

Probabilistic risk modeling at the wildland urban interface: the 2003 Cedar Fire

D. R. Brillinger^{1*,†}, B. S. Autrey² and M. D. Cattaneo³

¹ *Statistics Department, University of California, Berkeley, CA 94720-3860, U.S.A.*

² *Fire Cause Analysis, Inc., 935 Pardee St., Berkeley, CA 94710, U.S.A.*

³ *Economics Department, University of Michigan, Ann Arbor, MI 48109-1220, U.S.A.*

SUMMARY

The October 2003 Cedar Fire in San Diego County was a tragedy involving 15 deaths, the burning of some 280 000 acres of land, the destruction of approximately 2227 homes, and costs of suppression near \$30 million. It was the largest fire in California history. The data associated with the fire, however, do provide an opportunity to carry out probabilistic risk modeling of a wildland-urban interface (WUI) event. WUI's exist where humans and their development interface with wildland fuel. As home building expands from urban areas to nearby forest areas, these homes become more likely to burn.

Wildfires are an exceedingly complex phenomenon with uncertainty and unpredictability abounding, hence a statistical approach to gaining insight appears useful. In this research, spatial stochastic models are developed. These relate risk probabilities and losses measures to a variety of available explanatory quantities. There is a consideration of economic aspects and a discussion of the difficulties that arose in developing the data and of carrying out the analyses. Purposes of the work include highlighting a statistical method, developing variates associated with a destruction probability and employing the fitted risk probability to estimate future and possible losses. Copyright © 2008 John Wiley & Sons, Ltd.

KEY WORDS: damage; forest fires; houses; random process; risk; simulation; wildland urban interface (WUI)

1. INTRODUCTION

This paper presents a story of deriving a dataset for statistical analysis in a complex context and employing that dataset to develop wildfire risk probabilities. The concern is the destruction of a house and to estimate corresponding losses as a function of expanatories. In preparing the dataset, a variety of sources were involved. The specific concern was the destruction of houses by wildfires at the wildland urban interface (WUI). The WUI is “the place where humans and their development meet or intermix with wildland fuel,” Federal Register (2004). Studying the fires at the WUI is important because the population of houses there is steadily growing as is the corresponding risk to lives. This research

*Correspondence to: D. R. Brillinger, Statistics Department, University of California, 367 Evans, Berkeley, CA 94720-3860, U.S.A.

†E-mail: brill@stat.berkeley.edu

concerns the particular case of the 2003 San Diego County Cedar Fire. There are various difficulties involved in working with the data available for it. One of these is that the County is made up of two jurisdictions, the City and an unincorporated part. Sometimes the data are not available in the same form for both. Yet other data sources are needed. The data sources employed are listed in the Appendix.

The paper begins with general discussion of the Cedar Fire and then turns to modeling the probability of an existing house being destroyed as a function of various explanatories. Foremost among these are location and vegetation type at the house locations. There are important variates missing, for example, the meteorology during the fire. A productive approach to risk studies often comes from asking what is an appropriate insurance premium, P , to cover the occurrence of a damaging event. Formulas that have been proposed for P include

$$(1 + \alpha)\mu_L, \mu_L + \beta\sigma_L, \mu_L + \gamma\sigma_L^2, \alpha\mu_L + \beta\sigma_L + \gamma\sigma_L^2$$

where L is the damage, μ_L and σ_L are its mean and standard deviation, and α, β, γ are weights chosen for the particular context. One sees that the expected loss $E\{L\}$, and a measure of its variability, for example, $\text{var}\{L\}$, are needed. These values will be estimated for the Cedar Fire in this work. However, because the work involves but one fire, annual premiums may not be developed from the available data. Some sort of estimate of the occurrence rate of fires is needed for that. For a discussion of premium formulas such as these, see Beard *et al.* (1984).

Difficulties arose for the destroyed properties because important explanatories, such as roofing material, were not available so the contributions of the work are limited. The explanatories employed here come from available Tax Assessor records, SanGIS (2006). These include: assessed value and size of a house and its parcel. A problem with the assessed values is that re-assessment occurs irregularly. What was done by officials during and after the fire was to estimate economic damage by assuming a cost of \$150.00 per square foot. Large amounts of data were collected, but there may be disagreements in coordinate systems, scale and accuracy, as well as missing values.

The goals of the work include: providing visual displays for insight and understanding, looking for associations with various explanatory variables, trying to understand costs, modeling destruction of property, and to provide some discussion of private and social costs in general. One question is: can one obtain reasonable estimates of the probability of a house's destruction given the available explanatories? Another is: what can be said about the statistical variability of the total estimated loss?

