



Towards Better Crash Frequency Modeling: Fusing Machine Learning & Econometric Methods

Presenter:

Behram Wali

Ph.D. Student

TSITE 2017 Summer Meeting

Morning Session

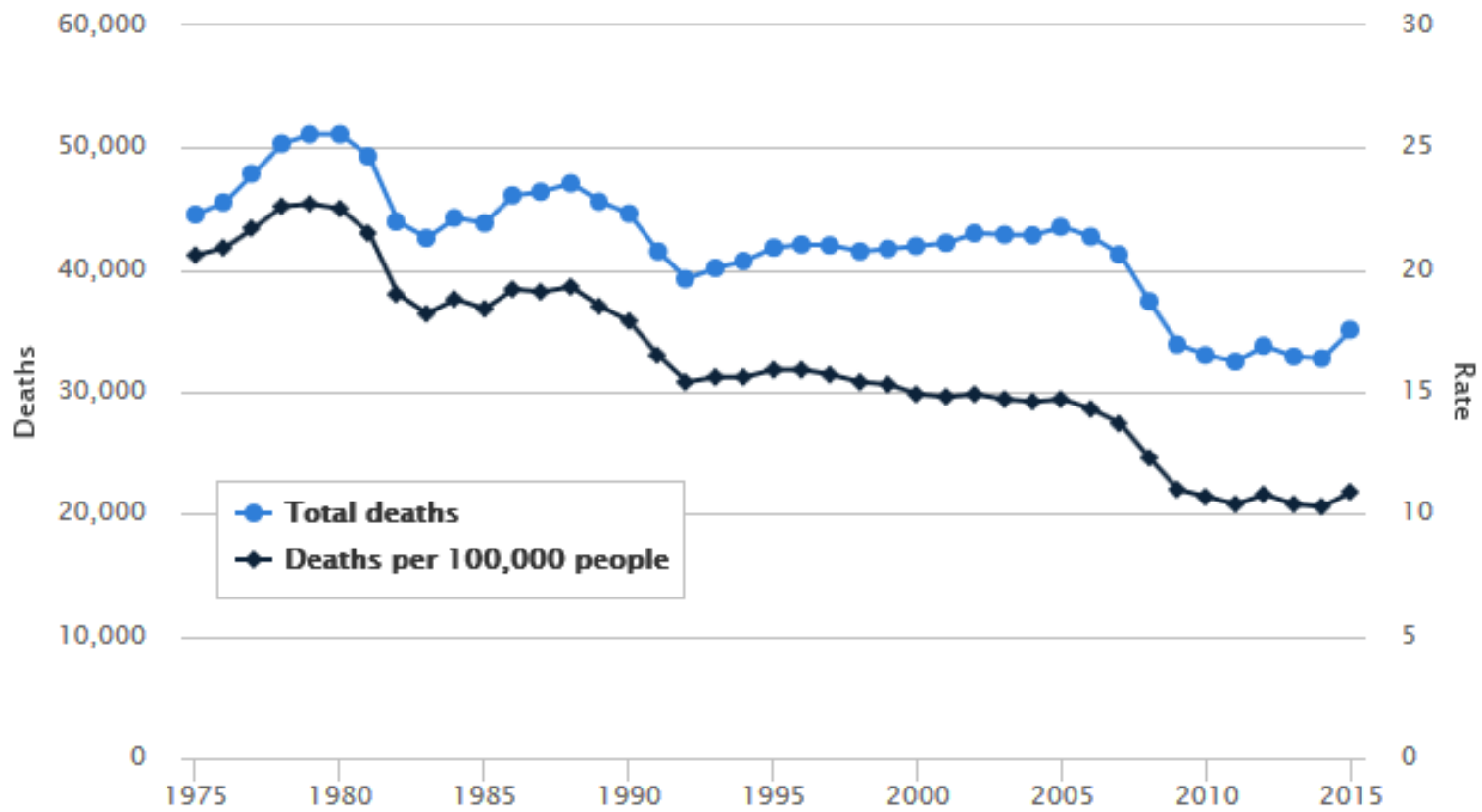
July 26, 2017



Contents

- Background/Challenges
- Conceptual Framework
- Crash Modeling: Methodological Frontiers
- State-of-the-art → State-of-the-practice
- Context: TN Rural TWTL Roads
- Take-Aways

Background



Source: [IIHS](#)



Background

Fortune

TRAFFIC ACCIDENTS

2016 Was the Deadliest Year on American Roads in Nearly a Decade

Kirsten Korosec
Feb 15, 2017



Lower gas prices and increased motor-vehicle mileage combined with risky activities like speeding and driving while texting is proving deadly for American drivers.

The New York Times

BUSINESS DAY

Biggest Spike in Traffic Deaths in 50 Years? Blame Apps

By NEAL E. BOUDETTE NOV. 15, 2016

The Washington Post

Transportation

Traffic deaths soared past 40,000 last year for the first time in a decade

Background

- Safety:

40,000/year

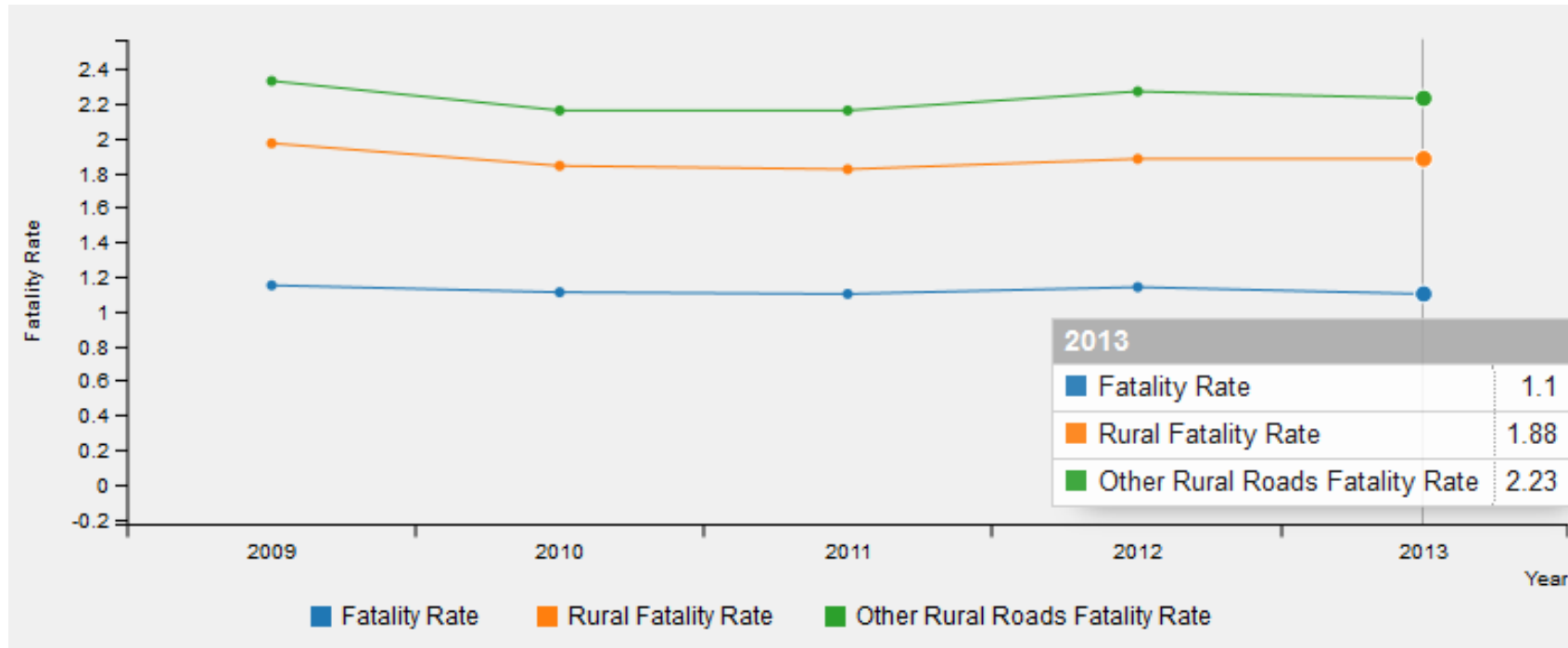
X

\$9.1 M/human life

\$364 billion/year

Serious Challenges

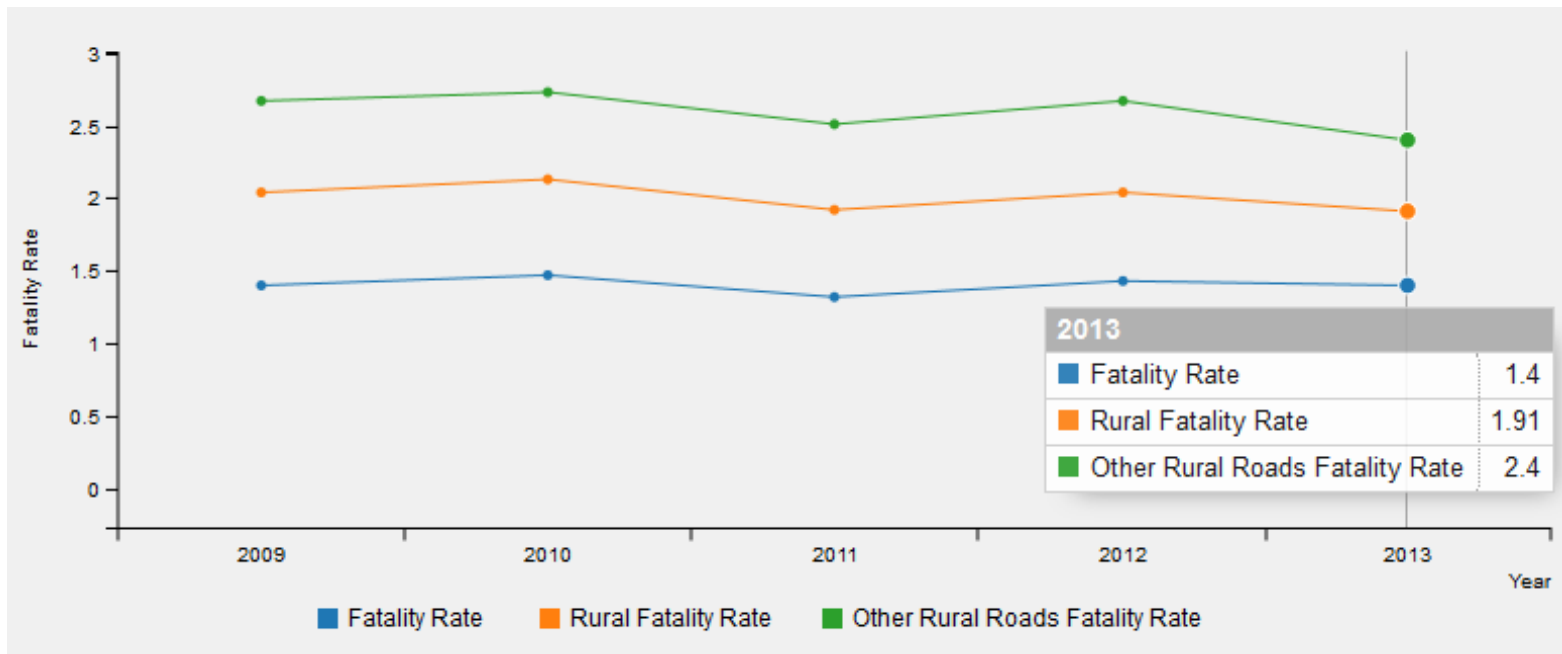
- Nationwide Fatality Rates



Source: fhwa.dot.gov

Serious Challenges

- Tennessee Fatality Rates:



Source: fhwa.dot.gov

Themes & trends: Emerging Hot Topics

Key Focus: Driver
& Technology

Driver behavior



(Sun & Yin, 2017)

Themes & trends: Emerging Hot Topics

Key Focus: Driver & Technology

Driver behavior



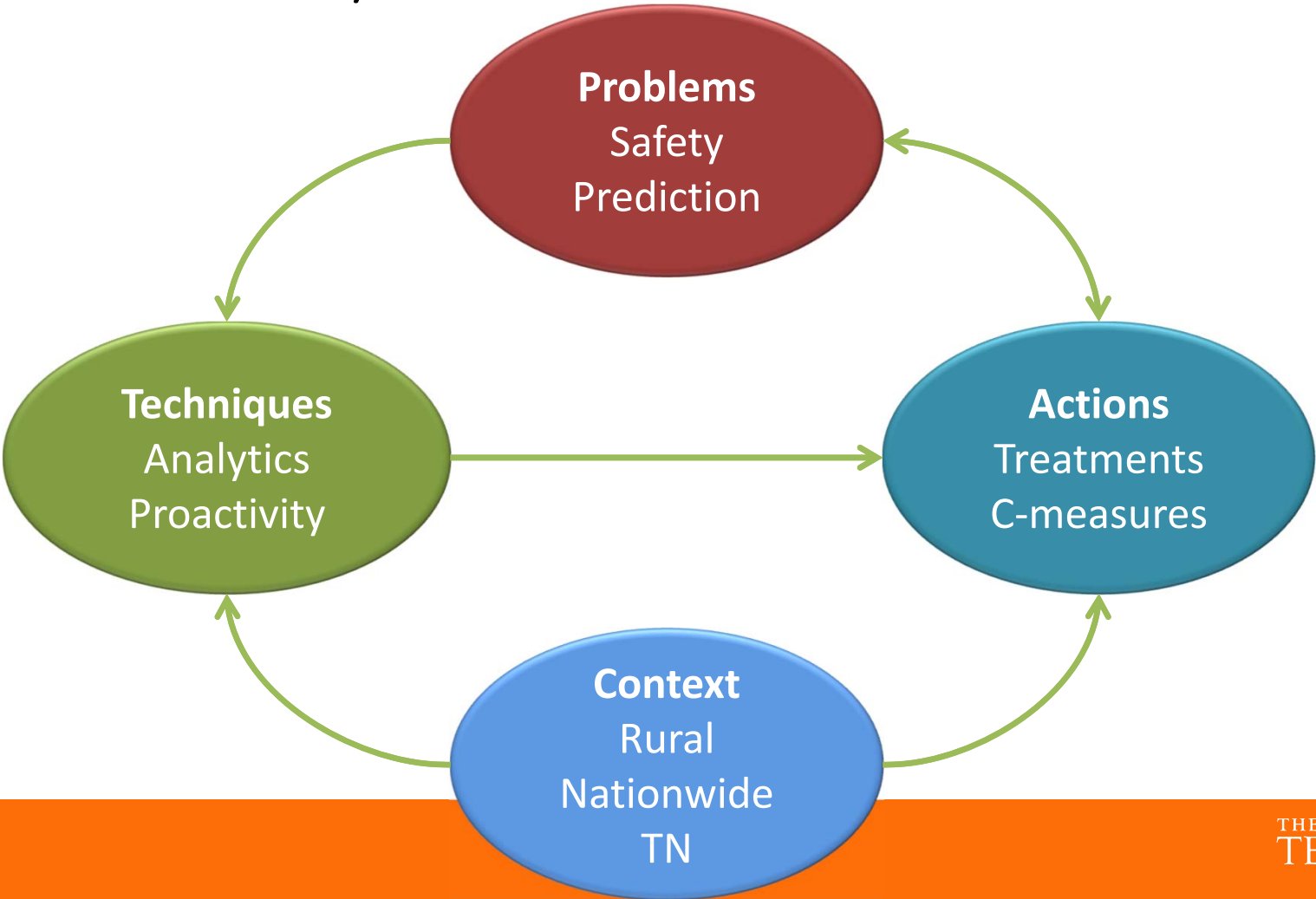
Key Targets: Safety

Safety

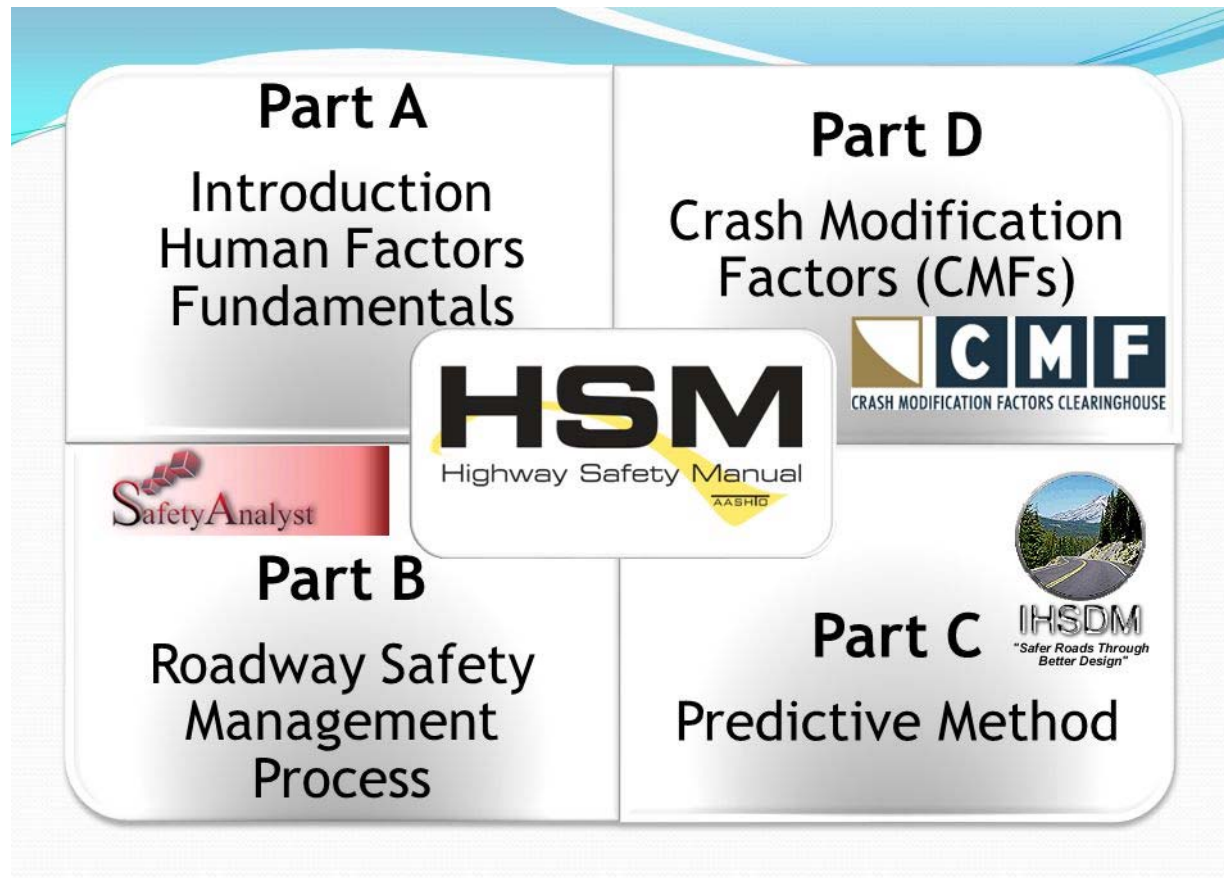


(Sun & Yin, 2017)

Framework – Learn
from success & failures/mistakes



Crash Frequency Models



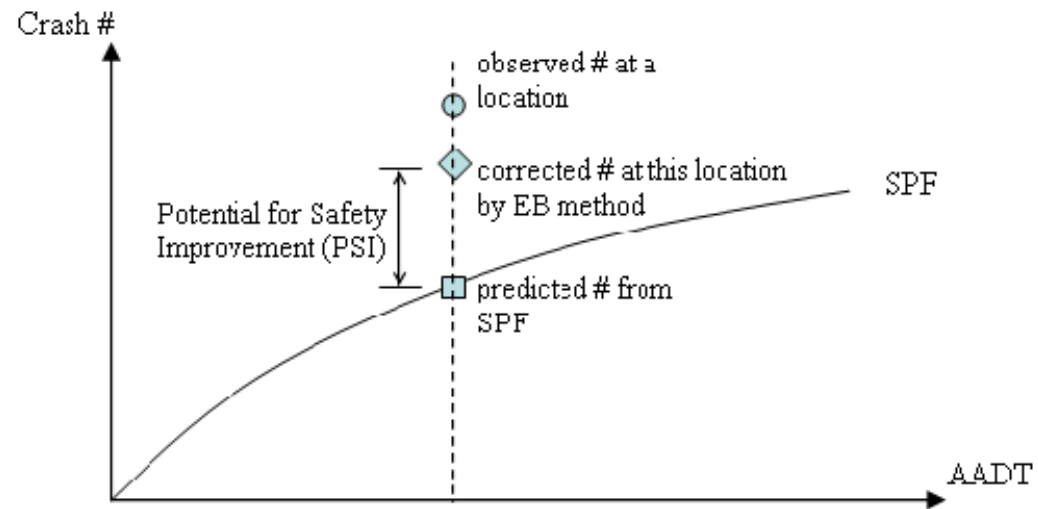
Source: [HSM](#)

