

---

## “Risk reference charts for speeding based on telematics information”

Montserrat Guillen, Ana M. Pérez-Marín and Manuela Alcañiz

---

The Research Institute of Applied Economics (IREA) in Barcelona was founded in 2005, as a research institute in applied economics. Three consolidated research groups make up the institute: AQR, RISK and GiM, and a large number of members are involved in the Institute. IREA focuses on four priority lines of investigation: (i) the quantitative study of regional and urban economic activity and analysis of regional and local economic policies, (ii) study of public economic activity in markets, particularly in the fields of empirical evaluation of privatization, the regulation and competition in the markets of public services using state of industrial economy, (iii) risk analysis in finance and insurance, and (iv) the development of micro and macro econometrics applied for the analysis of economic activity, particularly for quantitative evaluation of public policies.

IREA Working Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. For that reason, IREA Working Papers may not be reproduced or distributed without the written consent of the author. A revised version may be available directly from the author.

Any opinions expressed here are those of the author(s) and not those of IREA. Research published in this series may include views on policy, but the institute itself takes no institutional policy positions.

Reference charts are widely used as a graphical tool for assessing and monitoring children's growth given gender and age. Here, we propose a similar approach to the assessment of driving risk. Based on telematics data, and using quantile regression models, our methodology estimates the percentiles of the distance driven at speeds above the legal limit depending on drivers' characteristics and the journeys made. We refer to the resulting graphs as risk reference charts for speeding and illustrate their use for a sample of drivers with Pay-How-You-Drive insurance policies. We find that percentiles of distance driven at excessive speeds depend mainly on total distance driven, the percentage of driving in urban areas and the driver's gender. However, the impact on the estimated percentile for these covariates is not constant. We conclude that the heterogeneity in the risk of driving long distances above the speed limit can be easily represented using reference charts and that, conversely, individual drivers can be scored by calculating an estimated percentile for their specific case. The dynamics of this risk score can be assessed by recording drivers as they accumulate driving experience and cover more kilometres. Our methodology should be useful for accident prevention and, in the context of Manage-How-You-Drive insurance, reference charts can provide real-time alerts and enhance recommendations for ensuring safety.

*JEL classification:* C21, G22

*Keywords:* Motor insurance, Speed, Telematics, Quantile regression, Reference curves, Risk score

Montserrat Guillen: Dept. Econometrics, Riskcenter-IREA, Universitat de Barcelona, Av. Diagonal, 690, 08034 Barcelona, Spain. Email: [mguillen@ub.edu](mailto:mguillen@ub.edu)

Ana M. Pérez-Marín: Dept. Econometrics, Riskcenter-IREA, Universitat de Barcelona, Av. Diagonal, 690, 08034 Barcelona, Spain. Email: [amperez@ub.edu](mailto:amperez@ub.edu)

Manuela Alcañiz: Dept. Econometrics, Riskcenter-IREA, Universitat de Barcelona, Av. Diagonal, 690, 08034 Barcelona, Spain. Email: [malcaniz@ub.edu](mailto:malcaniz@ub.edu)

## *Acknowledgements*

---

The authors thank Fundación BBVA Big Data grants and ICREA Academia.

## 1. Introduction

Growth reference charts are used worldwide to provide a simple graphical tool for monitoring the evolution in children's height and weight. As such, they enable doctors and parents to track a child's estimated percentile path and observe his or her position with respect to that of their corresponding reference population of either boys or girls. Here, we seek to design a similar tool for assessing driving risk, based on the distance driven above the posted speed limit as an indicator of peril. A driver's risk evolution is then analysed with respect to total distance driven and other circumstances that need to be taken into consideration, primarily driving zone. The tool developed is both highly informative and simple, and can be directly used to communicate driving risk.

Speeding increases both the risk and severity of an accident (see Dissanayake and Lu, 2002; Ossiander and Cummings, 2002; Jun et al., 2007, 2011; Vernon et al., 2004), but drivers are not homogeneous with respect to their level of risk and driving style. Specifically, men present riskier driving patterns, driving more kilometres per day, during the night and at speeds above the limit than women (Ayuso et al., 2014, 2016a, 2016b). All these factors have been shown to be associated with a greater number of accidents (Gao et al., 2019a; Gao and Wüthrich, 2019; Guillen et al., 2019). Moreover, Paefgen et al. (2014) report that the risk of accident is higher on urban roads, during weekends, at nightfall and at low- (0–30 km/h) or high-range speeds (90–120 km/h). Indeed, Pérez-Marín and Guillen (2019) concluded that if excess speeds could be eliminated, the expected number of accident claims would be reduced by half. Interestingly, Pérez-Marín et al. (2019a) showed that young drivers tend to reduce posted speed limit violations after an accident, probably because they are more aware of the risk.

Speed and driving distance have been exhaustively analysed in transport research (see, for example, Hewson, 2008 or Plötz et al., 2017). Moreover, analyses of speeding in traffic safety research have focused not only on the average speed, but also on its quantiles. Specifically, Hewson (2008) explored the benefits of using quantile regression to evaluate whether or not an intervention is able to significantly modify the 85<sup>th</sup> percentile speed. Recently, Pérez-Marín et al. (2019b) applied quantile regression to an analysis of the effects of telematics information (location and time of driving and the total distance driven) on a range of percentiles of the distance driven at speeds above the limit by using a sample of drivers covered by a Pay-How-You-Drive (PHYD) insurance policy. In PHYD policies, the premium is calculated based on the customer's driving pattern (such as speeding, harsh acceleration, sudden braking or hard cornering). Based on these patterns, a driver's risk score can be obtained and used to calculate his or her premium (see a survey in Arumugam and Bhargavi, 2019).