2. THE CEDAR FIRE

There were various large wildfires raging near San Diego in the fall of 2003. Their burn scars may be seen in Figure 1. The fires developed in the forests and were driven by winds to move and destroy homes in their path. The Cedar Fire began in the Cleveland National Forest near San Diego. Initial ignition occurred when a lost hunter set off a signal flare. It landed at the point of fire origin indicated in Figure 3. The fire lasted from 25 October to 4 November 2003. It started in the unincorporated area of the county but reached part of the City of San Diego. There were 15 deaths, 6000 firefighters involved, approximately 2227 homes destroyed, 280 000 acres burnt over, and evacuations implemented.

Wind, particularly its strengths and directions, was much involved in the development and extinguishing of the Cedar Fire, Mutch (2007). So-called Santa Anna conditions were present part of the time. The fire moved exceedingly quickly at the beginning because of the winds and the presence of dead scrub. In the end, the fire had spread out in all directions. The estimated point of ignition and final

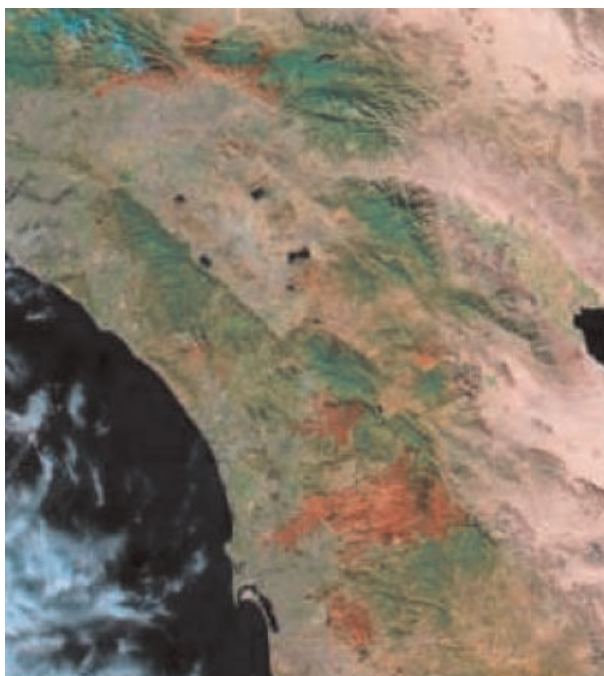


Figure 1. The burn areas for the Southern California fires as seen by satellite 5 November 2003. The Cedar Fire's perimeter, as sketched in Figures 2 and 3, may be seen in the figure from a picture in Clark *et al.* (2003).

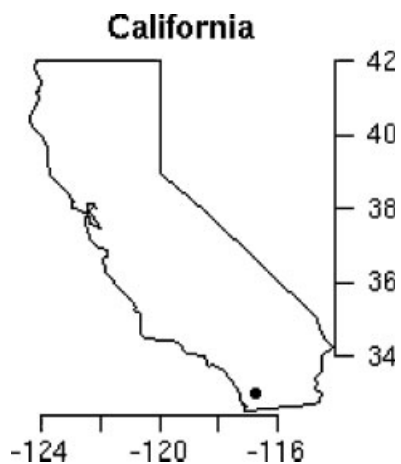


Figure 2. Map of California with the Cedar Fire point of origin indicated by the black dot.

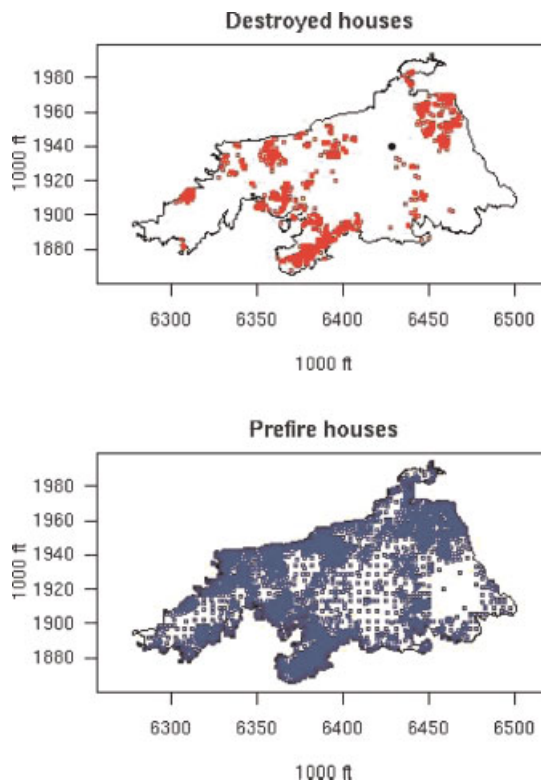


Figure 3. Top panel shows the locations of houses destroyed in red and the bottom those existing just before the fire. The black dot in the top panel is the estimated point of start of the fire.

perimeter may be seen in Figure 3. The two panels of that figure show the locations of the destroyed houses in grey/red, and of the houses existing just before the fire in black/blue. The house locations just before the fire come from the Tax Assessor records given in SanGIS (2004). The destroyed house locations come from coordinating details in an excel file provided by San Diego County with those in the Tax Assessor records. One notes clustering in both panels of the figure. The top panel also shows the estimated point of origin as a black dot. In the dataset studied, there were 2227 houses destroyed from a total of 19 560.

