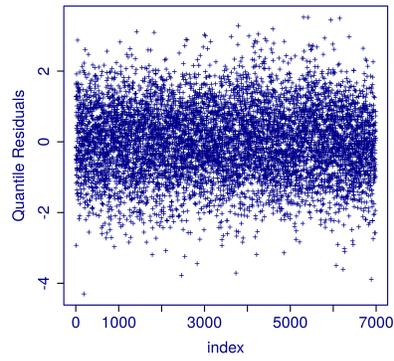
Mitigation: Energy efficient computing
Y02D.10



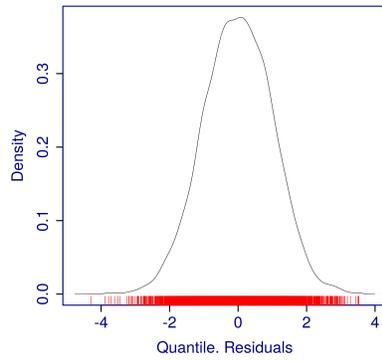
Against Fitted Values



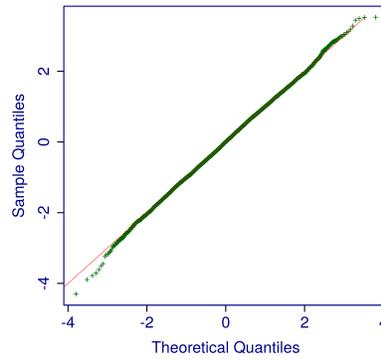
Against index

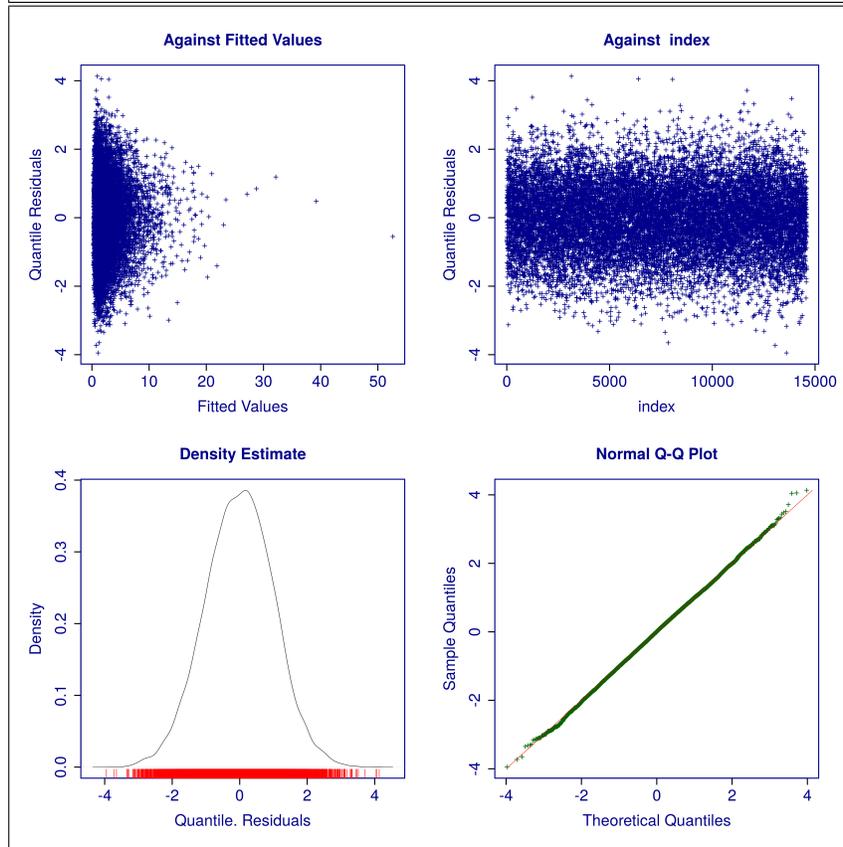
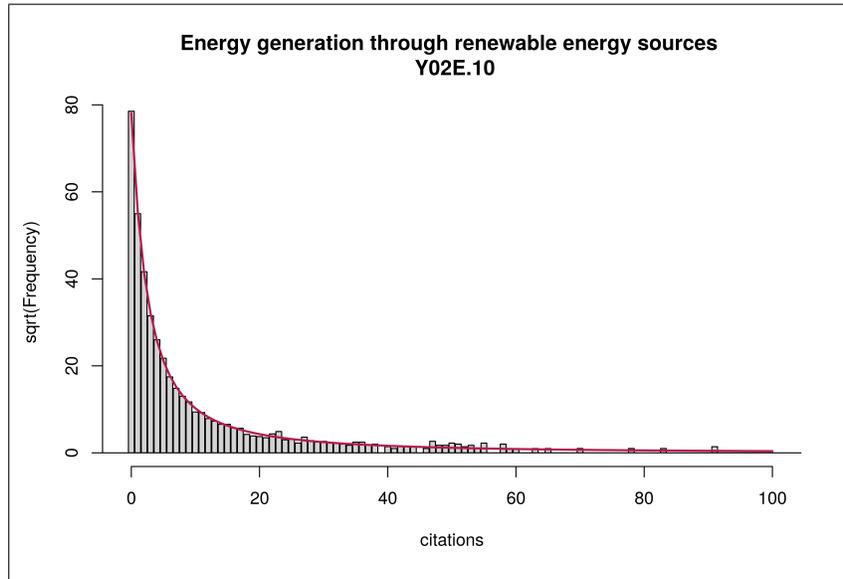