In this paper, we propose a methodology for displaying percentiles that allows us to quantify a driver's risk score. To do so, we use a graphical representation of the percentiles of distance driven at speeds above the limit, depending on specific driver characteristics and on the sort of trips they make. Employing charts similar to the well-known reference curves for child growth, we develop a new methodology in the context of speeding that should prove useful when a large number of covariates can influence a driver's behaviour on the road and, hence, their risk profile.

Specifically, we call our graphs *risk reference charts* for speeding, as they provide each driver with their corresponding percentile of distance driven at speeds above the legal limits, given all available information on that driver. This proves to be a straightforward risk score for the driver. We take the article by Perez-Marín et al. (2019) as our starting point, and use the methodology proposed by Wei et al. (2006) in the context of growth charts (based on quantile regression) to produce risk reference charts for speeding. We use the same data as presented in Perez-Marín et al. (2019b) and explore alternative model formulations in the context of generalized linear models (GLMs) and quantile regression. In particular, we investigate in-depth the relationship between distance driven at speeds above the legal limits (the dependent variable in our regression models) and total distance driven. We conclude that their relationship is not linear, but exponential. This exponential relationship determines the shape of the risk reference charts for speeding. As a result, we also observe that our methodology substantially improves the initial results obtained in Perez-Marín et al. (2019b).

The rest of this paper is organized as follows. In section 2, the quantile regression model and the database used in our study are presented. In section 3, the main results of the regression models are summarized and the risk reference charts are provided. Finally, in section 4, the main results are discussed.

## 2. Material and Methods

### 2.1. Methods

Risk reference charts for speeding are obtained by means of quantile regression, where each curve corresponds to a percentile level. This type of regression analysis is flexible enough to incorporate many covariates, both qualitative and quantitative. Moreover, a web application is easily designed, so that when a user enters his or her covariate information and observed mileage above the speed limit, a graph is displayed, locating the specific driver on the chart. In this paper, we also fit a GLM model prior to quantile regression; specifically, we fit a gamma model because the dependent variable, which is mileage above the speed limit, is expected to be asymmetric. That is, while a large number of drivers can be expected not to exceed the speed limit over a certain number of kilometres, only a few are expected to exceed the limit over a high percentage of the distance driven.

The  $\tau$ -quantile of a continuous random variable  $Y$  is the value  $c_\tau$  for which  $P(Y \leq c_\tau) = \tau$ . In the financial and actuarial industries, the  $\tau$ -quantile, or the percentile at the level  $\tau$ , is known as the value-at-risk at level  $\tau$ . Quantile regression is used in order to estimate conditional quantiles, as the model assumes that the  $c_\tau(Y)$  depends on certain explanatory variables. Specifically,

$$c_\tau(Y_i | X_{1i}, \dots, X_{ki}) = \beta_0^\tau + \beta_1^\tau X_{1i} + \beta_2^\tau X_{2i} + \dots + \beta_k^\tau X_{ki}, \quad (1)$$

where  $Y_i$  is the dependent variable for the  $i$ -th individual, with  $i=1, \dots, n$ , and  $X_{ji}$  are the observations of the explanatory variables, with  $j=1, \dots, k$ . It can be proved (Koenker and Bassett, 1978) that

$$\widehat{\beta}^\tau = \underset{b}{\operatorname{argmin}} \left[ \sum_{Y_i \geq X_i' b} \tau |Y_i - X_i' b| + \sum_{Y_i < X_i' b} (1 - \tau) |Y_i - X_i' b| \right]. \quad (2)$$

The objective function (2) corresponds to the sum of  $n$  components, called  $\rho_\tau(Y_i - X_i' b)$  that are expressed as follows:

$$\begin{aligned} \rho_\tau(Y_i - X_i' b) &= \tau(Y_i - X_i' b)I_{\{Y_i \geq X_i' b\}} + (\tau - 1)(Y_i - X_i' b)I_{\{Y_i < X_i' b\}} = \\ &= (Y_i - X_i' b)(\tau - I_{\{Y_i < X_i' b\}}), \end{aligned} \quad (3)$$

where  $I_{\{\cdot\}}$  is an indicator function equal to 1 if the condition in the subindex is fulfilled, and 0 otherwise. A quantile regression model can be easily fitted, for example in R, by using the function *qr* of the *quantreg* R package (Koenker et al., 2018).

Koenker and Machado (1999) proposed an expression to measure the goodness-of-fit of the quantile regression based on a comparison of the values of the objective functions of the estimated model and of the constrained model that only includes an intercept term. Specifically, let

$$\widehat{V}(\tau) = \sum_{i=1}^n \rho_\tau(Y_i - X_i' \widehat{\beta}^\tau) \quad (4)$$

be the value of the objective function of the estimated model and

$$\widetilde{V}(\tau) = \sum_{i=1}^n \rho_\tau(Y_i - \beta_0^\tau) \quad (5)$$

be the value of the objective function of the constrained model that only includes the intercept term. Then, the goodness-of-fit measure proposed by Koenker and Machado (1999) is

$$R^1(\tau) = 1 - \widehat{V}(\tau) / \widetilde{V}(\tau) \quad (6)$$

which is similar to the  $R^2$  in the multiple linear regression model. Additional details of quantile regression implementation in R can be found in Uribe and Guillen (2020).

