

# Machine Learning Tools for Informing Transportation Technology Design

Missy Cummings, Ph. D.



HAL  
Humans and Autonomy Lab

Duke  
UNIVERSITY



# Using AI in Transportation Data Analytics

---

- AI in the form of machine, aka deep, learning is a very popular analytic technique
  - Best use is for large data sets which are common in transportation settings
- Machine learning reasoning is a “black box” and how models generate output is not apparent to either engineers of such algorithms or to users
  - These issues have driven the “Explainable AI” or “Interpretable AI” research thrusts
- Transportation engineers need to understand what the strengths and weaknesses of such approaches are

# What does explainable AI really mean?

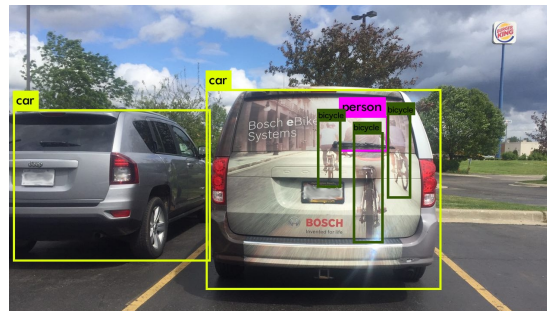
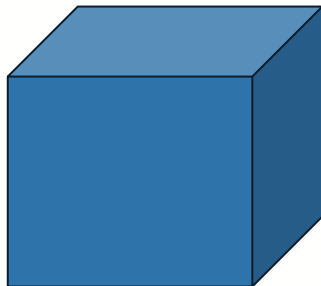
- Which algorithmic approach should I choose?
- How do I set parameters?
- What labels should I choose & where are my thresholds?



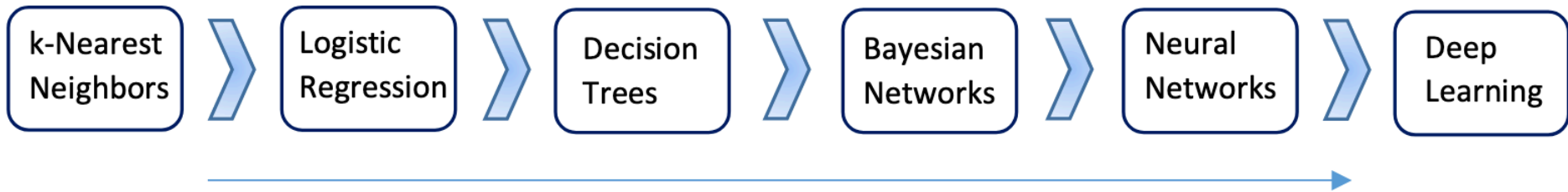
- How do I interpret the results?
- How do I adjust various parameters for the “best” sensitivity?

- How do I certify this system as safe?

1011  
1101  
1001  
1010  
...