Safety Performance Functions

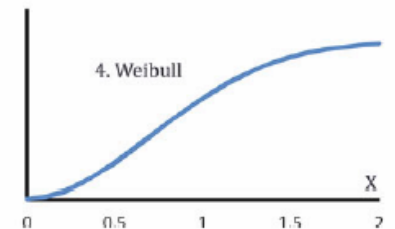
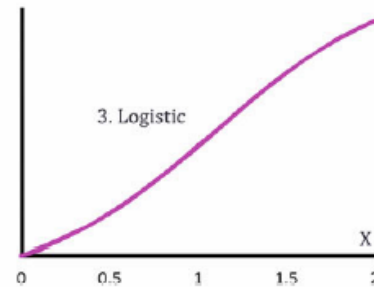
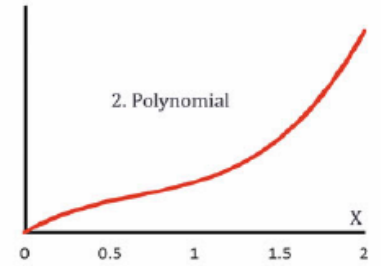
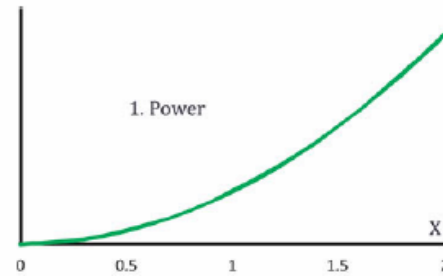
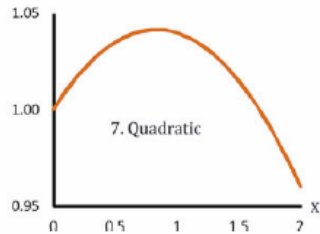
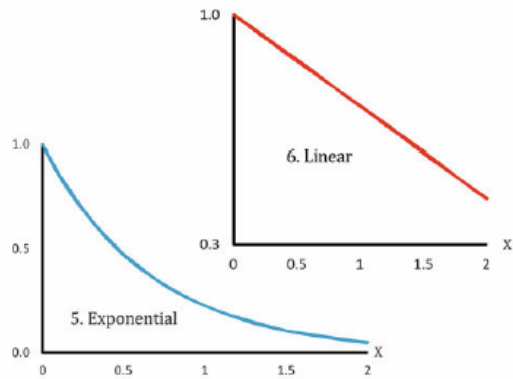
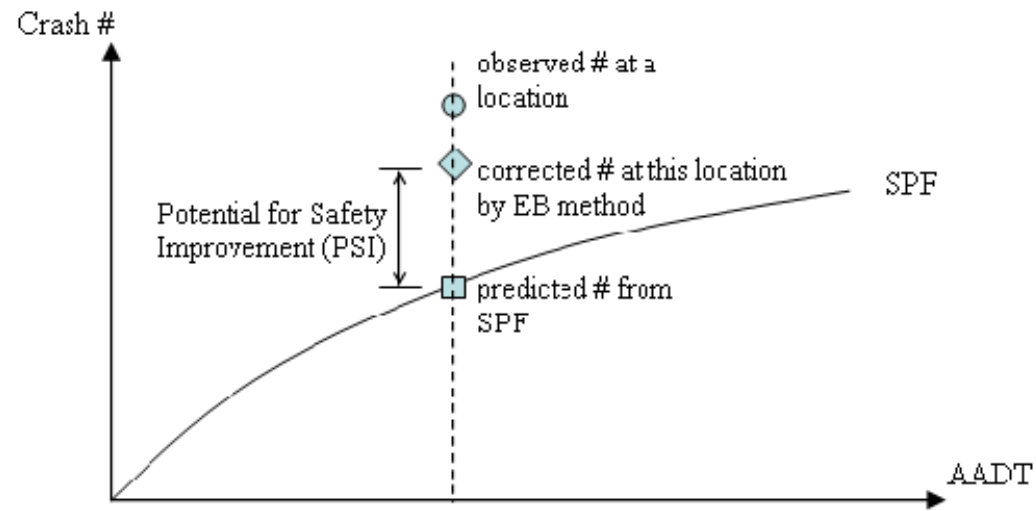
- $N_{SPF}(rural\ 2W2L) = AADT * L * 365 * 10^{-6} * e^{-0.312}$
- **Calibration done for:**
 - Base case conditions (AADT & SL only), assuming all other CMFs equal 1
 - Adjusting HSM base condition (with AADT & SL) predictions with appropriate CMFs

Source: [HSM](#)

Methodological Issues



Methodological Issues

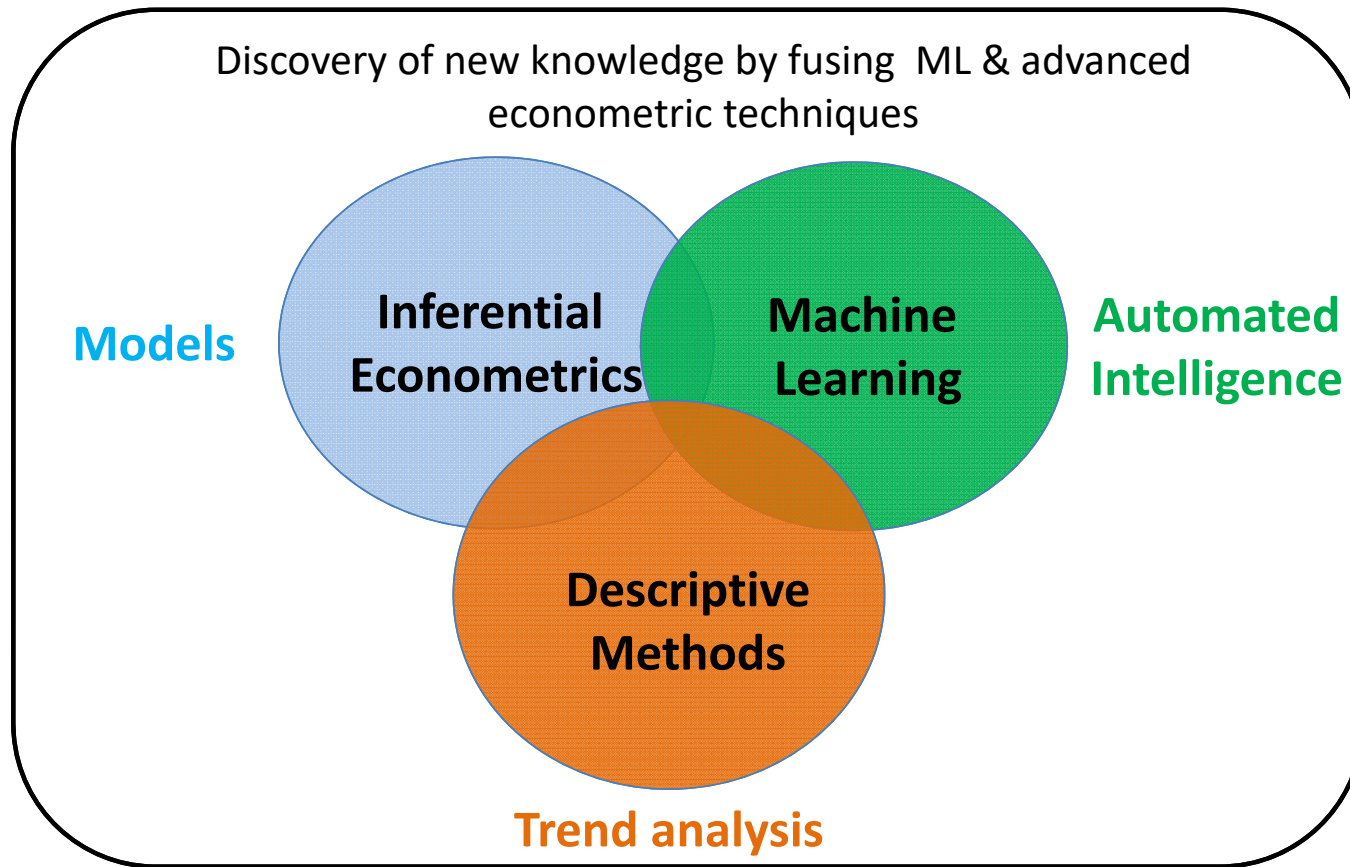


Key Issue: *How to correctly capture the complex non-linear dependencies in SPF development?*

Goal: *To enhance real-world crash prediction accuracy*

Key Challenge: *Connect advanced empirical methods to state-of-the-practice*

Methodological Frontier



Data Assembly

- ETRIMS
- **Crash data for segments**
- Rural 2W2L (seg length \geq 0.10 miles) <https://e-trims.tdot.tn.gov>
- N = 14, 777 roadway segments (total 22,000+)
- Random sample: 336 homogenous roadway segments
- Five years (2011-2015) **crash summary reports** (total and by crash severity)

Data Assembly

- **ETRIMS Exposure Data**
 - AADT for 2015 & segment length extracted
 - **Linked 2011-2014 AADT with 336 segments**
- <https://www.tdot.tn.gov/APPLICATIONS/traffichistory>

Traffic History reflects the Annual Average Daily Traffic (AADT) count along specific locations on Tennessee's road network

View stations on map: Non-Map Record Search: Station Number:

Station Information

To view Traffic History using the map:

1. Either zoom into the map to find the location or pick the "County" then zoom in

To view Traffic History without using the map:

1. Select a "County" from the list
2. Type in a "Station Number"
3. Pick "Search"

Download File: [KML](#) | [ESRI Geodatabase](#) | [ESRI Shapefile](#) | [Database Table](#)
Open With: [Google Earth](#) | [ArcGIS Explorer](#) | [MS Access or Excel](#)

Data Assembly

- **ETRIMS**-Inventory Image Viewer Web Applications
- **Detailed geometric data** manually extracted and coded
- **Data elements:**

<p>Response Variables:</p> <ol style="list-style-type: none"> 1. Total crashes <p>Key Correlates:</p> <ol style="list-style-type: none"> 1. AADT 2. Segment length 	<p>Additional Correlates:</p> <ol style="list-style-type: none"> 1. Lane width 2. Shoulder type 3. Combined shoulder width 4. Presence or absence of centerline rumble strips 5. Presence or absence of passing lane 6. Presence or absence of short four lane section 7. Presence or absence of two way left turn lane 8. Presence or absence of roadway lighting
---	---