## 2.2. Data

The dataset used in this article is the same as that employed in Pérez-Marín et al. (2019b). Our sample consists of 9,585 drivers aged 35 years or less, with PHYD coverage during the whole of 2010. Data were provided by a Spanish insurer. The description of the variables is presented in Table 1. We know the gender (variable *Gender*) and age of the driver at the beginning of 2010 (variable *Age*). Additionally, we also know the total number of kilometres driven during 2010 (*Km*), the number of kilometres driven at speeds above the posted limit (*Tolerkm*, which is our dependent variable), the percentage of kilometres driven on urban roads (*Urban*) and, finally, the percentage of kilometres driven at night (*Night*). In order to fit the gamma model, note that 29 observations with zero kilometres driven at speeds above the posted limit – 0.3% of the sample size – had to be removed from the original sample.

Table 1. Description of variables used in the insurance dataset

Variable	Description
<i>Tolerkm</i>	Number of kilometres driven at speeds above the posted limit during 2010
<i>Km</i>	Total number of kilometres driven during 2010*
<i>Urban</i>	% of kilometres driven on urban roads during 2010*
<i>Night</i>	% of kilometres driven at night (between midnight and 6 am) during 2010
<i>Age</i>	Age of the driver at the beginning of 2010
<i>Gender</i>	1 = male, 0 = female

\*Power transformations were used in the gamma model,  $Km\_tg = Km^{0.1}$  and  $Urban\_tg = Urban^{0.7}$ , and in the quantile regression models,  $Km\_tqr = Km^{1.7}$  and  $Urban\_tqr = Urban^{0.1}$

As shown in Table 2, *Tolerkm* presents a positive asymmetry (skewness coefficient = 3.64), with a long tail. The sample comprises 49% women and 51% men. The average age of drivers is 24.78 years. The average number of kilometres driven during the observed year was 13,099.91 (standard deviation of 7,698.98). On average, drivers travelled 26.2% of kilometers on urban roads, 7.02% of kilometers at night and 1,402.44 kilometers at speeds above the limit.

Table 2. Descriptive statistics of the insurance data set

	Min	1 <sup>st</sup> Quartile	Median	Mean	3 <sup>rd</sup> Quartile	Max	St. Dev.	Skewness
<i>Tolerkm</i>	0.03	285.78	692.92	1,402.44	1,710.44	23,500.19	1,996.90	3.64
<i>Km</i>	27.79	7,575.15	11,719.83	13,099.91	17,350.12	57,756.98	7,698.98	1.08
<i>Urban</i>	0.00	15.59	23.36	26.20	34.25	96.41	14.04	0.99
<i>Night</i>	0.00	2.49	5.32	7.02	9.85	78.56	6.12	1.67
<i>Age</i>	18.11	22.66	24.63	24.78	26.88	35.00	2.82	0.11

### 3. Results

In order to predict *Tolerkm*, we employed a gamma regression model<sup>1</sup>, and used different transformations of *Km* and *Urban* (including logarithmic and power transformations). The other two continuous explanatory variables (*Night* and *Age*) were not transformed, as this had almost no impact on the Akaike information criterion (AIC). The transformation of *Km* and *Urban* that produced the lowest AIC score (equal to 149,299.5) was the combination of the following power transformations<sup>2</sup>:  $Km\_tg = Km^{0.1}$  and  $Urban\_tg = Urban^{0.7}$ .

The parameter estimates of the corresponding gamma regression model are shown in Table 3. Coefficient estimates with a p-value lower than 1% correspond to gender, the transformed total number of kilometres driven (*Km\_tg*) and the transformed percentage of kilometres driven in urban areas (*Urban\_tg*). Age effect is only significant at the 10% level (p-value=0.0807), probably because the insurance policies were sold exclusively to young drivers. Likewise, the positive effect of percentage of kilometres driven at night (*Night*) is only significant at the 10% level (p-value=0.0987), which would indicate that drivers with a higher percentage of night time driving tend to have an average excess speed distance greater than those with a lower percentage of night time driving. *Km\_tg* has a positive parameter estimate, indicating that an increase in the total number of kilometres driven contributes to increasing the expected number of kilometres driven at speeds above the posted limits. In contrast, *Urban\_tg* presents the opposite effect: the higher the percentage of kilometres driven on urban roads, the lower the expected number of kilometres driven at speeds above the posted limit. Finally, gender (baseline reference: female) has a positive parameter estimate, indicating that men seem to drive more kilometres at speeds above the posted limit than women.

Table 3. Results of the gamma regression model for the insurance data set. Dependent variable is the number of kilometres driven above posted speed limits

	Parameter estimate (p-value)
<i>Intercept</i>	-5.126659 (<0.0001)
<i>Km_tg</i>	4.966361 (<0.0001)
<i>Urban_tg</i>	-0.065209 (<0.0001)
<i>Night</i>	0.002475 (0.0987)
<i>Age</i>	-0.005587 (0.0807)
<i>Gender</i>	0.207654 (<0.0001)

To estimate the quantile regressions, we also tried using other transformations of *Km* and *Urban*, selecting those that minimize the AIC score of the regression model. These transformations were  $Km\_tqr = Km^{1.7}$  and  $Urban\_tqr = Urban^{0.1}$ .

The parameter estimates and goodness-of-fit of the quantile regression models at different levels ( $\tau = 0.5, 0.75, 0.90, 0.95, 0.975, 0.99$ ) are shown in Table 4. We see that *Km\_tqr* has a significant effect, with a positive parameter estimate, for all levels of the quantile. This means that, for a specific quantile, increasing the total number of kilometres driven increases the quantile of the number of kilometres driven at speeds above the posted limits, *ceteris paribus*. In contrast, while *Urban\_tqr* also has a significant effect, it has a negative

<sup>1</sup> We also used lognormal (but it provided a higher AIC score) and inverse Gaussian regressions (but it was eventually discarded because of convergence problems in the algorithm).