3. ANALYSES

3.1. *Vegetation and proportions destroyed*

A first analysis, provided in Figure 4, simply plots the proportion destroyed against the vegetation class in which the houses are situated. It also includes approximate 95% confidence intervals. The vegetation class comes from the SANDAG 1995 Vegetation Dataset (www.sandag.cog.ca.us). It is based on the predominant vegetation class in each of a set of polygons covering the burn area. The class names are provided in Figure 4. The horizontal dashed line in the figure is the overall proportion destroyed, $2227/19\,560 = 11.39\%$. The figures in brackets after the named classes in the figure are the counts of the

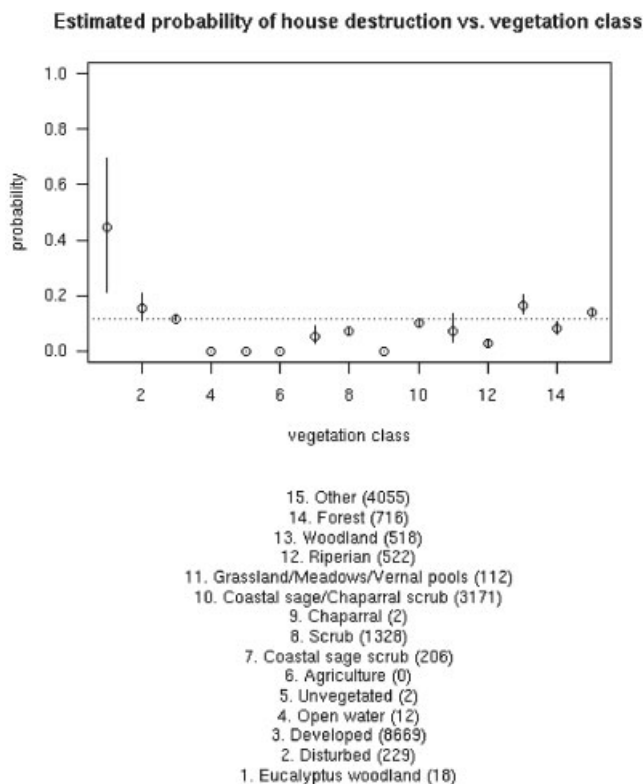


Figure 4. The top panel plots proportions destroyed with approximate 95% confidence intervals. The bottom panel lists the vegetation classes.

original houses falling in the class. One notes the claim of 12 houses in open water. This provides a warning of the errors present in the data. Incidentally, none of the open water houses were listed as destroyed in that dataset. Other classes with none destroyed are five, six, and nine. Having in mind the analyses that follow a pertinent model here is

$$\text{logit}[\text{Prob}\{\text{house destroyed}|\text{vegetation class } i\}] = \alpha_i \quad (1)$$

with i running through the vegetation classes. The proportions in the figure are then $\exp\{\hat{\alpha}\}/(1 + \exp\{\hat{\alpha}\})$. Working within the same class allows the various model fits to be compared.

Continuing with this discussion, Figure 5 shows that much of the fire area was estimated as covered by coastal sage/chaparral scrub before the fire. This may be compared with the Cedar burn scar in Figure 1. It appears as if virtually all the vegetation has been burnt by the fire. This may have resulted from several years of drought in the region.

3.2. Location and the generalized linear model

A goal of this work is to understand the probability of destruction of a house as a function of pertinent explanatories. Vegetation has just been considered. Next is location. The vegetation classes were listed

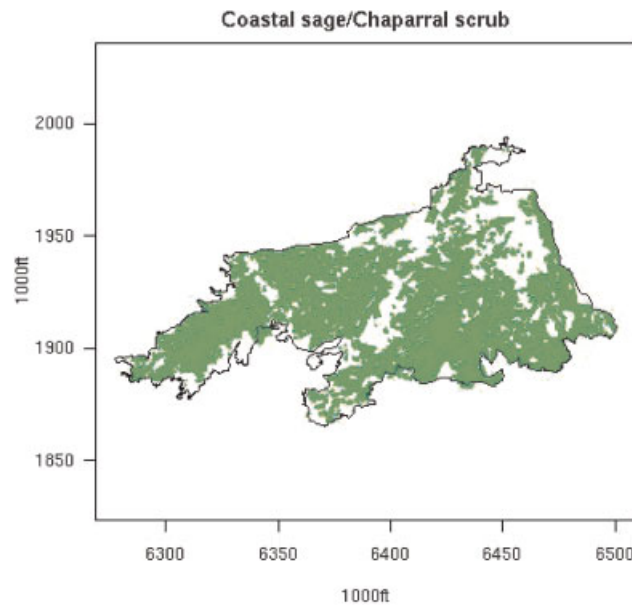


Figure 5. The region covered by coastal sage/chaparral scrub before the fire.