SLD Chart

AADT	5888	2080
Roadnames	STATE HWY. 6	
Speed Limit	55	
Pavement Width	20	LM: 24.688 ²
No. Lanes	2	

Roadway Composition

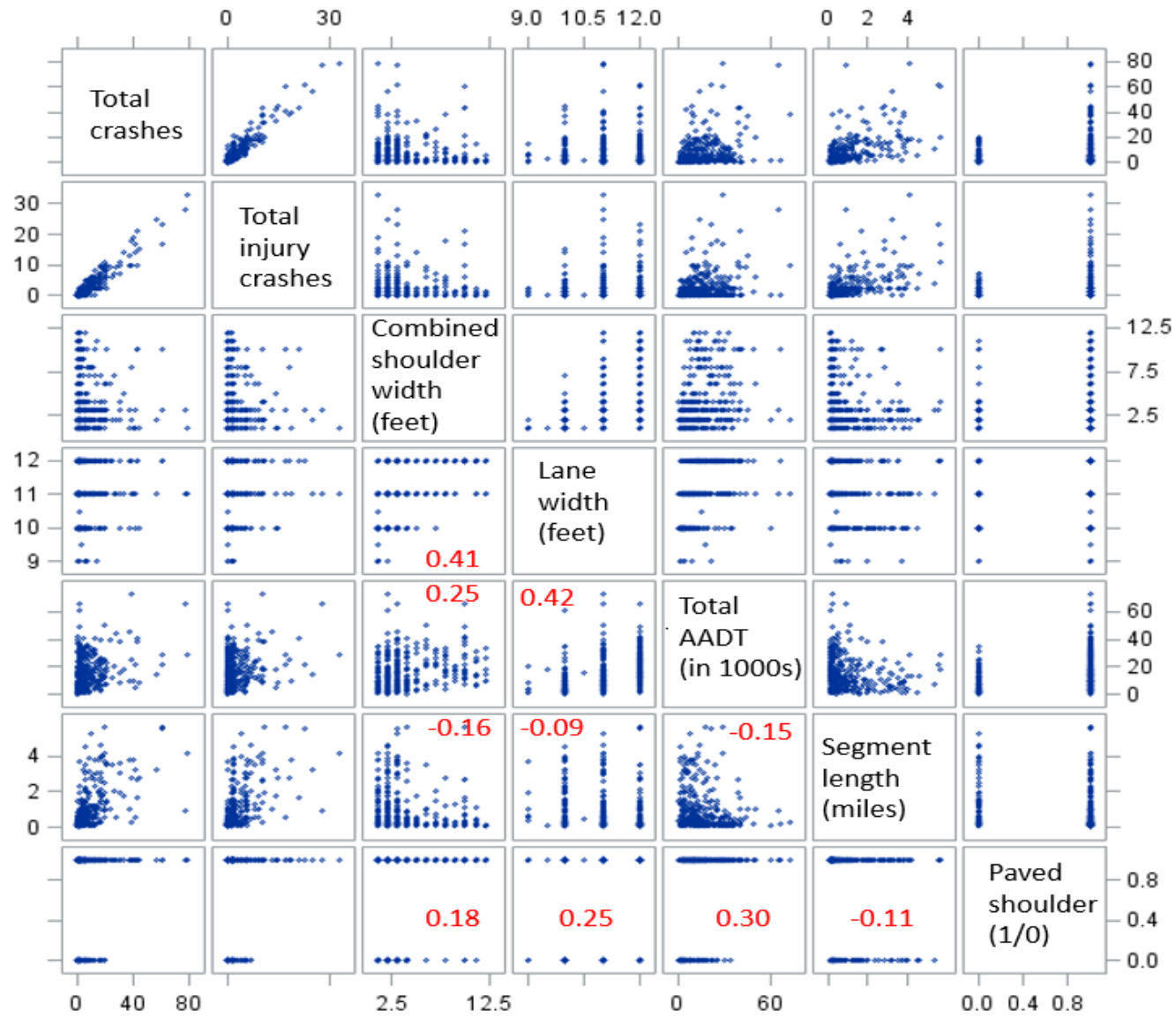
From: 24.490 To: 24.800

3'	22'	5'
Asphalt Concrete	Asphalt Concrete	Asphalt Concrete
Shoulder (Outside)	Pavement	Shoulder (Outside)

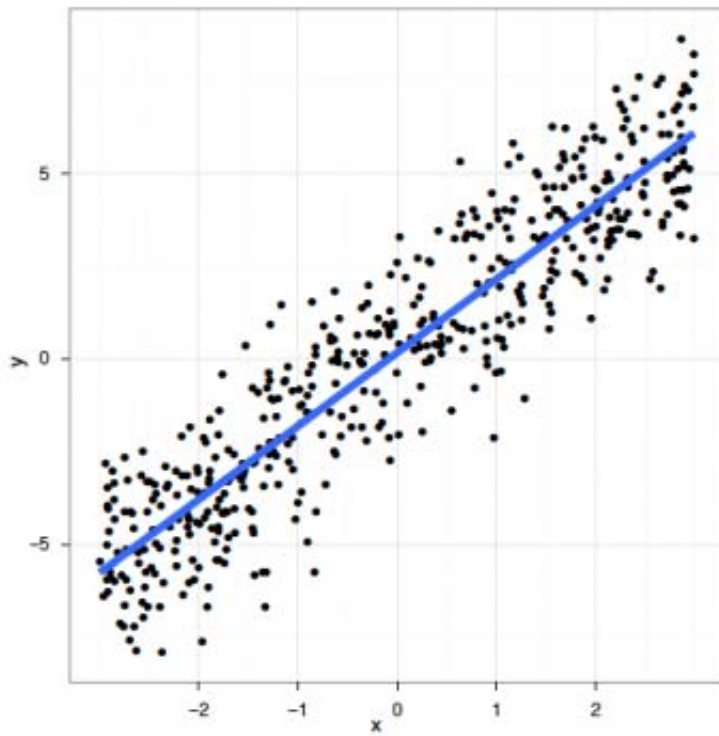
Descriptive Statistics

	Variable	N	Mean	SD	Min	Max
Key variables	Total crashes (5 years)	336	7.7	11.4	0.0	79.0
	Total injury crashes (5 years)	336	2.6	4.4	0.0	33.0
	Average AADT/Year	336	3101	2451	74	14610
	Total AADT (5 years)	336	15505	12256	368	73051
	Total AADT (5 years) in 1000s	336	15.0	12.3	0.4	73.1
	Segment length	336	0.93	1.14	0.10	5.66
Additional variables	Presence of passing lane	336	0.39	0.49	0	1
	Lane width	336	11.04	0.83	9	12
	Combined shoulder width	336	3.90	3.00	1	12
	Gravel	336	0.07	0.26	0	1
	Paved	336	0.76	0.42	0	1
	Turf	336	0.16	0.37	0	1
	Lighting	336	0.26	0.44	0	1
	Speed Limit	336	46	9	20	55

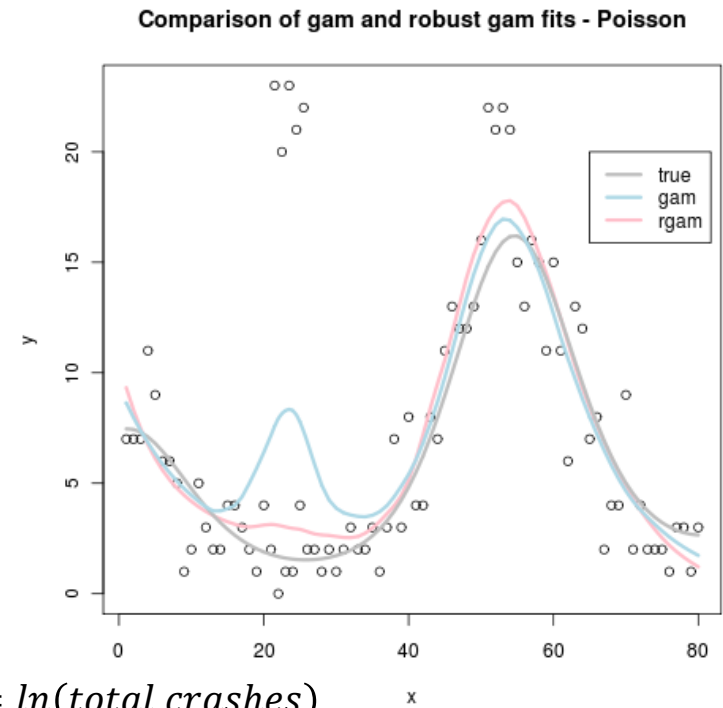
Matrix Plot



Applied Generalized Additive Models



$$\gamma_i = \ln(\mu_i) = \beta_0 + \sum_{i=1}^n \beta_j x_{ij}$$



$$\begin{aligned} \gamma_i &= \ln(\text{total crashes}) \\ &= \beta_0 + \sum_{i=1}^n \beta_i(x_{ij}) + f_i(\text{AADT}) + f_i(\text{Segment length}) \end{aligned}$$

Selected Results: Category 1 NBGAMs

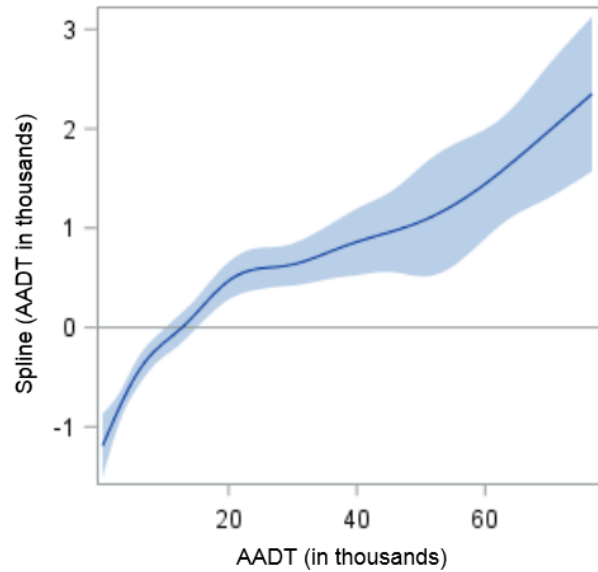
	Category 1 NBGAM		
Variables	Parameter estimate	t-statistic/F-statistic	p-value
Models for total crashes			
Intercept	1.53	38.25	< 0.0001
Spline (AADT)	DF = 6.63	F-value = 191.32	< 0.0001
Spline (Segment length)	DF = 5.52	F-value = 432.15	< 0.0001
Paved shoulder	---	---	
Combined Shoulder Width	---	---	
Lane width	---	---	
Dispersion parameter	0.35	1.41	---
Model for injury crashes			
Intercept	0.39	6.5	< 0.0001
Spline (AADT)	DF = 4.93	F-value = 124.17	< 0.0001
Spline (Segment length)	DF = 5.40	F-value = 300.29	< 0.0001
Paved shoulder	---	---	
Combined Shoulder Width	---	---	
Lane width	---	---	
Dispersion parameter	0.36	1.31	---