<sup>2</sup> We also tried other combinations of power transformations on *Km* and *Urban*, specifically,  $Km^i$  and  $Urban^j$  where  $i = 0.05$  to  $0.5$  increasing by  $0.05$ , and  $j = 0.1$  to  $1$  increasing by  $0.1$ .

parameter estimate. Thus, as the percentage of kilometres driven in urban areas increases, the quantile of the number of kilometres driven at speeds above the limits decreases. *Night* has a significant effect only when estimating the median of the kilometres driven at speeds above the limits, but for other levels of the quantile, it has no significant effect. In the case of the median, the parameter estimate is positive, indicating that increasing the percentage of kilometres driven at night increases the median kilometres driven at speeds above the limits. *Age* has a significant effect only when estimating the quantiles at the 95<sup>th</sup> and 97.5<sup>th</sup> levels. In both cases, the corresponding parameters are positive; thus, increasing the driver's age also increases the corresponding percentiles of the distance driven at speeds above the limits. Finally, gender (baseline reference: female) has a significant parameter for all levels of the quantiles up to the 95<sup>th</sup>. The coefficient is positive; thus, men have higher percentile values of distance driven at speeds above the limits than women. In the case of the goodness-of-fit criterion, it is apparent that the contribution explaining the quantiles of the model with covariates vs. the model without increases with the increase in percentile level, reaching 61.22% at the 99<sup>th</sup> level. Additionally, in Figure A1 in the Appendix we also provide the marginal effect (estimated parameter) of each explanatory variable in the quantile regression models, as a function of the level of the estimated quantile, showing that the impact of covariates on different percentile levels is not always constant, which highlights the great utility of reference charts as graphical tools.

Figure 1 shows the risk reference charts for speeding for males and females, respectively, together with the sample data. The plots show *Tolerkm* vs. *Km*, and additionally the grey lines represent the estimated quantiles at different levels. The red line represents the conditional of *Tolerkm* estimated using the gamma regression model in Table 3. In Figure 1, the values of *Urban*, *Night* and *Age* have been fixed at the mean values in the sample for men and women, respectively. Note that the transformed variable  $Km\_tqr = Km^{1.7}$  introduced in the quantile regression models captures the shape of the scatter plot correctly. Table 5 provides various examples of percentiles obtained when using the speed reference curves in Figure 1. For example, if a male driver drives 2,000 km per year at speeds above the limits, he is in the 90<sup>th</sup> percentile curve if he drives 10,000 km per year. On the other hand, the same driver is in the 54<sup>th</sup> percentile curve if he drives 20,000 km per year, and finally, he is in the 29<sup>th</sup> percentile curve if he drives 30,000 km per year. The corresponding percentiles for women are also shown in Table 5, and are very similar if just a little higher, indicating that women seem to drive at speeds above the posted limit speed less than men.

Table 4. Parameter estimates of the quantile regression model for different percentiles of mileage above the speed limit

	50 <sup>th</sup>	75 <sup>th</sup>	90 <sup>th</sup>	95 <sup>th</sup>	97.5 <sup>th</sup>	99 <sup>th</sup>
	percentile (p-value)	percentile (p-value)	percentile (p-value)	percentile (p-value)	percentile (p-value)	percentile (p-value)
<i>Intercept</i>	1077.1401 (<0.00001)	3359.33609 (<0.00001)	6581.18447 (<0.00001)	8506.75790 (<0.00001)	10173.74522 (<0.00001)	11288.44387 (<0.00001)
<i>Km_tqr</i>	0.00008 (<0.00001)	0.00013 (<0.00001)	0.00020 (<0.00001)	0.00024 (<0.00001)	0.00028 (<0.00001)	0.00032 (<0.00001)
<i>Urban_tqr</i>	-739.33632 (<0.00001)	-2285.83883 (<0.00001)	-4529.95580 (<0.00001)	-5938.64364 (<0.00001)	-7176.65122 (<0.00001)	-7960.61577 (<0.00001)
<i>Night</i>	2.36200 (0.00224)	1.08333 (0.43992)	-0.94645 (0.64086)	4.77960 (0.18409)	9.13793 (0.26052)	-0.57912 (0.95843)
<i>Age</i>	-1.80286 (0.20585)	-1.22676 (0.70373)	6.18758 (0.25224)	17.61657 (0.02532)	28.01504 (0.00894)	37.93405 (0.09753)
<i>Gender</i>	104.81103 (<0.00001)	167.33072 (<0.00001)	167.22760 (<0.00001)	140.47488 (0.00360)	101.91653 (0.16106)	189.44456 (0.23436)
<i>Goodness of fit (%)</i>	23.43	33.59	44.56	50.89	55.55	61.22

As discussed above, the speeding reference curves shown in Figure 1 have been obtained by assuming that the other explanatory variables (*Urban*, *Night* and *Age*) are equal to the corresponding sample means for men and women, respectively. In Figure 2 we show how these reference curves for the 95<sup>th</sup> percentile change for different values of *Urban* (which is the most relevant explanatory variable, apart from *Km*, for explaining *Tolerkm*). Specifically, for men we considered values of *Urban* equal to 8.70, 23.45 and 52.47% (5<sup>th</sup>, 50<sup>th</sup> and 95<sup>th</sup>