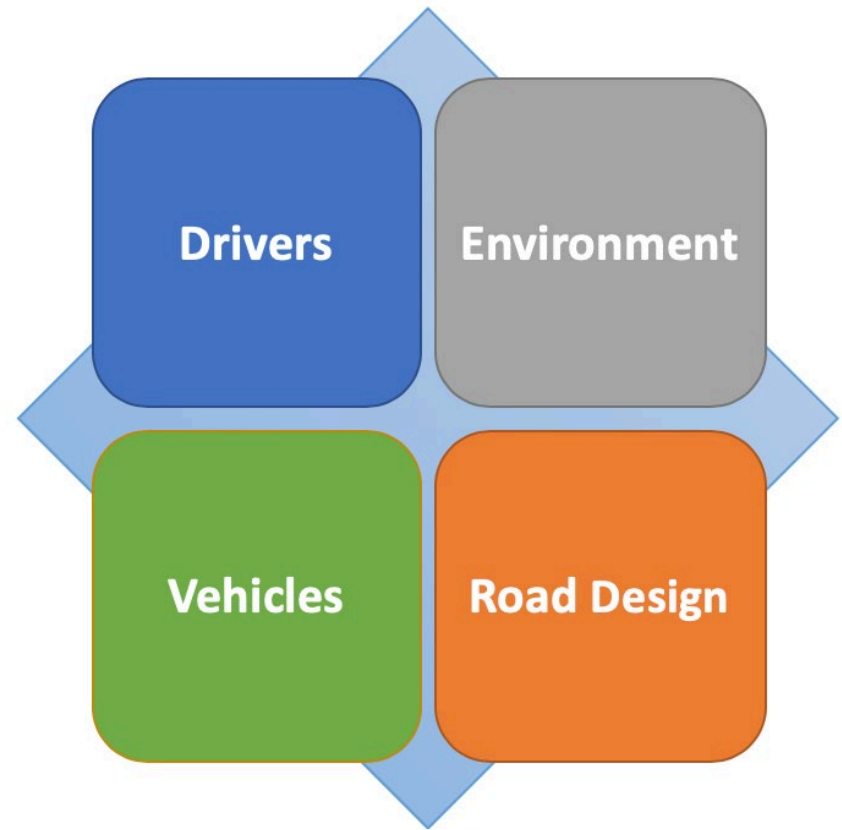
# Neural Net/LR Analyses for Explainability



- Neural nets and logistic models were chosen to illustrate model variability
- Two different data sets
  - Fatal/serious injury driver predictors (HSIS, Highway Safety Information System)
  - Pedestrians fatality predictors (NASS, National Automotive Sampling System )
- How to balance model accuracy with model utility?
  - How would different algorithm choices inform policy

# Data Set #1: Large set

- Representative research question: Are there roadway elements that contribute to driver fatalities?
- First major subjective choice: Which features/predictors should I use?
  - Choices should be informed by literature
  - Four groups of variables were determined to account for system variability

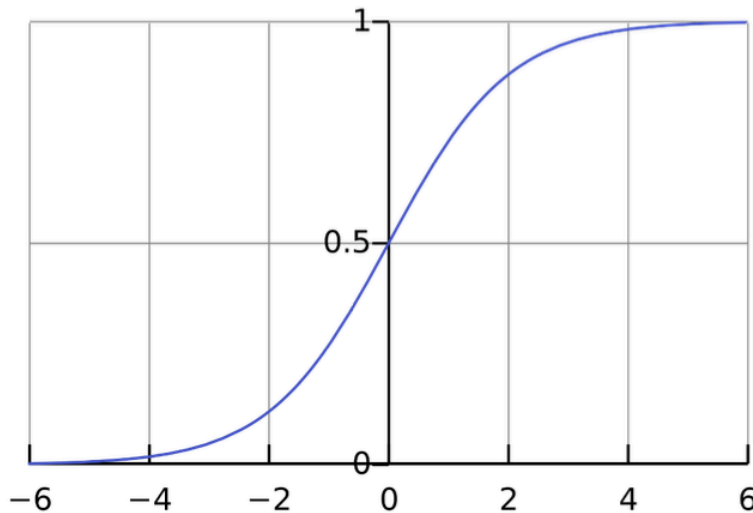


# 18 Selected Features

1	2	3	4	5	6	7	8	9
Speed limit	Average annual daily traffic	Access control type	Left shoulder width	Left shoulder width 2	Right shoulder width	Right shoulder width 2	Number of lanes	Median width
10	11	12	13	14	15	16	17	18
Section length	Light	Weather	Max driver age involved	Min driver age involved	Vehicle type	Sobriety	Urban / Rural	Lane width