Selected Results: Category 1 NBGAMs

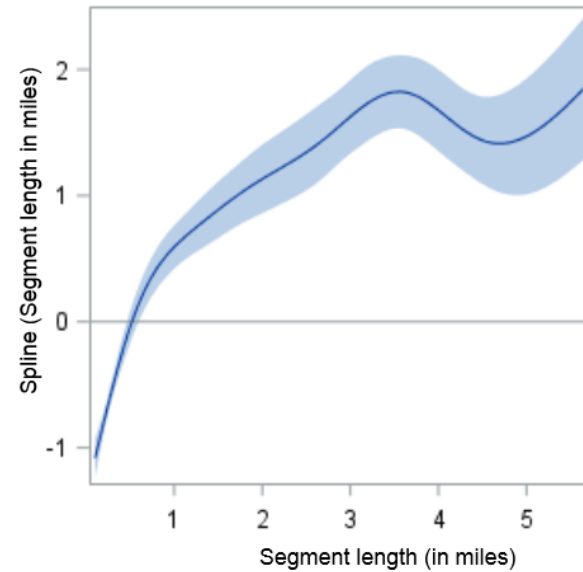
Smoothing Components for Total Crashes (Category 1 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 6.63



DF = 5.52

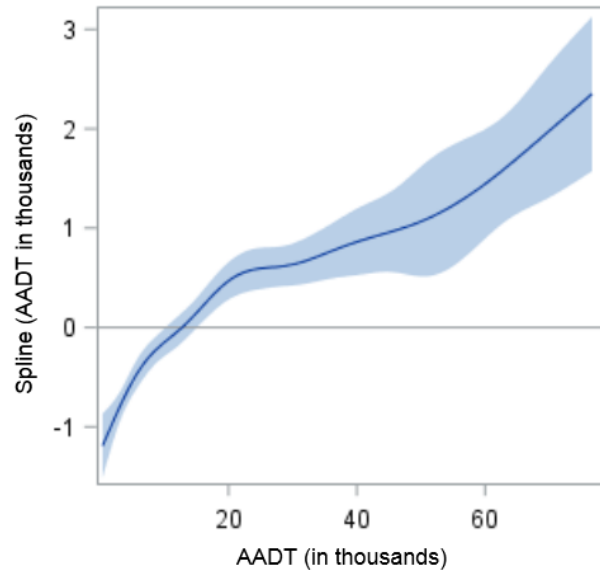


Selected Results: Category 1 NBGAMs

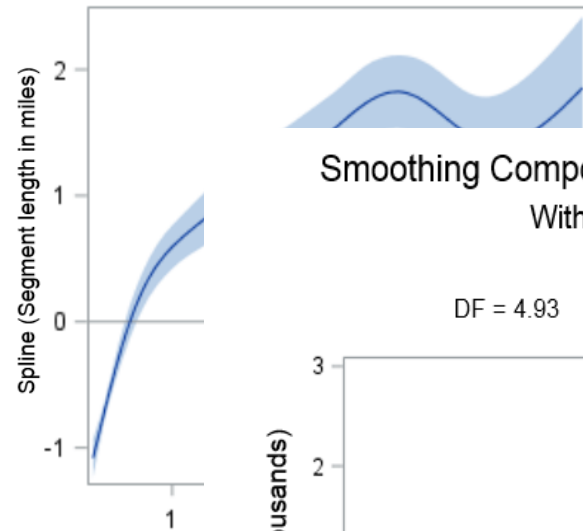
Smoothing Components for Total Crashes (Category 1 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 6.63



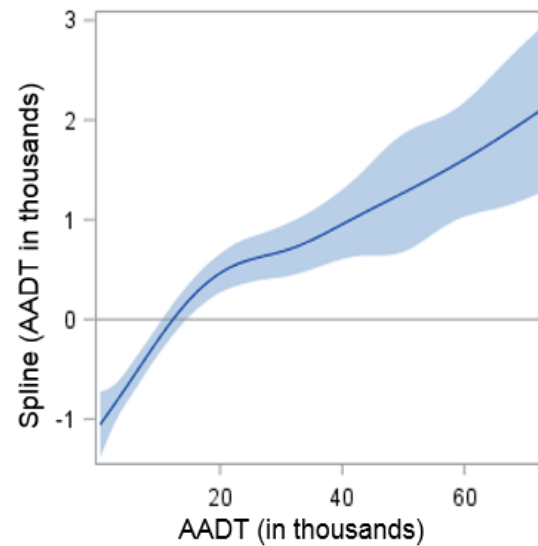
DF = 5.52



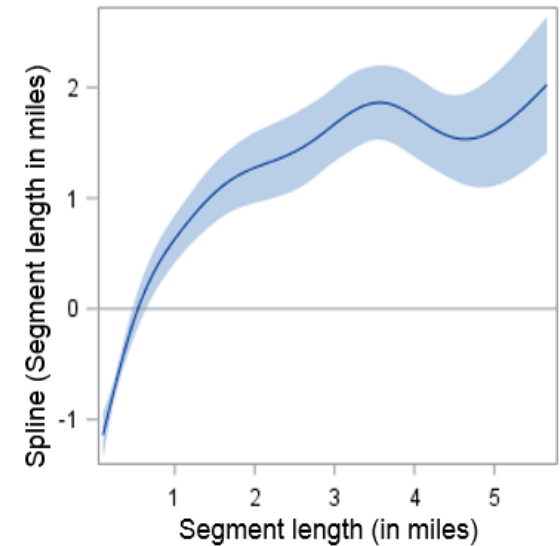
Smoothing Components for Total Injury Crashes (Category 1 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 4.93



DF = 5.40



Selected Results: Category 2 NBGAMs

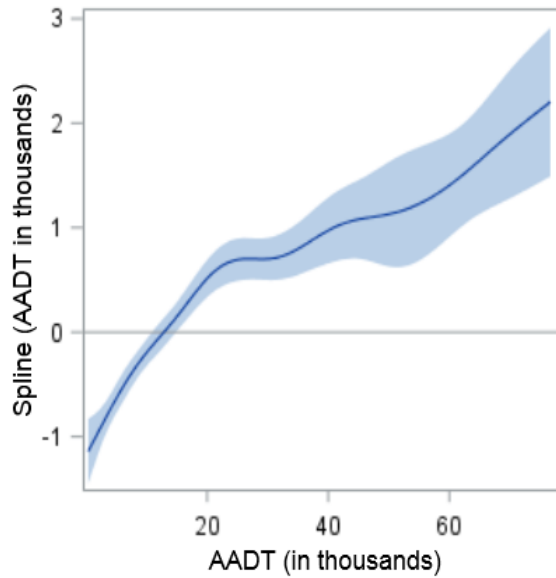
	Category 2 NBGAM		
Variables	Parameter estimate	t-statistic/F-statistic	p-value
Models for total crashes			
Intercept	2.74	4.08	< 0.0001
Spline (AADT)	DF = 6.33	F-value = 167.52	< 0.0001
Spline (Segment length)	DF = 5.04	F-value = 447.08	< 0.0001
Paved shoulder	0.41	3.72	0.0003
Combined Shoulder Width	-0.05	-5.02	0.0067
Lane width	-0.12	-2.03	0.0152
Dispersion parameter	0.3	0.97	---
Model for injury crashes			
Intercept	0.86	0.81	0.3016
Spline (AADT)	DF = 4.55	F-value = 103.07	< 0.0001
Spline (Segment length)	DF = 5.44	F-value = 312.66	< 0.0001
Paved shoulder	0.41	2.85	0.0096
Combined Shoulder Width	-0.07	-3.51	0.0018
Lane width	-0.01	-0.91	0.5353
Dispersion parameter	0.29	1.19	---

Selected Results: Category 2 NBGAMs

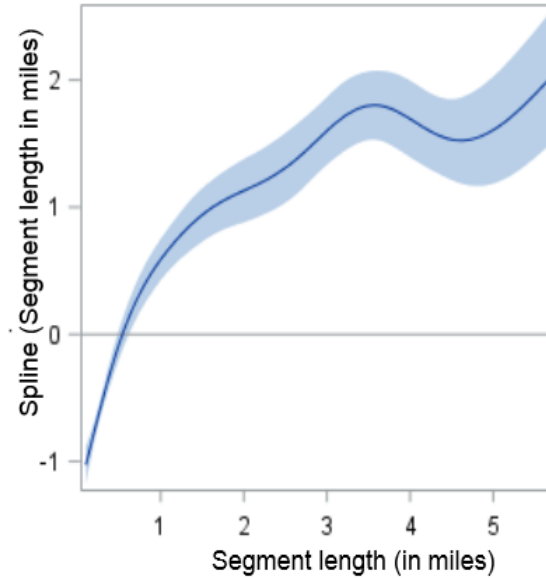
Smoothing Components for Total Crashes (Category 2 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 6.33