Density Estimate

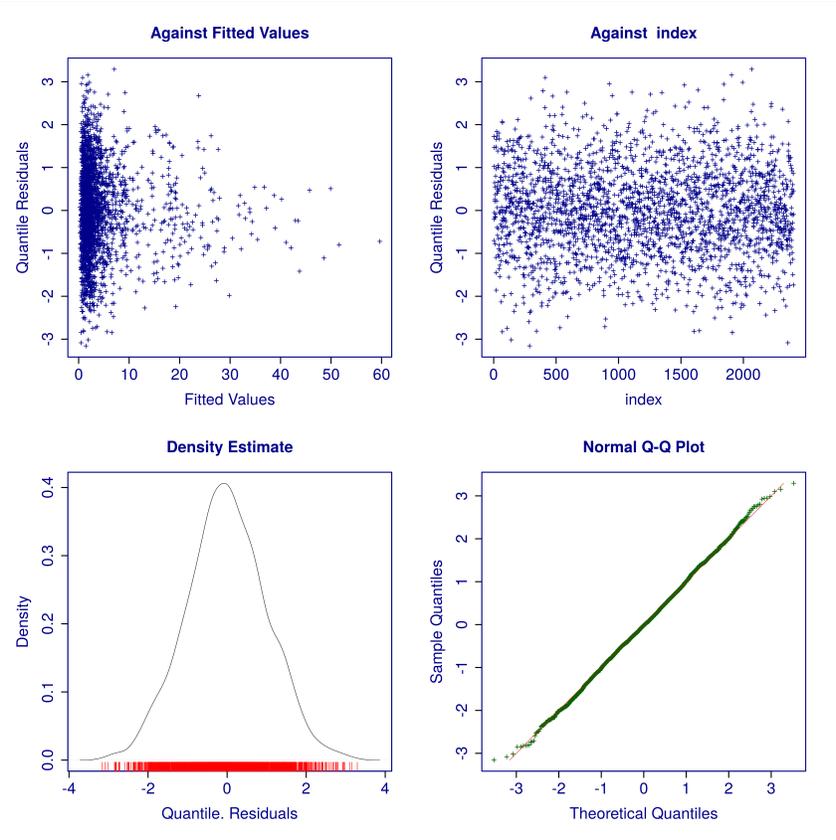
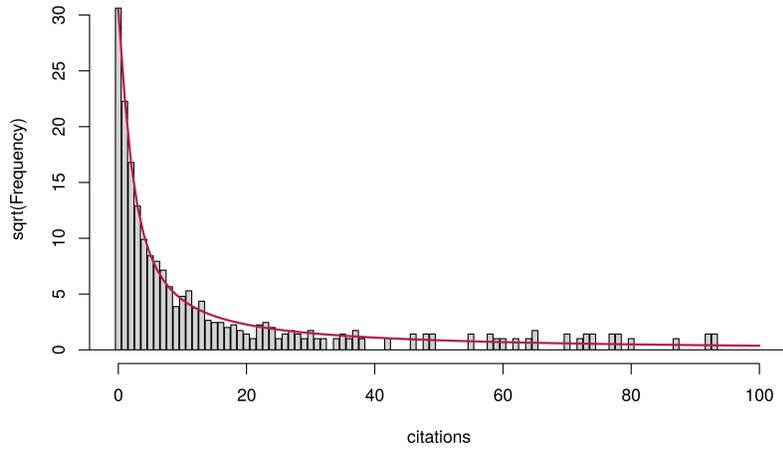