percentiles of *Urban* in the male sample, respectively), and we refer to these values as low, median and high levels of urban driving. Similarly, for women we considered values of *Urban* equal to 8.37, 23.01 and 53.81% (5<sup>th</sup>, 50<sup>th</sup> and 95<sup>th</sup> percentiles of *Urban* in the female sample, respectively), and similarly we refer to them as low, median and high levels of urban driving. The corresponding reference speed curves for the 95<sup>th</sup> percentile of *Tolerkm* are represented in Figure 2 for men and women, respectively, where the red lines are used to represent the corresponding curves for the average values obtained with the gamma regression model (Table 3). We observe that, as the percentage of urban driving increases, all curves move downwards, as *Urban\_tqr* has a negative coefficient. Specifically, in Table 6 we show some examples of the 95<sup>th</sup> percentile of *Tolerkm* for certain values of *Urban* and *Km*. For a male driver driving 10,000 km per year the 95<sup>th</sup> percentile of *Tolerkm* is equal to 1,753.22 km if he has a high percentage of urban driving, 2,466.03 km if he has a median percentage and 3,234.56 km if he has a low percentage. When the distance driven by the male driver increases to 20,000 km per year, then the 95<sup>th</sup> percentile of *Tolerkm* is equal to 5,120.66, 5,803.48 and 6,572.01 km for high, median and low percentages of urban driving, respectively. Table 6 also shows the results corresponding to women drivers, and we observe that they are slightly lower than those for male drivers.

Table 5. Percentiles obtained using the risk reference charts for examples of speeding (*Tolerkm*) and total distance driven (*Km*) by gender.

			<i>Km</i>		
			10,000	20,000	30,000
<i>Tolerkm</i>	Men	2,000	90 <sup>th</sup>	54 <sup>th</sup>	29 <sup>th</sup>
		3,000	98 <sup>th</sup>	74 <sup>th</sup>	44 <sup>th</sup>
	Women	2,000	92 <sup>th</sup>	57 <sup>th</sup>	30 <sup>th</sup>
		3,000	98 <sup>th</sup>	75 <sup>th</sup>	45 <sup>th</sup>
Other explanatory variables (Urban, Night and Age) are equal to their sample means for men and women, respectively.					

Figure 1. Risk reference chart for speeding for male drivers (left) and female drivers (right). *Tolerkm* vs. *Km*, where grey lines represent the estimated quantiles at different levels. The red line represents the mean of *Tolerkm* estimated using the gamma regression model in Table 3. Age, urban and night driving are fixed at the mean level by gender.

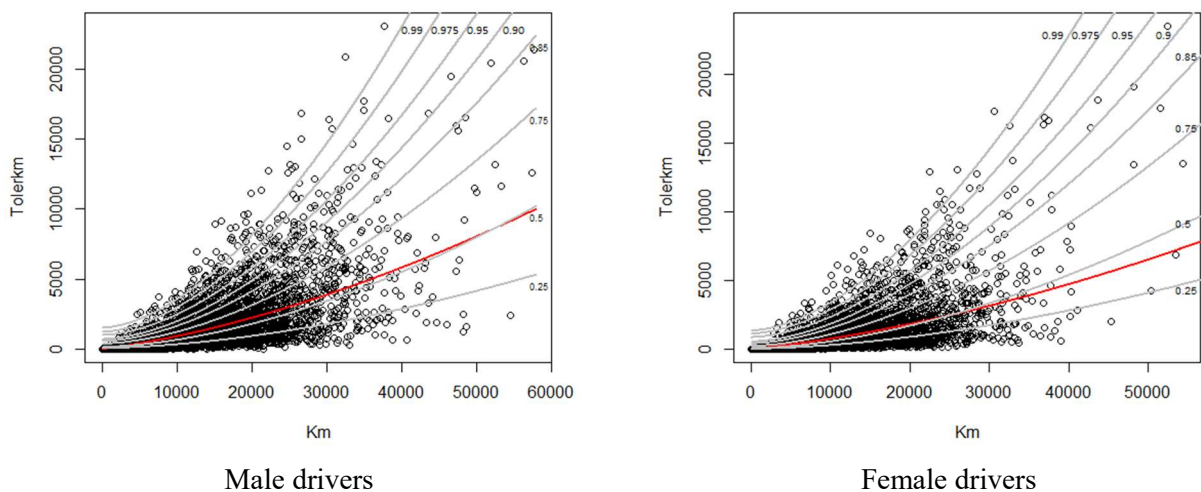
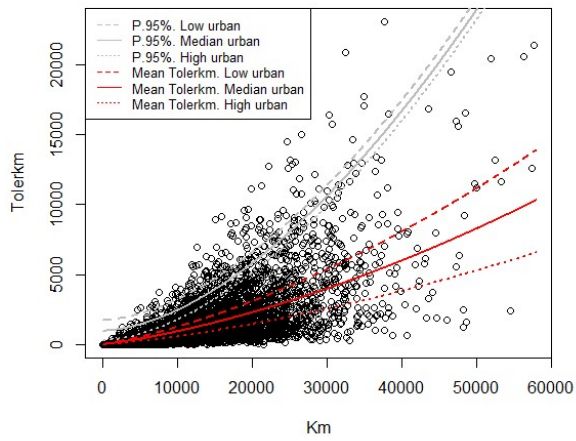
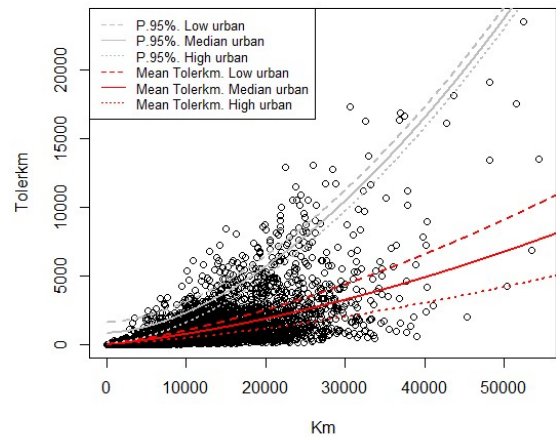


Figure 2. Risk reference chart for speeding for male drivers (left) and female drivers (right) at the 95<sup>th</sup> level. *Tolerkm* vs. *Km*, where grey lines represent the estimated 95<sup>th</sup> percentile and red lines represent the mean of *Tolerkm* estimated using the gamma regression model, for different values of *Urban* (dashed = low level of urban driving, solid = median level, and dotted = high level). Age and night driving are fixed at the sample mean level by gender.