DF = 5.04

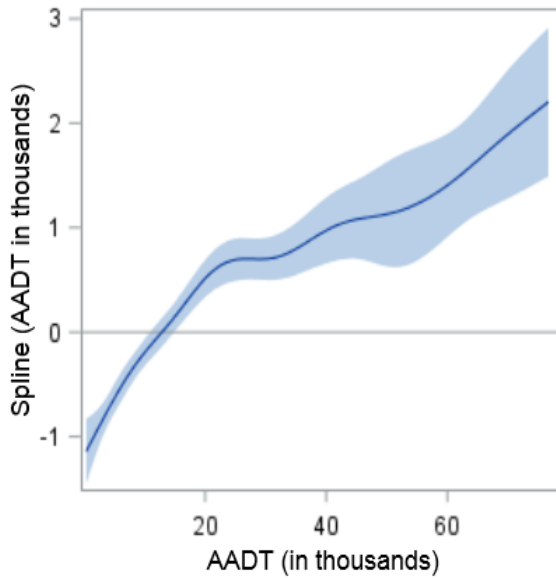


Selected Results: Category 2 NBGAMs

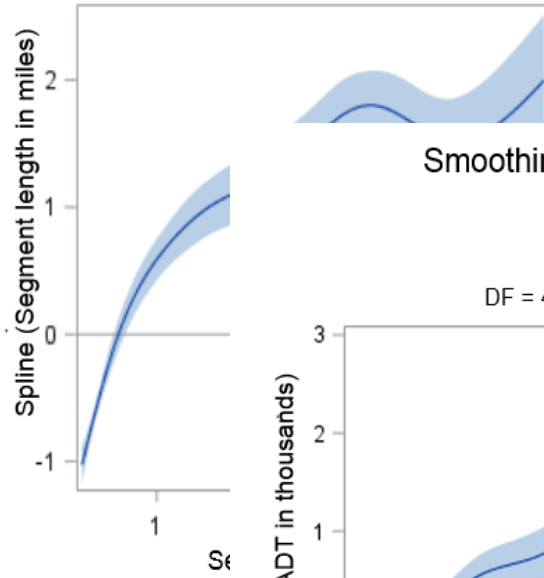
Smoothing Components for Total Crashes (Category 2 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 6.33



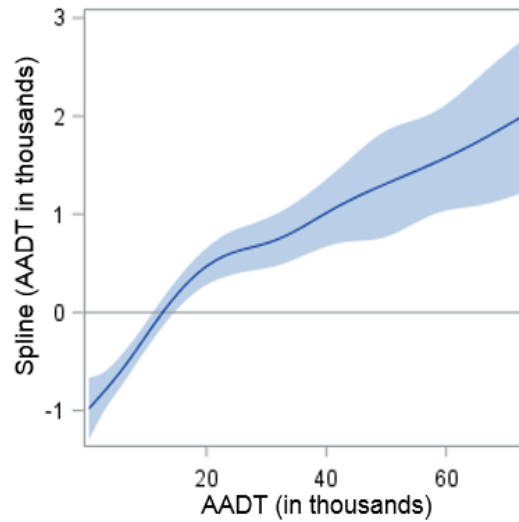
DF = 5.04



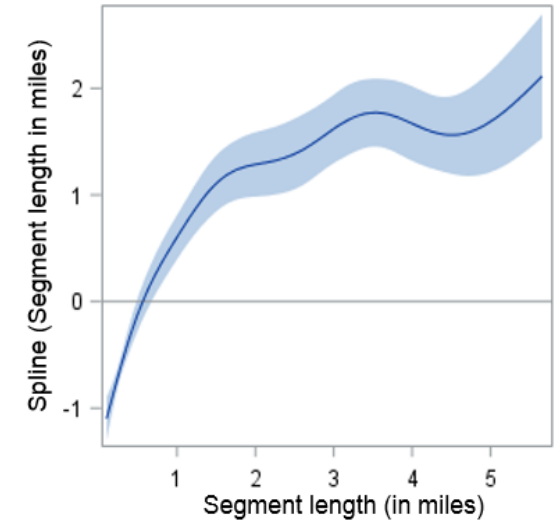
Smoothing Components for Total Injury Crashes (Category 2 NBGAMs)

With 95% Bayesian curve-wise Confidence Bands

DF = 4.55



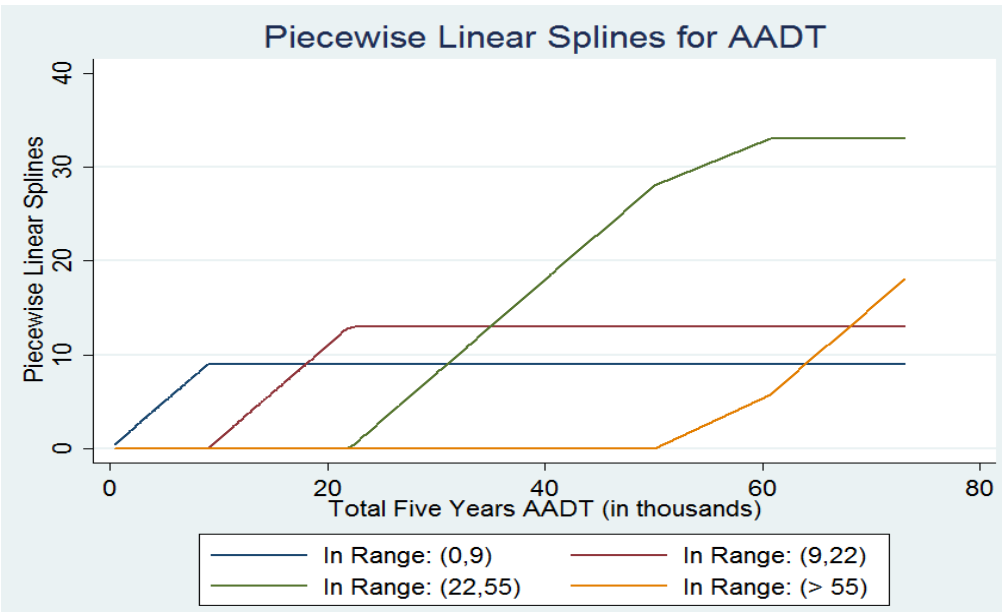
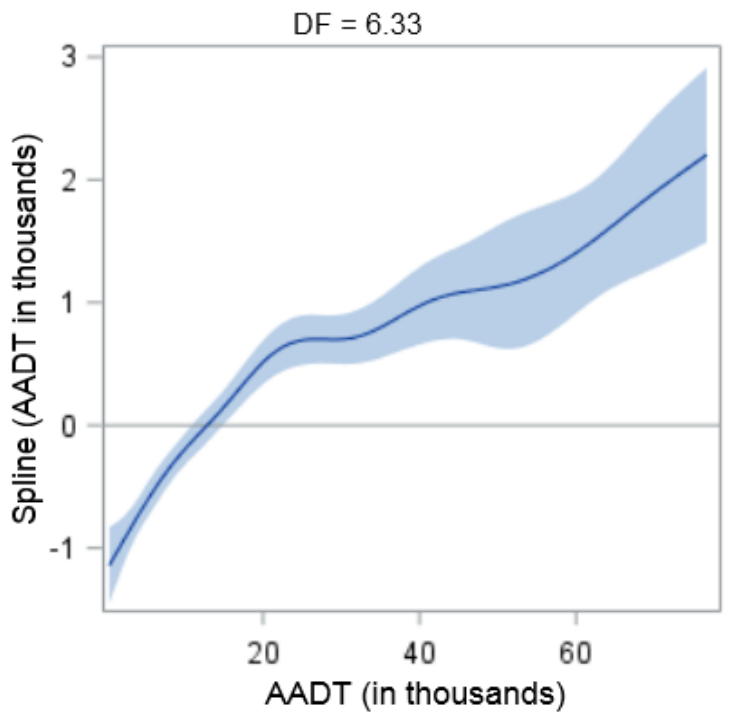
DF = 5.44



Connecting the method to practice...

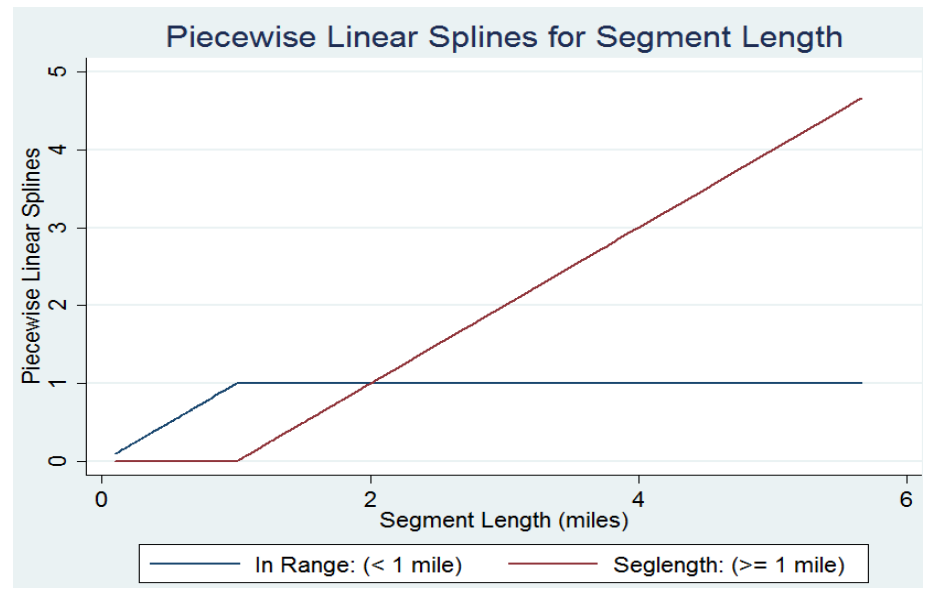
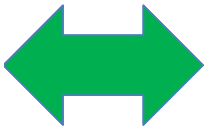
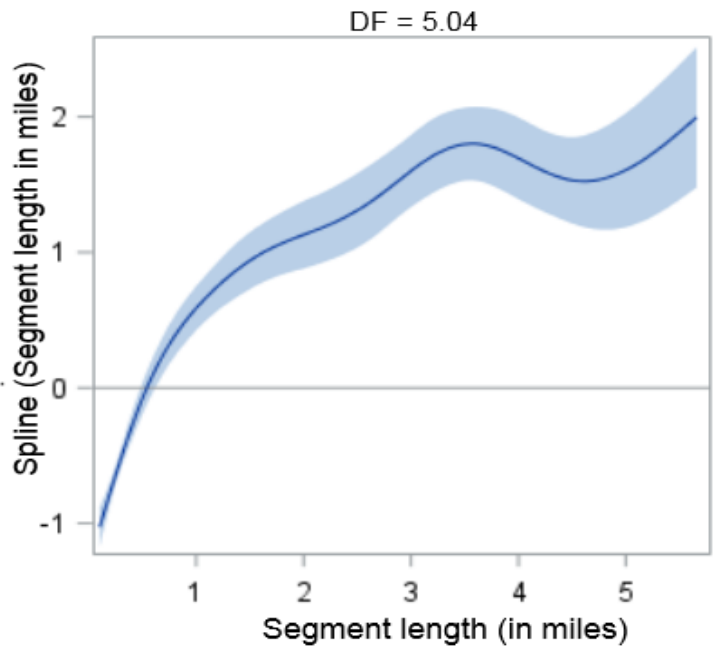
Generalized Additive Models → Piecewise Linear Count Data Models