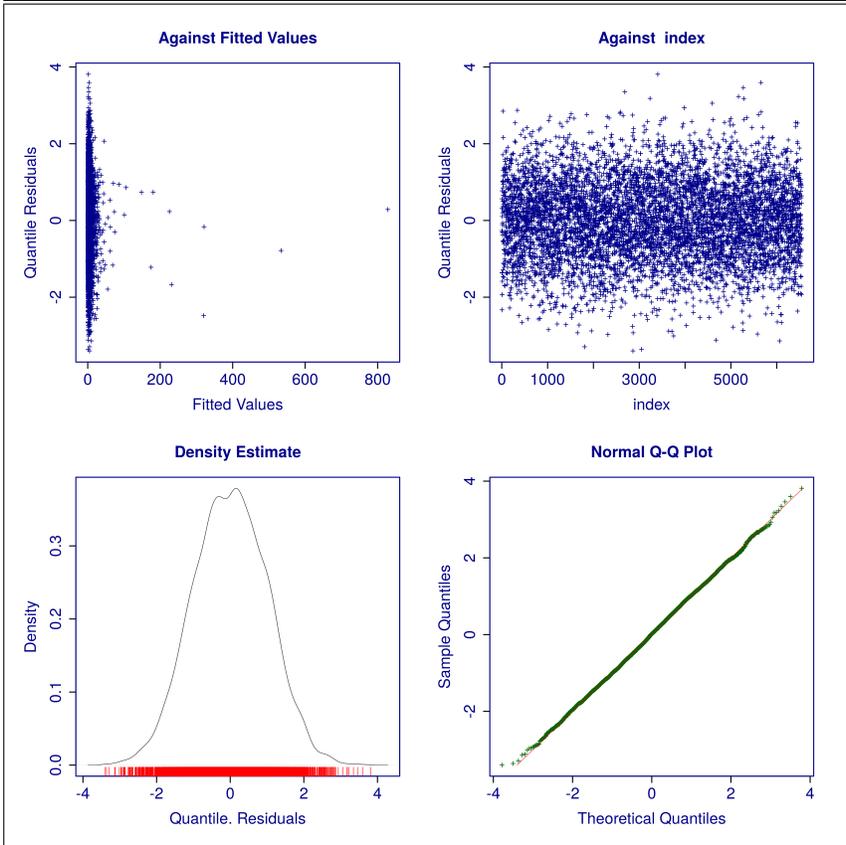
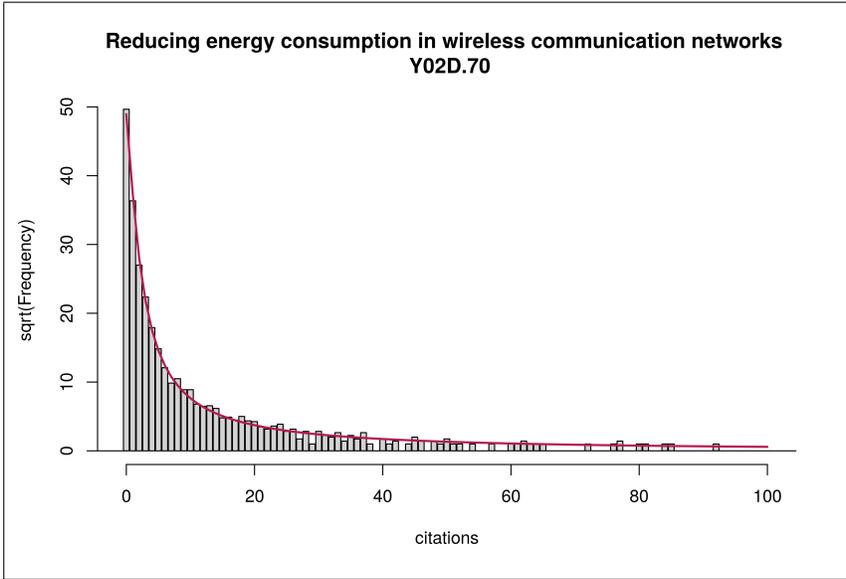
Normal Q-Q Plot