Male drivers



Female drivers

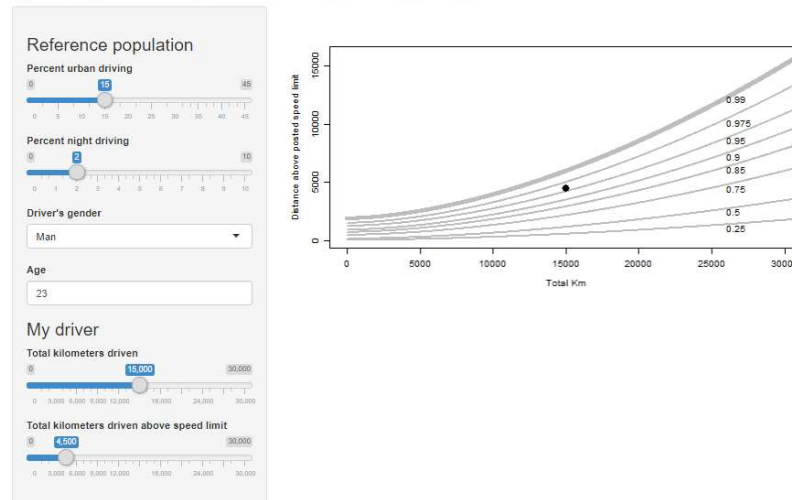
Table 6. Estimated *Tolerkm* for the 95<sup>th</sup> reference risk charts for different values of *Urban* and *Km* for men and women. Age and night driving are fixed at the sample mean level.

		<i>Km</i>		
		10,000	20,000	30,000
<i>Men</i>	Low <i>Urban</i>	3,234.56	6,572.01	11,356.31
	Median <i>Urban</i>	2,466.03	5,803.48	10,587.78
	High <i>Urban</i>	1,753.22	5,120.66	9,904.96
<i>Women</i>	Low <i>Urban</i>	3,110.06	6,447.51	11,231.81
	Median <i>Urban</i>	2,328.48	5,665.93	10,450.22
	High <i>Urban</i>	1,607.98	4,945.43	9,729.72

An interactive graphical tool that displays the evolution of a driver's speeding risk percentile as a function of total distance driven, night-time driving, gender and principal driving zone can be seen in Figure 3 and it can also be accessed online<sup>3</sup>.

Figure 3. Example of interactive speeding risk reference chart that locates a particular driver (black dot), given total distance driven, total speeding kilometres and all other reference characteristics stated in the left panel.

Speeding risk reference charts based on insurance data



<sup>3</sup> [https://riskcenter.shinyapps.io/speeding\\_risk\\_reference\\_chart/](https://riskcenter.shinyapps.io/speeding_risk_reference_chart/)

#### 4. Discussion

We have found that the most relevant variables explaining the number of kilometres driven at speeds above the limits are: total distance driven, percentage of urban driving and gender. In most of the models for these data, age and night-time driving do not have a significant impact. In both cases, this appears to be due to the lack of variability in the PHYD policies, which in our sample were sold exclusively to young drivers. We analysed the relationship between the distance driven at speeds above the limits and the total distance driven, and found this relationship not to be linear, but rather exponential. This means that as the total distance driven increases, the number of kilometres driven at speeds above the limits also increases, but at an ever-increasing rate. This might be due to the driving experience gained or to an excess of confidence on the part of the driver.

The exponential relationship introduced in the covariates of quantile regressions determines the shape of the reference risk curves for driving at excess speeds. Such models allow the factors associated with higher quantile values to be identified and, hence, for risky drivers to be detected. Our results contribute to calculating the percentile risk score for each driver by controlling for their specific characteristics (and not for the whole population of drivers). Based on these quantile regression models, risk reference curves have been obtained. These graphical tools provide, for each driver, the corresponding percentile of the distance driven at speeds above the limits (which constitutes that driver's risk score), as a function of the total distance driven. Moreover, these curves can be easily obtained for particular types of driver, depending on their characteristics (gender, percentage of urban driving, etc.).

One limitation of the analysis reported here, and which should be pointed out, is that the degree to which drivers exceeded the posted speed limit was not recorded and, therefore, we do not know the severity of the speed violation.

We consider this methodology of risk quantification to be very useful in application with Manage-How-You-Drive (MHYD) insurance products, where the premium is calculated using the same procedure as that used in PHYD insurance, but, in addition, drivers are provided with real-time alerts and recommendations for guaranteeing their safety (Arumugam and Bhargavi, 2019). As such, MHYD insurance improves both customer service and protection in the sector. In this context, the methodology presented here is able to deliver valuable graphical information in terms of preventive early warnings. Estimating just how a driver ranks with respect to distance driven above the posted speed limit is personalized information that should constitute interesting feedback for policy holders (Pérez-Marín et al., 2019b). Here, it should be stressed that excess speed is perhaps the only feature a driver can easily modify, given that other factors, such as percentage of urban driving, are largely determined by external circumstances and drivers are essentially unable to change them. Indeed, Pérez-Marín et al. (2019a) report that young drivers have a tendency to reduce speed limit violations after an accident, probably because of their greater awareness of the associated risks. As speed is the main cause of severe accidents, those who present lower risk scores (a lower percentile on the risk reference curve) should probably have lower insurance premiums. As to how this ranking should be translated into an insurance price is a question we leave for further research, but there is no doubt that direct bonuses rewarding careful drivers could easily be introduced.