Piecewise Linear SPFs



AADT Spline Transformations

Piecewise Linear SPFs



Segment Length Spline Transformations

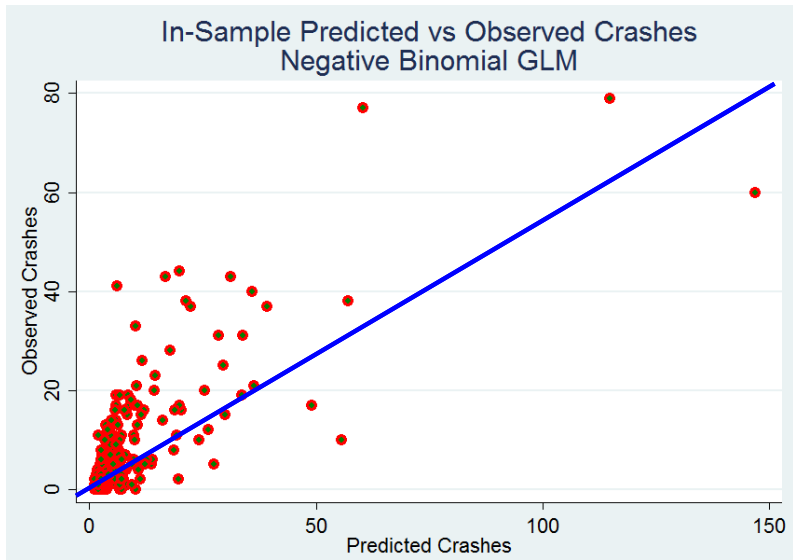
Results:
PLNB SPFs

Total Crashes

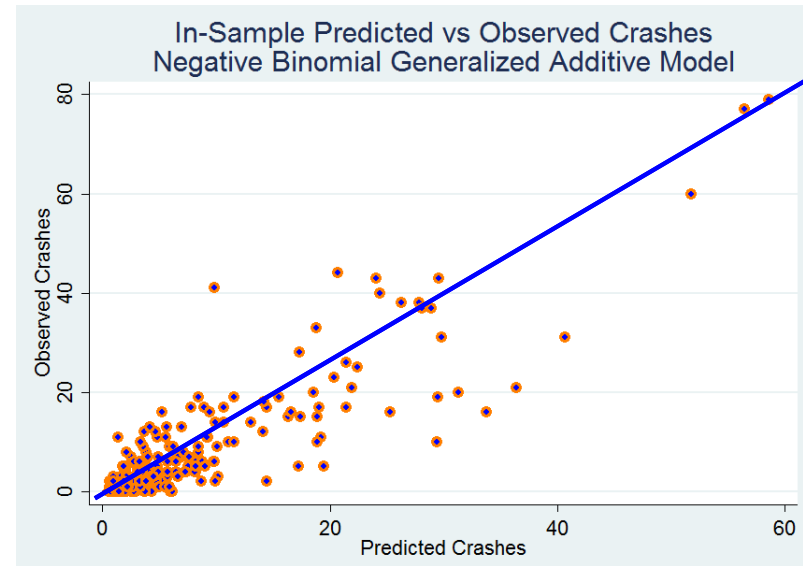
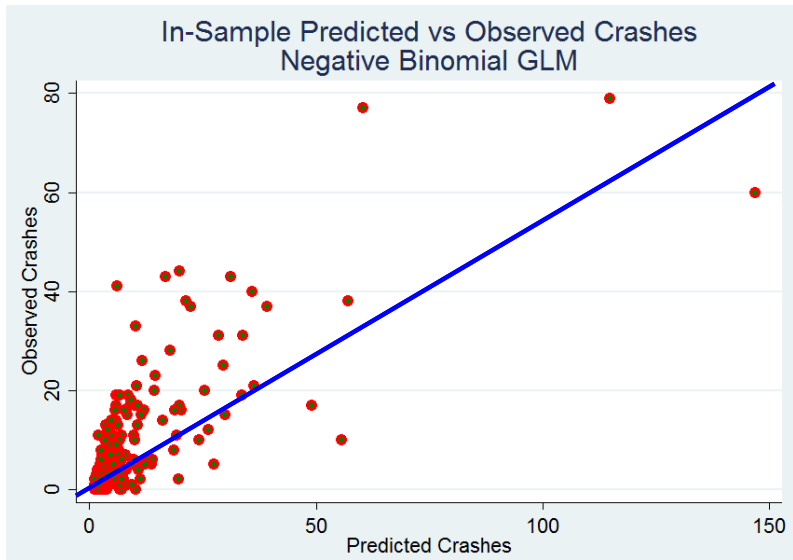
Parametric Coefficients	Category 1 PLNB			Category 2 PLNB		
	Parameter estimate	std. error	z-value	Parameter estimate	std. error	z- value
Models for total crashes						
Intercept	-0.790	0.188	-4.190	0.555	0.687	0.810
AADT1	0.110	0.024	4.590	0.103	0.024	4.260
AADT2	0.054	0.012	4.440	0.066	0.012	5.430
AADT3	0.012	0.007	1.714	0.011	0.006	1.831
AADT4	0.080	0.030	2.590	0.069	0.029	2.380
SL1	1.979	0.155	12.730	1.951	0.149	13.140
SL2	0.299	0.054	5.490	0.301	0.051	5.950
Paved	---	---	---	0.432	0.113	3.830
Lane width	---	---	---	-0.135	0.067	-2.020
Shoulder width	---	---	---	-0.050	0.017	-2.970
Dispersion Parameter	0.39	0.04	9.75	0.33	0.05	6.60

So What Test.....

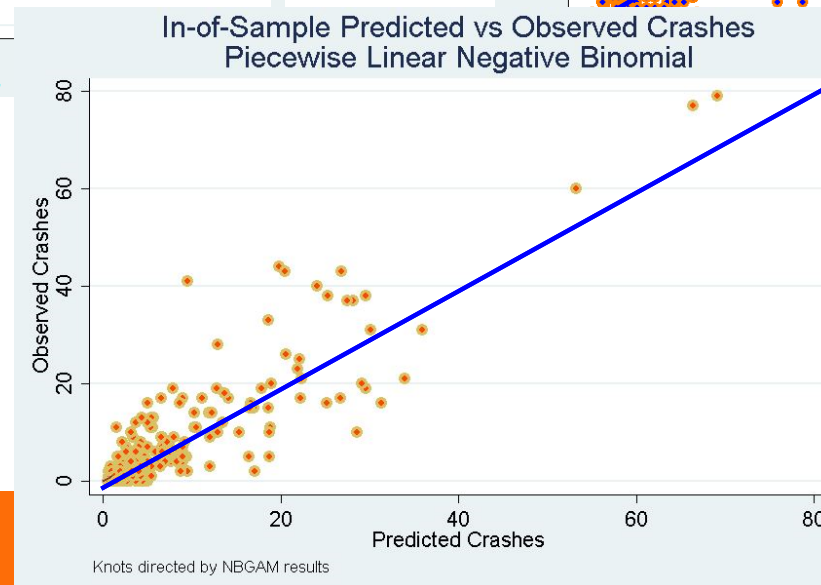
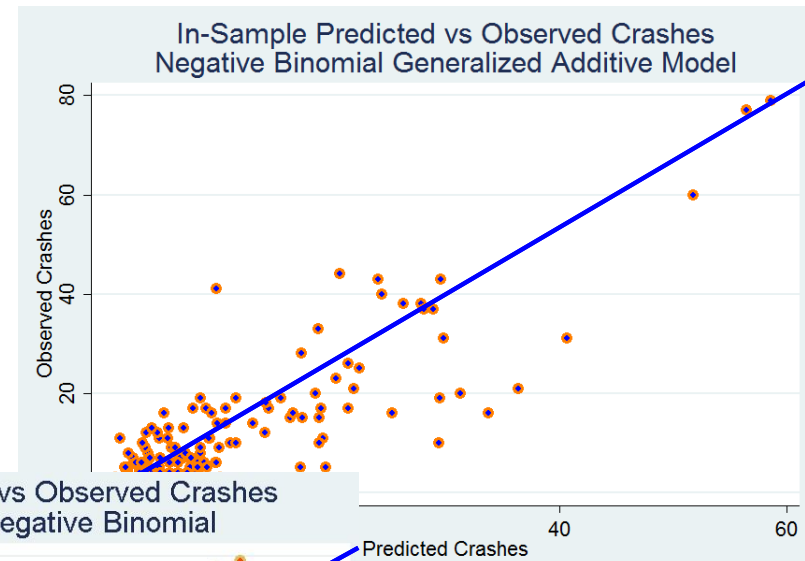
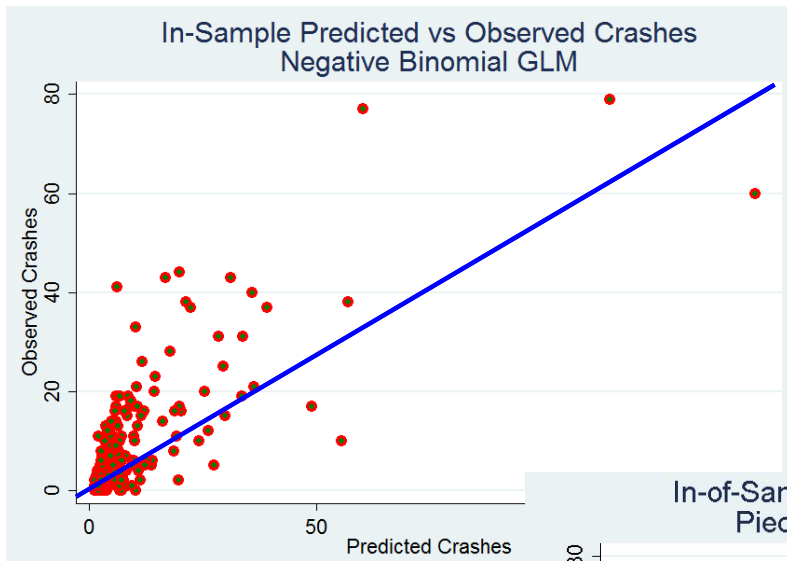
In-sample forecasts