 Drivers    Environment    Vehicles    Road Design

# Logistic Regression (LR) Model



$$p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$x$  = feature

$\beta_0$  = value of the criterion when features = 0

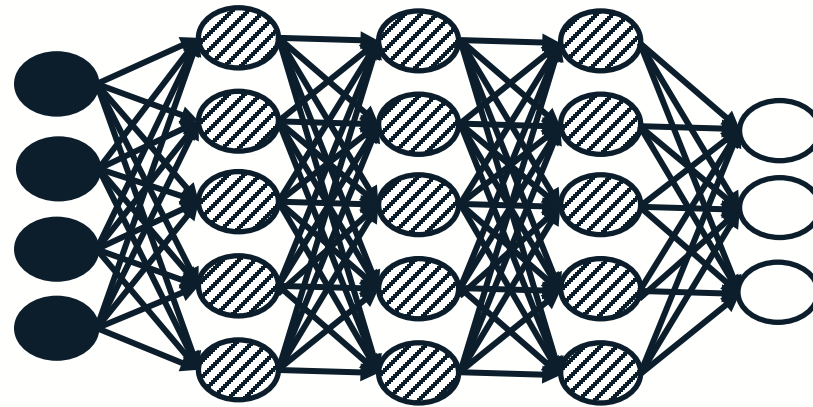
$\beta_1$  = regression coefficient

$p(x)$  = probability of success, i.e., 1

- 0 – no/minor injury, 1 – major/fatal injury
- Skewed dataset (2%)
- Model accuracy = 76%

Selection Criteria	Important Features									
Statistically significant	2	3	9	10	11	12	14	15	16	17
Odds Ratios > 2	2	10	11	15	16					

# Neural Network (NN) Model



Forward propagation (prediction)

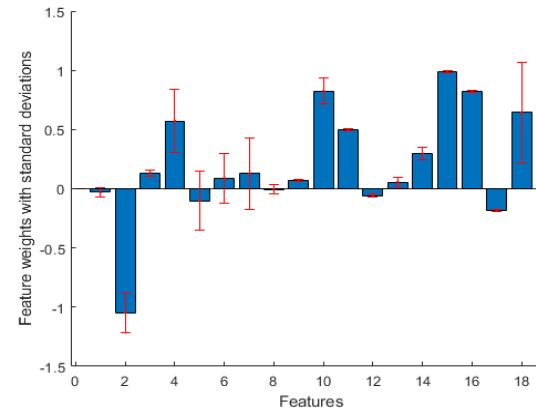
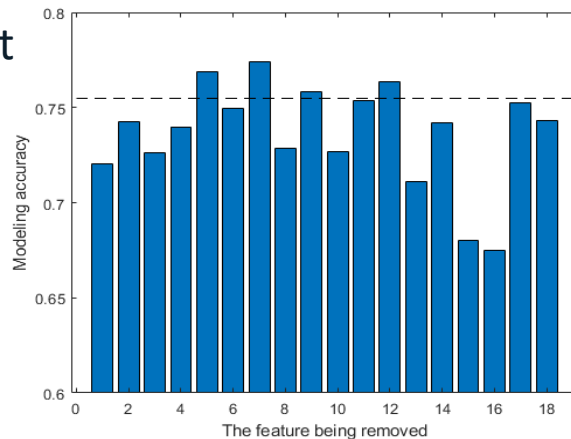
Back propagation (learning)

- 0 = no/minor injury, 1 = major/fatal injury
- 10 hidden layers
- 68% training, 12% validation, 20% testing
- Model accuracy = 75.2%