In this paper, we have presented the design for a prototype graphical tool that displays the evolution of a driver's speeding risk percentile as a function of total distance driven, night-time driving, gender and principal driving zone. We show that this interface produces a personalized percentile that provides immediate feedback to the user. By measuring a driver's current speeding based on telemetry, that driver can see their evolution over distance driven and they can be provided with a score that is based on their peers' driving records. The methodology described is also applicable to different scenarios and for benchmarking drivers accordingly. We firmly believe such reference charts are set to become a standard in the visualization of driving risk.

## References

- Ayuso, M., Guillen, M. and Pérez-Marín, A.M. 2014. Time and distance to first accident and driving patterns of young drivers with pay-as-you-drive insurance. *Accident Analysis and Prevention* 73: 125–31. DOI: 10.1016/j.aap.2014.08.017.
- Ayuso, M., Guillen, M. and Pérez-Marín, A.M. 2016a. Telematics and gender discrimination: some usage-based evidence on whether men's risk of accident differs from women's. *Risks* 4:2: 10. DOI: 10.3390/risks4020010.
- Ayuso, M., Guillen, M. and Pérez-Marín, A.M. 2016b. Using GPS data to analyse the distance travelled to the first accident at fault in pay-as-you-drive insurance. *Transportation Research Part C Emerging Technologies* 68: 160–7. DOI: 10.1016/j.trc.2016.04.004.
- Arumugam, S. and Bhargavi, R. (2019). A survey on driving behavior analysis in usage based insurance using big data, *Journal of Big Data*, 6, 86, 1-21. DOI: 10.1186/s40537-019-0249-5.
- Dissanayake, S. and Lu, J.J. 2002. Factors influential in making an injury severity difference to older drivers involved in fixed object-passenger car crashes. *Accident Analysis and Prevention* 34: 5: 609–18. DOI: 10.1016/S0001-4575(01)00060-4.
- Gao, G., Meng, S. and Wüthrich, M.V. 2019. Claims frequency modeling using telematics car driving data. *Scandinavian Actuarial Journal* 2019: 2: 143-62. DOI: 10.1080/03461238.2018.1523068.
- Gao, G. and Wüthrich, M.V. 2019. Convolutional neural network classification of telematics car driving data. *Risks* 7: 1: 6. DOI: 10.3390/risks7010006.
- Guillen, M., Nielsen, J.P., Ayuso, M. and Pérez-Marín, A.M. 2019. The use of telematics devices to improve automobile insurance rates. *Risk Analysis*, 39: 3: 662-72. DOI: 10.1111/risa.13172.
- Hewson, P.J. 2008. Quantile regression provides a fuller analysis of speed data. *Accident Analysis and Prevention* 40: 502–10. DOI: 10.1016/j.aap.2007.08.007.
- Jun, J., Ogle, J. and Guensler, R. 2007. Relationships between crash involvement and temporal-spatial driving behavior activity patterns: use of data for vehicles with global positioning systems. *Transportation Research Record* 2019: 246–55. DOI: 10.3141/2019-29.
- Jun, J., Guensler, R. and Ogle, J. 2011. Differences in observed speed patterns between crash-involved and crash-not-involved drivers: application of in-vehicle monitoring technology. *Transportation Research Part C Emerging Technologies* 19: 4: 569–78. DOI: 10.1016/j.trc.2010.09.005.
- Koenker, R. and Bassett, G. (1978). Regression quantiles. *Econometrica* 46, 1, 33-50.
- Koenker, R. and Machado, J.A.F. 1999. Goodness of fit and related inference processes for quantile regression. *Journal of the American Statistical Association* 94: 448, 1296-310. DOI: 10.1080/01621459.1999.10473882.
- Koenker, R., Portnoy, S., Ng, P.T., Zeileis, A., Grosjean, P. and Ripley, B.D. 2018. Package 'quantreg'. R Package Version 5.38, <https://cran.r-project.org/web/packages/quantreg/quantreg.pdf>.
- Ossiander, E.M. and Cummings, P. 2002. Freeway speed limits and traffic fatalities in Washington State. *Accident Analysis and Prevention* 34: 13–8. DOI: 10.1016/S0001-4575(00)00098-1.
- Paefgen, J., Staake, T. and Fleisch, E. 2014. Multivariate exposure modeling of accident risk: insights from pay-as-you-drive insurance data. *Transportation Research Part A Policy and Practice* 61: 27–40. DOI: 10.1016/j.tra.2013.11.010.

Pérez-Marín, A.M., Ayuso, M. and Guillen, M. 2019a. Do young insured drivers slow down after suffering an accident? *Transportation Research Part F: Traffic Psychology and Behaviour* 62: 690-99. DOI: 10.1016/j.trf.2019.02.021.

Pérez-Marín, A.M., Guillen, M., Alcañiz, M. and Bermúdez, L. (2019b) “Quantile regression with telematics information to assess the risk of driving above the posted speed limit”, *Risks* 2019, 7(3), 80. DOI: 10.3390/risks7030080.

Pérez-Marín, A.M. and Guillen, M. 2019. Semi-autonomous vehicles: Usage-based data evidences of what could be expected from eliminating speed limit violations. *Accident Analysis and Prevention* 123: 99–106. DOI: 10.1016/j.aap.2018.11.005.