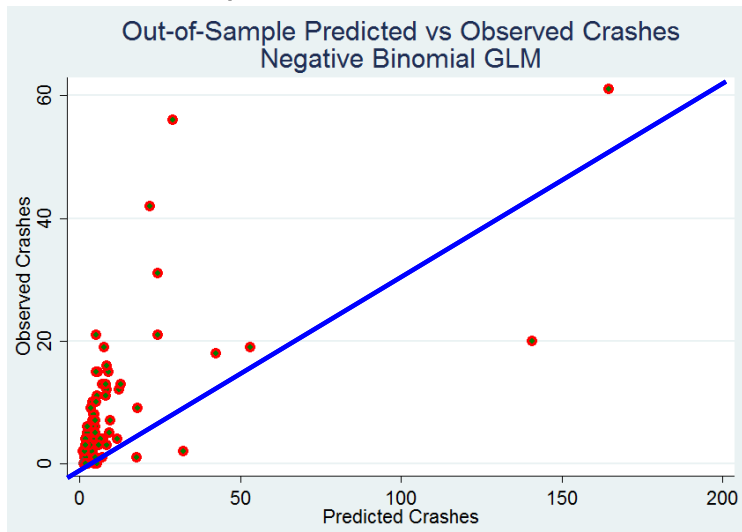
In-sample forecasts



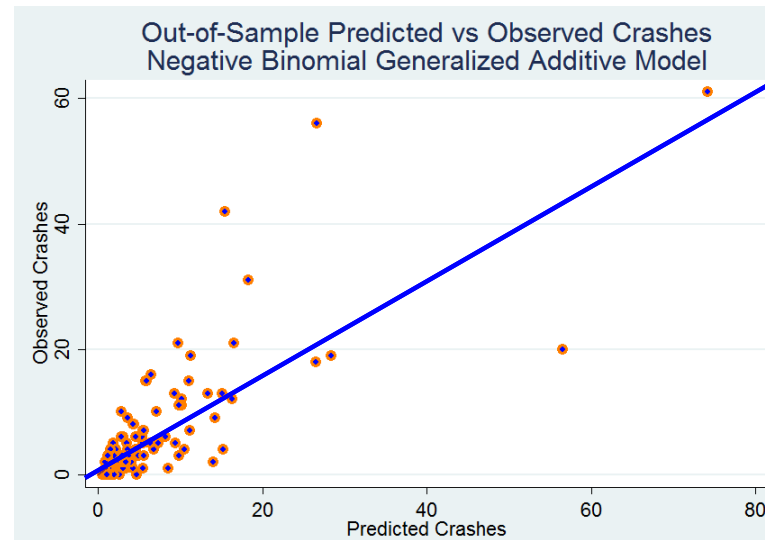
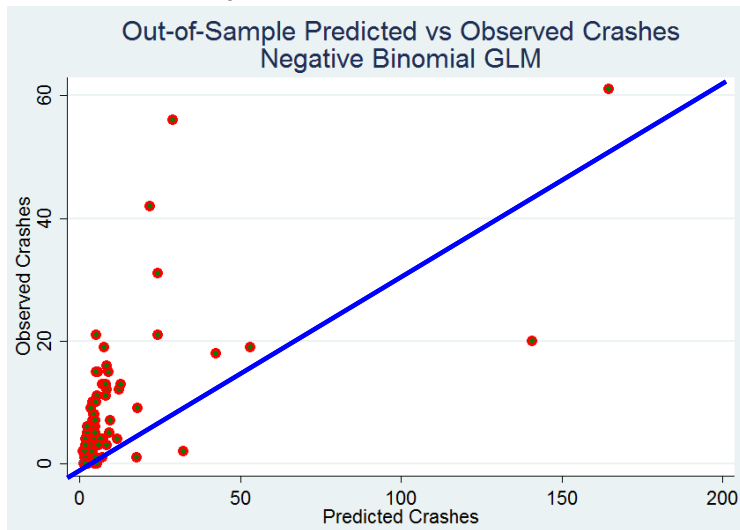
In-sample forecasts



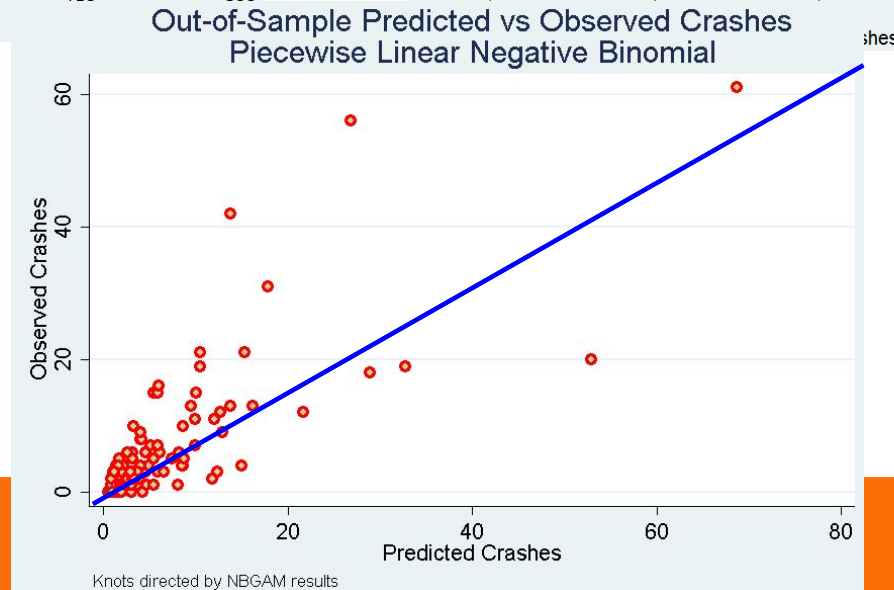
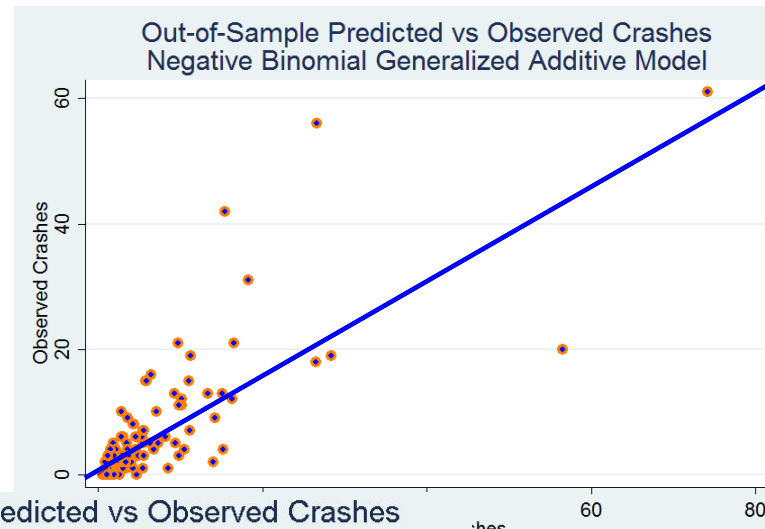
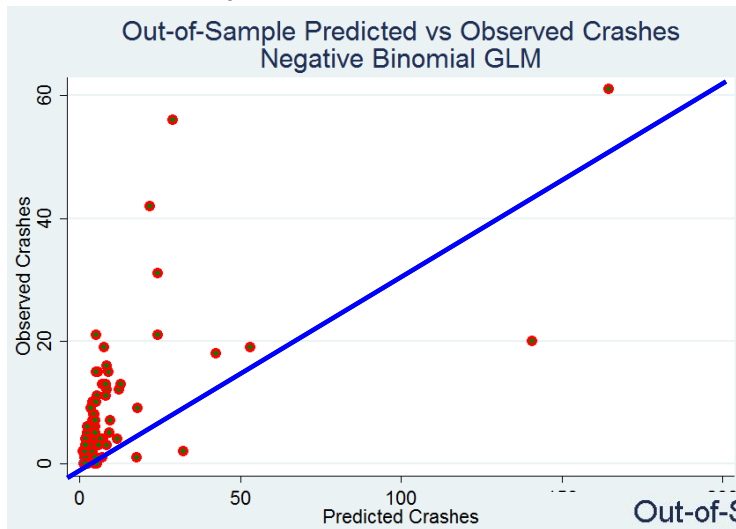
Out-of-sample forecasts



Out-of-sample forecasts



Out-of-sample forecasts



So What....?

Prediction Accuracy

		Model Comparisons						
		AADT + Segment length only						
		NBGLM		NBGAM		PLNB		
Total Crashes	P-Index	Training	Testing	Training	Testing	Training	Testing	
		MAE	5.8	6.29	3.79	3.56	3.91	3.82
		RMSE	15.2	18.34	6.36	6.36	6.36	7
		AIC	1299.47		1246.78		1242.92	
		AICC	1299.64		1248.29		1246.12	
		BIC	1313.3		1289.7		1270.49	

So What....?

Prediction Accuracy

		Model Comparisons						
		AADT + Segment length only						
		NBGLM		NBGAM		PLNB		
Total Crashes	P-Index	Training	Testing	Training	Testing	Training	Testing	
		MAE	5.8	6.29	3.79	3.56	3.91	3.82
		RMSE	15.2	18.34	6.36	6.36	6.36	7
		AIC	1299.47		1246.78		1242.92	
		AICC	1299.64		1248.29		1246.12	
		BIC	1313.3		1289.7		1270.49	
Total Injury Crashes	MAE	2.25	2.45	1.65	1.59	1.63	1.55	
	RMSE	5.52	5.95	2.82	2.72	2.77	2.75	
	AIC	869.8		831.92		826.13		
	AICC	869.98		833.04		829.25		
	BIC	883.64		868.81		854.38		

So What....?

Percentage reductions in out-of-sample prediction (testing) errors

	Models	PR	% reduction
Total Crashes	NBGAM	MAE	43
		RMSE	65
	PLNB	MAE	39
		RMSE	62
Total Injury Crashes	NBGAM	MAE	35
		RMSE	54
	PLNB	MAE	37
		RMSE	54

Take-Aways

- Quantification of non-linear dependencies → **Fusing machine learning & statistical frontiers**
- Methodological advances to **improve HSM procedures**
- More accurate predictions → **Help TDOT in screening and implementation of countermeasures**
- NBGAMs accurate but hard to interpret
- Feed knowledge from NBGAMs to PLNBs for **friendly but more accurate practical use**



Study sponsored by TDOT/US-DOT

Thank YOU

Behram Wali

bwali@vols.utk.edu

bwali.weebly.com



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