in Figure 4. The particular case of coastal sage/chaparral scrub is shown in Figure 5. As noted, it covers most of the fire region. In the case of location (x, y) alone, the model employed will be written,

$$\text{logit}[\text{Prob}\{\text{house destroyed}|\text{located at } (x, y)\}] = \beta(x, y) \quad (2)$$

with (x, y) location. The function β will be assumed smooth. To this end, a thinplate spline was employed, see Wahba (1990). In it, the function β is represented by

$$\beta(x, y) = \sum_{j=1}^J \gamma_j r_j^2 \log r_j$$

where for the nodes (x_j, y_j) the variable $r_j^2 = (x - x_j)^2 + (y - y_j)^2$. The nodes are taken on a lattice. As the γ_j appear linearly, the function $\text{glm}(\cdot)$ of R may be employed in the analysis.

The estimated $\beta(x, y)$ is displayed in Figure 6 as both an image and a contour plot.

The figure provides evidence for dependence of the probability of destruction on location. One sees a hot spot in dark green. The contour plot shows a highest level of 0.3. There are a number of destroyed houses clustered all around the boundary of the level 3 dark green region. The regions with lowest estimated chance of house destruction are shown in level 0 red. They make up about a third of the area. The smoothing has spread out the probability values.

There is always a need to assess the uncertainty of inferences and the validity of models generally. A direct way to do this is via a synthetic plot, Neyman *et al.* (1953), Brillinger (2008). In the construction of a synthetic plot, one simulates data from a fitted model and then compares these pseudo-observations to a display of the actual ones. In this present case, the fitted probabilities of Figure 6 were applied to thin out the original population of houses. This was done three independent times and Figure 7 provides the results. The upper left panel shows the actual houses destroyed. The other three panels are the

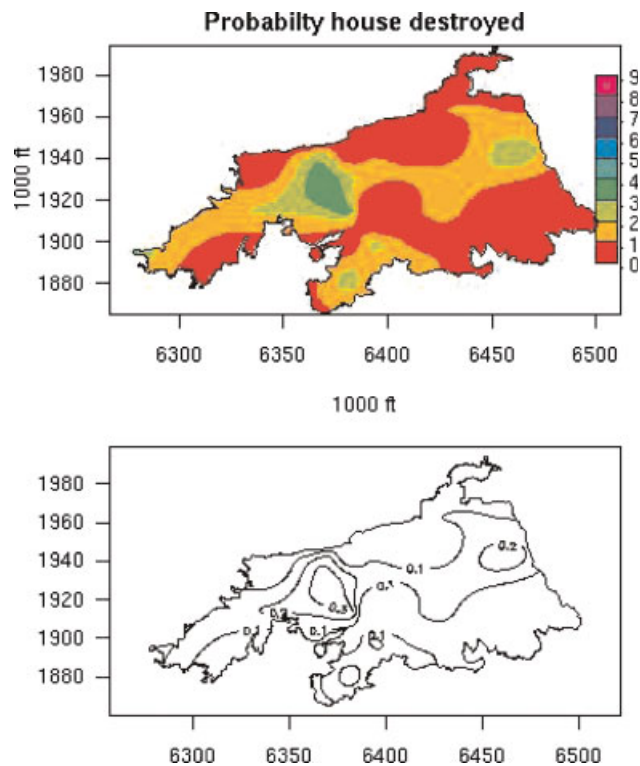


Figure 6. Estimated destruction probability as a function of location. The estimate is displayed in both perspective and contour form.

results of the independently random thinnings. The results are not unfavorable to the model. In this analysis, the only explanatory variable included in the model was location.

3.3. Other explanatories

In this section, generalized linear model analyses are carried out for the following cases of explanatories:

- location and vegetation,
- location and vegetation and size of parcel acreage (acres),
- location and vegetation and total living area (sqft),
- location and vegetation and assessor land value (\$),
- location and vegetation and assessor value of improvements (\$).

The model fits all have the form

$$\text{logit} [\text{Prob}\{\text{house destroyed}|\text{located at } (x, y), \text{ vegetation type } i, Z = z\}] = \alpha_i + \beta(x, y) + \gamma(z) \quad (3)$$