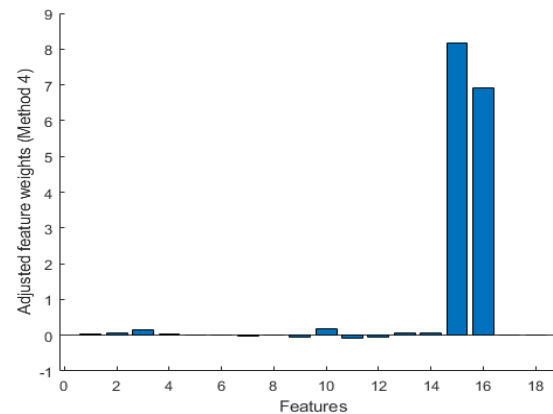
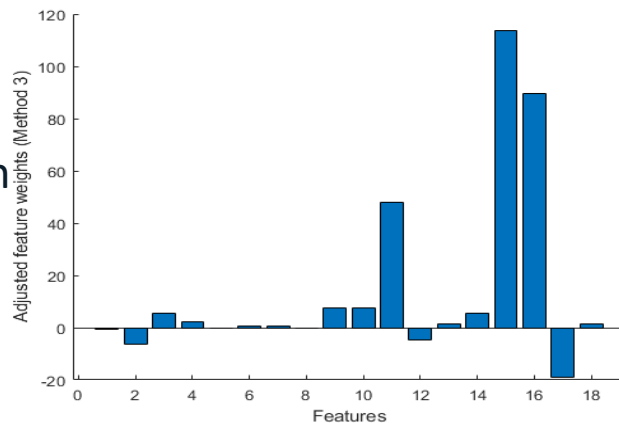
# Different interpretations of feature weights

Leave one out



Shallow NN w/ 0 hidden layers

Shallow NN weights/Stan Dev



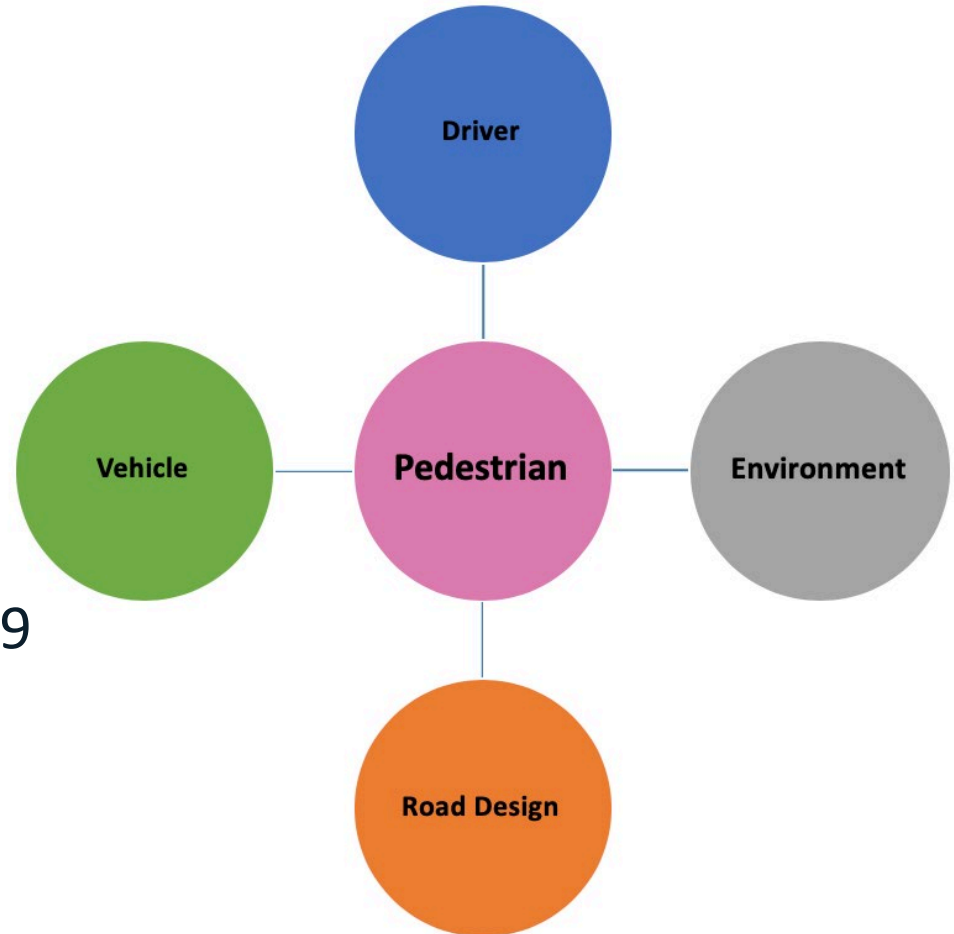
Shallow NN  
(weights/Stan Dev)\*  
accuracy drop