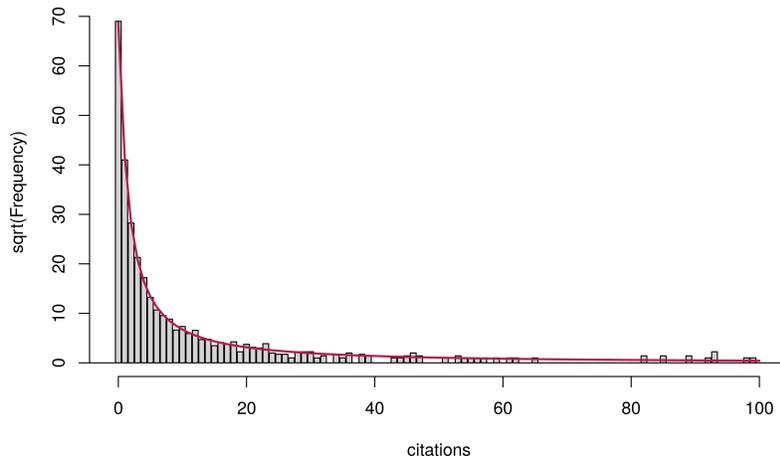


**Production: GHG emissions mitigation (factory control/smart factories/manufacturi
Y02P.90**

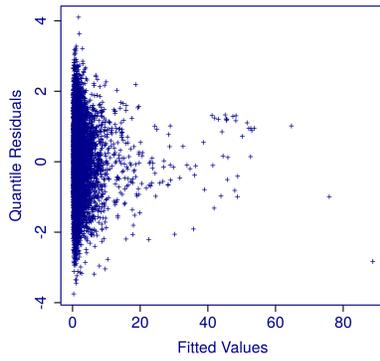




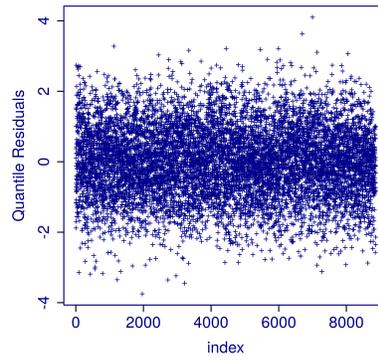
**Energy: GHG emissions mitigation (storage/hydrogen/fuel cells/power networks)
Y02E.60**



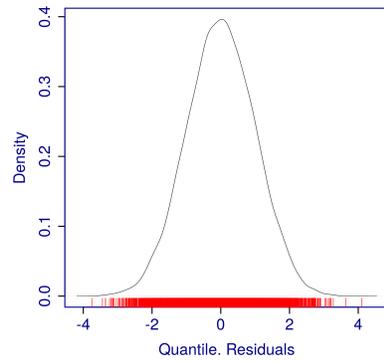
Against Fitted Values



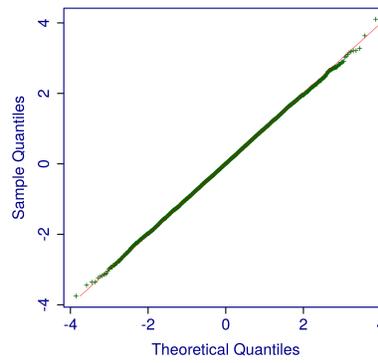
Against index



Density Estimate



Normal Q-Q Plot



3 Exploring the data and the covariates

The following data exploration examines whether the log forward citations is reasonably linear on the covariates or not. Where the covariates have been transformed based on the raw data, we found that the transformation indicated a better model fit to the covariate. We compare a linear relation to more flexible regression models to assess if the difference is large or not.

The following plots

- Show that the mean/var and skewness/kurtosis can vary much across subclasses: These will be important for modeling the scale and shape of the distributions.
- Allow roughly checking if the covariates are reasonably linear with the target on the log scale (log x shown as a log x plus 1 transformation).
 - To visually inspect how the average of citations on the log scale varies over the range of continuous/count covariates, we use local regression (local polynomials with smoothing): This is the *gam* smoother option in *ggplot2* which can visually indicate non-linearities.
 - In most cases the deviation from the linear trend is not substantial to rule out the covariate in the model.
 - The TCT is sometimes clearly non-linear: We handle this by discretizing it (a factor for different ranges, 4 levels).