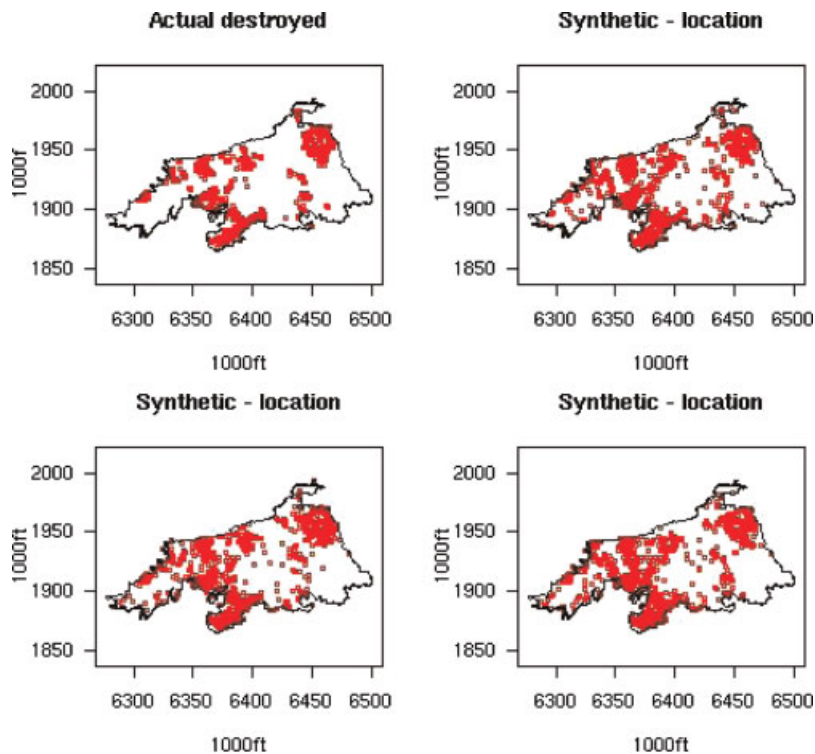


Figure 7. The locations of the original houses destroyed and three synthetic plots involving random thinnings employing the estimated probability function of Figure 6.

with i indexing vegetation type, Z a continuous and real-valued explanatory variable, and γ a smooth function. In the computations to estimate γ , Z will be taken to be a re-expression of variable values found in the assessor records. These values were taken from a SanGIS layer datasheet dated 8/18/2002, that is, about a year before the fire.

In the `glm()` computations, once again thinplate splines were employed in the estimation of β . In the case of the function γ , the quantity Z was taken to be the square root in the case of an area and the \log_{10} in the case of a dollar value. The function γ was represented as a bspline. The estimates of the respective γ s are given in Figure 8. The two area-based variates are given in the left column and the house-based ones in the right column. Approximately 95% marginal confidence intervals have been added in each case.

Examination of the plot of the top left panel provides evidence that the risk of destruction of a house is principally associated negatively with the increasing size of the parcel the house is located in. This is to be contrasted with the figure directly below which evidences an increase in risk with an increase in assessed value up to around \$100 000. The right column's figures do not suggest much dependence of the risk probability on the size of the house or its assessed value except in the case of the largest living areas. This last may represent a greater effort on the part of the firefighters to save such properties.

The dependence on the variables may also be studied via an analysis of the deviances obtained in the `glm()` fits. These are provided for the models indicated in the table below. The bracketed figures are

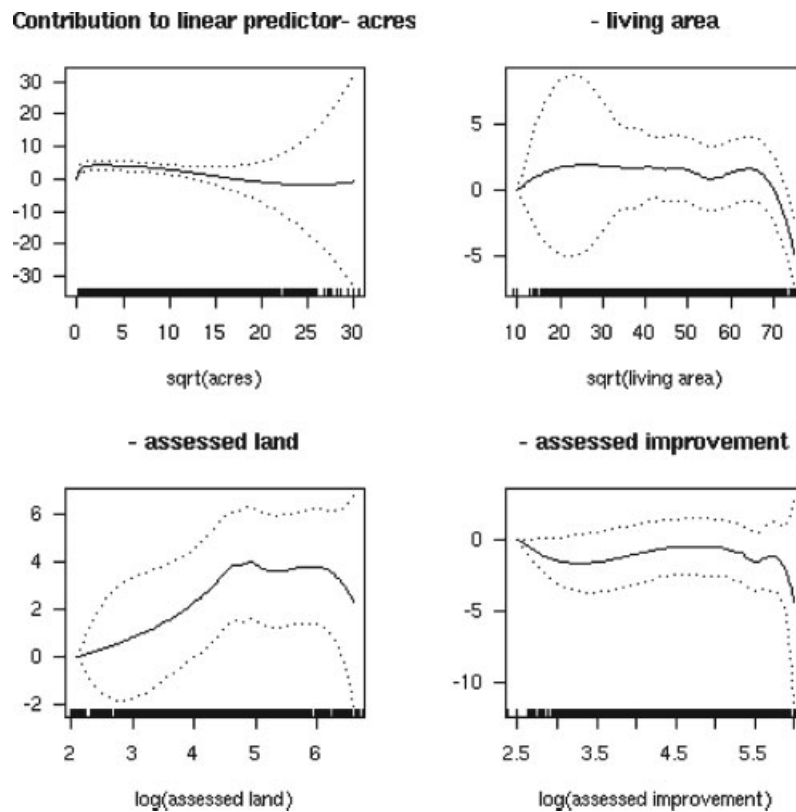


Figure 8. Estimated transforms, γ , of explanatories. The x -axis variable re-expresses the original variates by square roots and logs, respectively, in an attempt to improve the estimates.

degrees of freedom. The degrees of freedom vary because of missing values whose number changes from variate to variate.