Plötz, P., Jakobsson, N. and Frances Sprei, S. (2017). On the distribution of individual daily driving distances. *Transportation Research Part B* 101, 213–227. DOI: 10.1016/j.trb.2017.04.008.

Uribe, J. and Guillen, M. (2020) *Quantile Regression for Cross-Sectional and Time Series Data: Applications in Energy Markets Using R*. SpringerBriefs in Finance. Springer. DOI: 10.1007/978-3-030-44504-1.

Wei, Y., Pere, A., Koenker, R. and He, S. (2006). Quantile regression methods for reference growth charts, *Statistics in Medicine* 25, 8, 1369-1382. DOI: 10.1002/sim.2271.

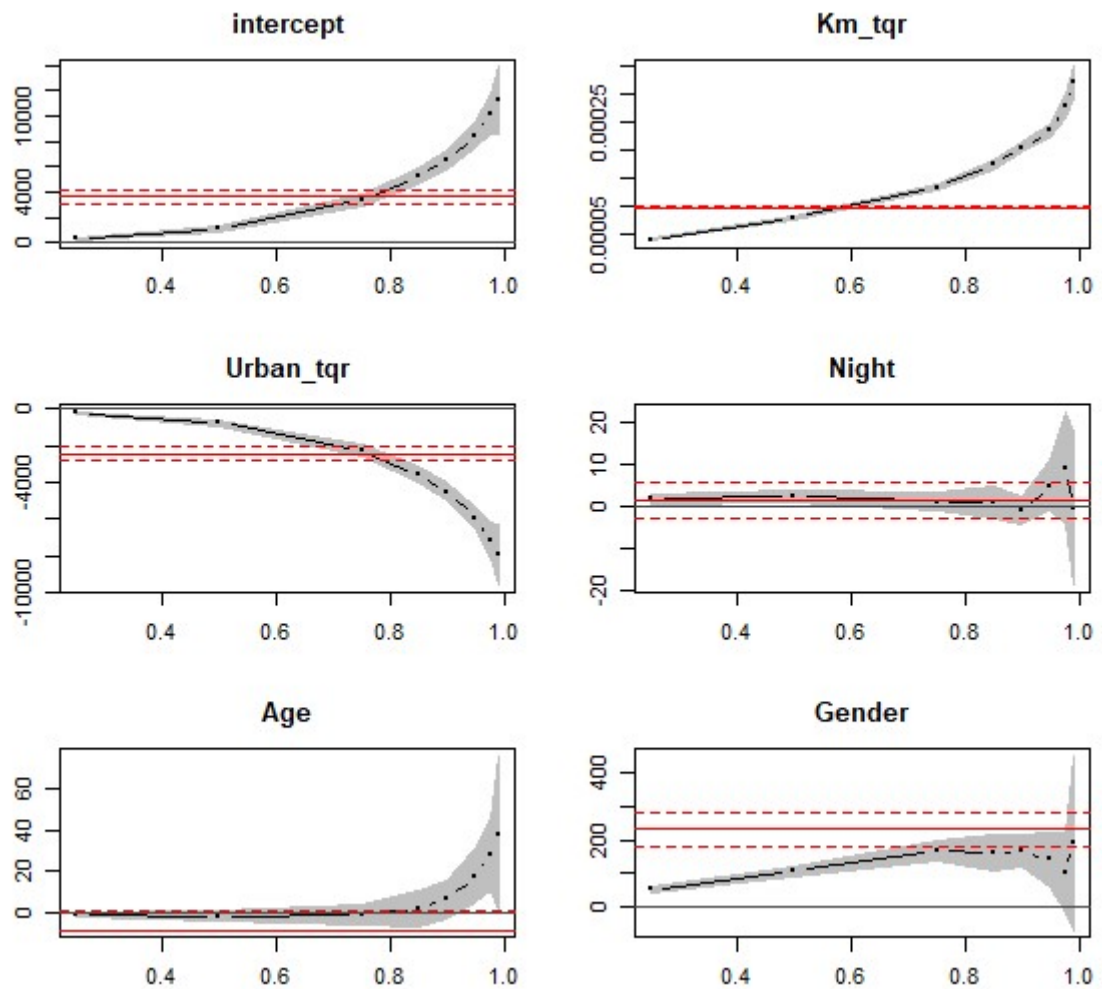
Vernon, D., Cook, L.J., Peterson, K.J., and Dean, J.M. 2004. Effect of the repeal of the national maximum speed limit law on occurrence of crashes, injury crashes, and fatal crashes on Utah highways. *Accident Analysis and Prevention* 36: 223–9. DOI: 10.1016/S0001-4575(02)00151-3.

## **Acknowledgements**

The authors thank Fundación BBVA Big Data grants and ICREA Academia.

## Appendix

Figure A1. Parameter estimates of quantile regression for total kilometres driven above the speed limit at different levels of the quantile. Confidence intervals at a 5% level of significance are shown as shaded bands. The horizontal red line represents the corresponding parameter estimate in a classical linear regression model.



The logo for UBIREA, featuring the text "UBIREA" in a bold, sans-serif font. The "UB" is in a light blue color, and "IREA" is in a darker blue. The logo is set against a white background that is part of a larger blue graphic element.

Institut de Recerca en Economia Aplicada Regional i Pública  
*Research Institute of Applied Economics*

**Universitat de Barcelona**

Av. Diagonal, 690 • 08034 Barcelona

---

**WEBSITE:** [www.ub.edu/irea/](http://www.ub.edu/irea/) • **CONTACT:** [irea@ub.edu](mailto:irea@ub.edu)

---