# Top 5 predictors for each model

Feature	LR	Method 1	Method 2	Method 3	Method 4
1 Speed limit	-	4	17	16	-
2 AADT	1	10	1	11	-
3 Access control	8	5	11	6	-
4 Left shoulder width	-	8	6	13	-
5 Left shoulder width 2	-	17	12	15	-
6 Right shoulder width	-	12	9	14	-
7 Right shoulder width 2	-	18	18	18	-
8 Number of lanes	-	7	16	17	-
9 Median width	10	15	14	10	-
10 Section length	4	6	4	9	-
11 Light	5	14	7	3	-
12 Weather	9	16	15	7	-
13 Maximum age	-	3	13	8	-
14 Minimum age	6	9	8	5	-
15 Vehicle type	2	2	2	2	2
16 Sobriety	3	1	3	1	1
17 Urban / Rural	7	13	10	4	-
18 Lane width	-	11	5	12	-

# Data Set #2: Small set

- Representative research question: Are there roadway elements that contribute to pedestrian fatalities?
- Problems with data set
  - 549 observations of pedestrian fatalities with 189 possible features
  - 310 observations, 26 fatalities
  - 18 features



# 18 Selected Features

1	2	3	4	5	6	7	8	9
Month	Time	Day of week	Pedestrian weight	Pedestrian age	Pedestrian gender	Pedestrian motion	Pedestrian action relative to vehicle	Pedestrian first avoidance action
10	11	12	13	14	15	16	17	18
Driver drinking	Speed limit	Vehicle curb weight	Driver attention to driving	Relation to junction	Traffic way flow	Number of travel lanes	Roadway surface condition	Traffic control device functioning

Pedestrian
  Drivers
  Environment
  Vehicles
  Road Design

**Model Accuracy: LR = 85%, NN = 67%**

# Top 5 predictors for each model

Feature	LR	Method 1	Method 2	Method 3	Method 4
1 Month	-	13	15	6	17
2 Time	-	5	18	18	3
3 Day of week	-	10	13	15	8
4 Pedestrian weight	-	14	10	11	13
5 Pedestrian age	2	1	3	5	1
6 Pedestrian sex	-	4	9	14	5
7 Pedestrian motion	-	6	12	12	7
8 Action relative to vehicle	-	15	11	9	15
9 First avoidance action	-	9	4	8	11
10 Driver drinking	-	3	1	3	4
11 Speed limit	1	2	2	1	2
12 Vehicle curb weight	-	18	16	13	12
13 Driver attention	-	8	5	7	10
14 Relation to junction	-	12	14	16	9
15 Traffic way flow	-	11	17	17	6
16 Number of travel lanes	-	17	6	4	18
17 Surface condition	-	7	7	2	16
18 Traffic light functioning	-	16	8	10	14

# Points of subjectivity

---

- Picking the model
- Picking which features should be included,
- Picking a p value for LR significance,
- Deciding numbers of neurons for hidden NN layers,
- Selecting data training and testing ratios,
- Picking the maximal training iteration for the NN training process,
- Picking stopping rules for training performance factors (software dependent)
- Picking threshold values between 1/0 values,
- Choosing thresholds between important/unimportant features,
- Deciding across models with slightly different results, what the actual important features were.

# Meta-Analysis

- Machine learning is just one tool in a bigger toolbox
  - It can be significantly subjective
  - Expertise is needed
- Size of data set is important
- Skewness in data can be problematic
- Descriptive vs. predictive modeling
  - Occam's Razor
- Cost/benefit considerations in results interpretation