The results for the models with explanatory vegetation, location, and both in turn are given in Table 1.

The deviance changes here are so large that it seems reasonable to include the variates vegetation and location in each of the following models and this will be done. Next the variates acres, living area, assessed land value, and assessed improvement value will be added to the model separately.

In an assessment of the remaining models, their deviance values will be considered. Working with deviances is problematic when the response is binary valued; however, we will take refuge in the

Table 1. ANODEV following model (3)

Vegetation	13 709.47 (19 546)	27 116 (19 560)
Location	12 566.96 (19 524)	27 116 (19 560)
Veg and loc	12 432.54 (19 510)	27 116 (19 560)

Second column: the final deviance.

Final column: original deviance. Bracketed numbers are degrees of freedom.

Table 2. ANODEV following model (3)

Acres	7671.4 (12 094)	7881.9 (12 104)	210.5	10	1.040e-39
Living	8102.1 (11 390)	8161.2 (11 400)	59.1	10	5.417e-09
Assessed improve	9219.9 (12 460)	9306.7 (12 470)	86.8	10	2.313e-14
Assessed land	11 515.1 (18 000)	11 985.7 (18 010)	470.6	10	8.234e-95

Second column: final deviance.

Third column: deviance with vegetation and location included. Fourth column: deviance change when indicated variable added. Fifth: degrees of freedom. Final column: χ^2 p -value.

following remarks on page 122 in McCullagh and Nelder (1989). “It is good statistical practice, however, not to rely on either D (deviance) or χ^2 (Pearson chi-squared) in these circumstances. It is much better to look for specific deviations from the model of a type that is easily understood scientifically. . . . The reduction in deviance thus induced is usually well approximated by a χ^2 distribution.”

The results obtained for the model (3) are shown in Table 2.

No matter which variate is added, the change in deviance is substantial. In terms of p -values, the assessed land value is the smallest. In future computations, all four of these variables will be added at the same time.

Figure 9 provides boxplots concerning the sizes of the houses in the dataset classified according to whether the house was destroyed or not. As evidenced by the interquartile ranges, the spread of the destroyed houses is somewhat larger than that of the non-destroyed houses. Also, the notches of the two



Figure 9. Squared feet destroyed and survived statistics displayed as notched box plots. The boxes' widths are proportional to the square roots of the respective number of data values.

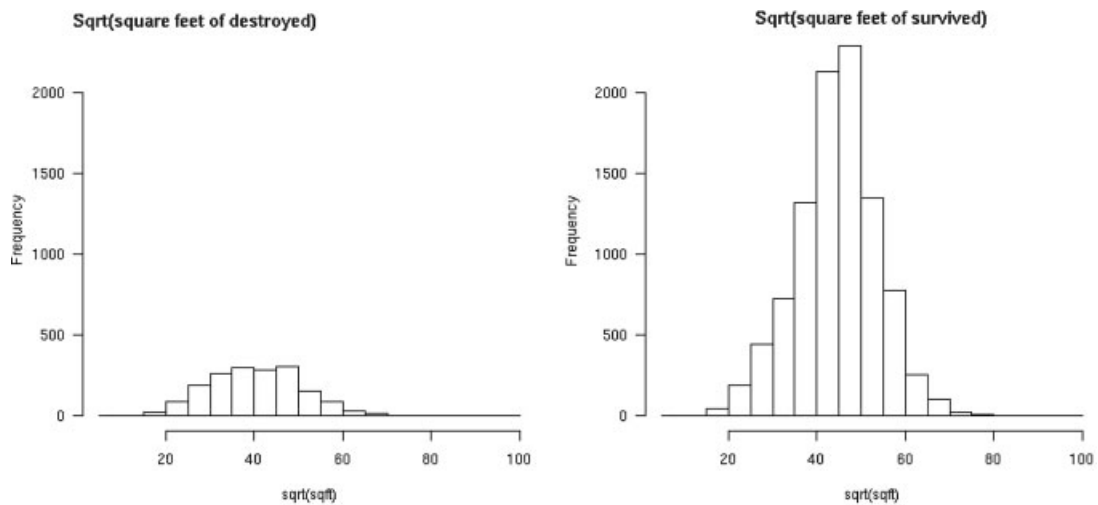


Figure 10. Histograms of the square roots of the square feet destroyed and survived.

boxes do not overlap providing evidence that the median size of the destroyed houses is smaller at the 5% level of significance.

The destroyed and survived houses may also be compared via the histograms of their square feet of living space. In preparing these, it proved more reasonable to work with the square root of the square feet. The histograms became more symmetric and took on shapes suggesting possible specific distributions for their description. The histograms appear in Figure 10.

Using the cases with available square feet in the records, there are 1757 houses destroyed and 9691 survived. The sample average and standard deviation for the destroyed houses are 1751.7 and 869.2 sqft, respectively. These are estimates of the μ_L and σ_L referred to in Section 1 as basic to the computation of premiums. One could multiply up by 150 to get 2003 dollar costs.

The fitted model may be employed to derive some interesting numbers via simulation. Suppose that one wishes to gain some appreciation of the fluctuations of the total loss, as an insurance company would. Suppose that one works with the losses in units of square feet. Considering all the houses existing before the fire, the estimated expected loss using the fitted function of Figure 4 is 3 247 666 sqft. This is obtained by summing for each house its fitted probability times its size in square feet. But now, one can generate synthetic samples and sum the sizes of those “destroyed.” An insurance company would be interested in the upper percentiles of the distribution of the difference, that is, the difference between the actual loss and the expected loss. It would assist them in setting the weights α , β , γ to avoid ruin. In a run of 10 000 “fires” the 90, 95, and 99 percentiles obtained for the loss were 99 620.32, 127 095.67, and 176 793.90 sqft.

These values can be converted to 2003 costs by multiplying with \$150. The assessor’s assessed values were available to work with, but have not been studied because similar houses can vary widely in their assessed values depending on just when the assessments were carried out.

Similar approaches to the modeling of this section were taken in Brillinger *et al.* (2003) and Preisler *et al.* (2004). Autrey (2005) is another pertinent reference.

4. DISCUSSION

Assessing the full economic impact of wildfires on society is a very challenging task. Forests provide several non-market goods and services which are especially difficult to quantify and value. Wildfires tend to have many different negative effects on society ranging from immediate physical destruction to long-term health and environmental deterioration. Unfortunately, many of these effects cannot be captured by standard methods of non-market valuation, which are often based on observed prices, indirect behavior, or self-reported valuations (for a review, see, e.g., Field and Field, 2005). Consequently, many economic price-based models that assess wildfire management and economic impact typically fail to adequately account for effects on non-marketed resources such as recreation, flora and fauna, air quality, soil, water quality, or cultural heritage. Furthermore, these models require a considerable amount of information, which in many cases is unavailable. Because of these restrictions, there still exists a substantial gap in the scientific understanding of the overall social cost associated with wildfires.

In the particular case of the Cedar Fire, lack of availability of information prevented us from using standard price-based models to assess its economic effects. Alternatively, using the available information, we are able to provide a set of basic descriptive statistics and some tentative estimates of its direct economic impact. Some specific economic values were employed in the previous section, namely assessed land and assessed improvement values. As mentioned before, these calculations underestimate its overall negative effect on society, although they do provide some interesting new results.

In terms of direct economic valuation, the 2003 Cedar Fire generated over 2200 residences destroyed alone which accounts for at least a direct estimated economic cost of \$7000 million according to the information in our database. In addition to these losses, the cost of fire suppression was estimated in \$32.5 million. Unfortunately, these figures do not account for the effects that the fire had on non-marketed resources and the corresponding long-run implications associated with them. Just to give an example in terms of vegetation, it is estimated a total loss of around 1/2 of the tree canopy population and 3/4 of both chaparral and shrub populations which, in turn, affected substantially the ecosystem services: retaining storm-water runoffs and removing air pollutants have an approximate estimated cost of \$25 million and \$1 million, respectively. These figures refer to the city of San Diego, see American Forests (2006). See also Figures 1 and 5 above, which show the near total disappearance of the green cover for the whole county.

It is important to remember that the data studied and the results have many limitations. It was noticed, for example, that some houses were classified as destroyed when they were not. There were errors in the estimated locations, both in definition and measurement. For example, there were 12 houses supposedly located in water. The limitations of deviance as a measure of fit in the case of 0-1 variables has been mentioned. There is a need to remember that these analyses are just getting at associations, for example, there may be lurking variables/proxies such as effort applied to save a building. There may be disagreements in coordinate systems, scale and accuracy, as well as missing values. Also, just one fire was studied in this research. There are other fires. There are other models. There are other known explanatories. Furthermore, there are both model and statistical uncertainties.

To properly assess the effect of the wildfire, it is important to estimate both social costs (e.g., vegetation lost or air pollution) and private costs (e.g., assets destroyed). Moreover, an appropriate analysis would account for both short-run and long-run effects. It is the case, however, that one needs really good data to proceed. This study provides but one dimension of private costs in the short run.

Figure 8 above depicted the role of different explanatory variables available to include in the model. These were plot size in acres, size of the living area in square feet, assessed land value, and assessed house improvement value. It is interesting to note that while the size of the living area and the assessed

house improvement value do not appear to be important in explaining the likelihood of the house being destroyed, the other two explanatory variables (plot size in acres and assessed land value) exhibit a considerable nonconstant, nonlinear effect on the probability of destruction. In particular, the statistical results suggest that (i) the larger the parcel size, the lower the probability of the house being destroyed, and (ii) the greater value of the land, up to a possible threshold, the greater the likelihood of destruction. These findings may have important implications in terms of risk premia and risk profiles since, for instance, houses located in larger plots should have lower risk premia. The second finding suggests an interesting risk profile namely houses located on lands relatively inexpensive (approximately less than \$100 000) have an increasing risk profile, while houses located on more valuable geographical areas have a flatter risk profile.

5. CONCLUSION

This has been each of a blue collar statistics study, a data analysis, and statistical model building. It is a work in progress, but some things have been learned. These include that: a generalized linear model provides a unified approach, it makes sense to work with square feet to deal with changing construction costs, one can estimate basic risk probabilities and limits of losses, and one can simply involve GIS files in an analysis using the *R* statistical package.

An attempt is made to assess the economic impact of the Cedar wildfire by employing a simple, often used, measure of house value namely its square feet. This was converted into dollars via a multiplier of \$150 in the official statements and reports at the time of the fire. Using the square feet, a one-dimensional quantification of the economic effect of the Cedar Fire was obtained. Further, using the results of the statistical model we also discuss the role of some observed characteristics of the houses affected, allowed us to examine the relative importance of these characteristics.

ACKNOWLEDGEMENTS

There were many people who helped in the assembly of the dataset and direction of research for this study. They include: A. Ager (USFS), J. Batchelor (SD County), J. Benoit (USFS), L. Campbell (Tierrasanta), M. Diaz (SD Foundation), F. Fujioka (USFS), D. Gilmore (SD County), P. Godden (SanGIS), A. Gonzalez-Caban (USFS), C. Hunter (Rancho Santa Fe), R. Lovett (UCB), J. Lyon (Poway), D. Martell (U of Toronto), R. Martin (SD County), H. Preisler (USFS), M. Rose (Tierrasanta), D. Sapsis (State of California), P. Spector (UCB), M-H. Tsou (SDSU), T. Westerling (UCM), C. Westling (SD County), K. Wright (USFS). This work was supported by grants NSF DMS-0504162 and DMS-0707157.

REFERENCES

- American Forests. 2006. *San Diego Urban Ecosystem Analysis After the Cedar Fire*. American Forest: Washington, DC.
- Autrey BS. 2005. Institutions, information, and catastrophes. *Honors Thesis*, UC Berkeley.
- Beard RE, Pentikainen T, Pesonen E. 1984. *Risk Theory*, 3rd Edition. Chapman & Hall London.
- Brillinger DR. 2008. The 2005 Neyman lecture: Dynamic indeterminism in science. *Statistical Science*. Neyman.
- Brillinger DR, Preisler HK, Benoit JW. 2003. Risk assessment: a forest fire example. *Science and Statistics*. In *Lecture Notes in Statistics 40, IMS*; pp 177–196.
- Clark J, Parsons A, Zajkowski T, Lannon K. 2003. Remote sensing imagery support for burned area emergency response teams on 2003 Southern California wildfires. *Report RSAC-2003-RPT1*, United States Department of Agriculture Forest Service—Engineering.
- Federal Register, Vol. 66, No. 3, January 2001, Notices, p. 753.

- Field BC, Field MK. 2005. *Environmental Economics*. McGraw-Hill/Irwin, New York.
- McCullagh P, Nelder JA. 1989. *Generalized Linear Models*, 2nd edn. Chapman and Hall, London.
- Mutch RW. 2007. Faces: the stories of the victims of Southern California's 2003 fire siege. Wildland Fire Lessons Learned Center. www.WildFiresLessons.net
- Neyman J, Scott EL, Shane P. 1953. On the spatial distribution of galaxies a specific model. *Astrophysical Journal* **117**: 92–133.
- Preisler HK, Brillinger DR, Burgan RE, Benoit JW. 2004. Probability based models for estimation of wildfire risk. *International Journal of Wildland Fire* **13**: 133–142.
- SanGIS. (2006). San Diego Firestorm 2003. Compact disk. San Diego, CA. www.sangis.org
- Wahba G. 1990. Spline models for observational data. *CBMS-NSF Conference Series in Applied Mathematics*.

6. APPENDIX

The sources of the data will now be recorded.

Much of the data came from the SanGIS CD, “San Diego Firestorm 2003.” That CD included the Tax Assessor Records for 8/12/2002 for the entire county, that is, before the 2003 Cedar Fire. Which in turn included the Assessor Parcel Number (APN) and the locations of the houses.

The data on the destroyed houses in the unincorporated part of the county came from the County of San Diego office. For the houses destroyed in the City, specifically Tierrasanta, Scripps, and Poway, use was made of the web and telephoning. The destroyed houses could then be matched up with the Tax Assessor's records to obtain additional details of use in the analyses.