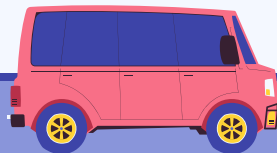


Predicting Driver At Fault In A Vehicle Collision

Julia Chen and Yuan Yin



Dataset

<https://catalog.data.gov/dataset/crash-reporting-drivers-data>

- Traffic Collisions within Montgomery County in Maryland
- Collected by Automated Crash Reporting System (ACRS) of the Maryland State Police
- 39 attributes and 184,897 instances of vehicle crashes

Report	Local Case	Agency	ACRS	Crash Date/Time	Route	F	R	G	N	Cross	SI	Off-Road	Manager	Related	Collision	Weather	Surface	Light	Vehicle	Driver	Other	Q	W	R	T	U	Y	Z	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR	AS	AT	AU	AV	AW	AX	AY	AZ	BA	BB	BC	BD	BE	BF	BG	BH	BI	BJ	BK	BL	BM	BN	BO	BP	BQ	BR	BS	BT	BU	BV	BW	BX	BY	BZ	CA	CB	CC	CD	CE	CF	CG	CH	CI	CJ	CK	CL	CM	CN	CO	CP	CQ	CR	CS	CT	CU	CV	CW	CX	CY	CZ	DA	DB	DC	DD	DE	DF	DG	DH	DI	DJ	DK	DL	DM	DN	DO	DP	DQ	DR	DS	DT	DU	DV	DW	DX	DY	DZ	EA	EB	EC	ED	EE	EF	EG	EH	EI	EJ	EK	EL	EM	EN	EO	EP	EQ	ER	ES	ET	EU	EV	EW	EX	EY	EZ	FA	FB	FC	FD	FE	FF	FG	FH	FI	FJ	FK	FL	FM	FN	FO	FP	FQ	FR	FS	FT	FU	FV	FW	FX	FY	FZ	GA	GB	GC	GD	GE	GF	GG	GH	GI	GJ	GK	GL	GM	GN	GO	GP	GQ	GR	GS	GT	GU	GV	GW	GX	GY	GZ	HA	HB	HC	HD	HE	HF	HG	HH	HI	HJ	HK	HL	HM	HN	HO	HP	HQ	HR	HS	HT	HU	HV	HW	HX	HY	HZ	IA	IB	IC	ID	IE	IF	IG	IH	II	IJ	IK	IL	IM	IN	IO	IP	IQ	IR	IS	IT	IU	IV	IW	IX	IY	IZ	JA	JB	JC	JD	JE	JF	JG	JH	JI	JJ	JK	JL	JM	JN	JO	JP	JQ	JR	JS	JT	JU	JV	JW	JX	JY	JZ	KA	KB	KC	KD	KE	KF	KG	KH	KI	KJ	KL	KM	KN	KO	KP	KQ	KR	KS	KT	KU	KV	KW	KX	KY	KZ	LA	LB	LC	LD	LE	LF	LG	LH	LI	LJ	LK	LM	LN	LO	LP	LQ	LR	LS	LT	LU	LV	LW	LX	LY	LZ	MA	MB	MC	MD	ME	MF	MG	MH	MI	MJ	MK	ML	MN	MO	MP	MQ	MR	MS	MT	MU	MV	MW	MX	MY	MZ	NA	NB	NC	ND	NE	NF	NG	NH	NI	NJ	NK	NL	NM	NO	NP	NQ	NR	NS	NT	NU	NV	NW	NX	NY	NZ	OA	OB	OC	OD	OE	OF	OG	OH	OI	OJ	OK	OL	OM	ON	OO	OP	OQ	OR	OS	OT	OU	OV	OW	OX	OY	OZ	PA	PB	PC	PD	PE	PF	PG	PH	PI	PJ	PK	PL	PM	PN	PO	PP	PQ	PR	PS	PT	PU	PV	PW	PX	PY	PZ	QA	QB	QC	QD	QE	QF	QG	QH	QI	QJ	QK	QL	QM	QN	QO	QP	QQ	QR	QS	QT	QU	QV	QW	QX	QY	QZ	RA	RB	RC	RD	RE	RF	RG	RH	RI	RJ	RK	RL	RM	RN	RO	RP	RQ	RR	RS	RT	RU	RV	RW	RX	RY	RZ	SA	SB	SC	SD	SE	SF	SG	SH	SI	SJ	SK	SL	SM	SN	SO	SP	SQ	SR	SS	ST	SU	SV	SW	SX	SY	SZ	TA	TB	TC	TD	TE	TF	TG	TH	TI	TJ	TK	TL	TM	TN	TO	TP	TQ	TR	TS	TU	TV	TW	TX	TY	TZ	UA	UB	UC	UD	UE	UF	UG	UH	UI	UJ	UK	UL	UM	UN	UO	UP	UQ	UR	US	UT	UU	UV	UW	UX	UY	UZ	VA	VB	VC	VD	VE	VF	VG	VH	VI	VJ	VK	VL	VM	VN	VO	VP	VQ	VR	VS	VT	VU	VV	VW	VX	VY	VZ	WA	WB	WC	WD	WE	WF	WG	WH	WI	WJ	WK	WL	WM	WN	WO	WP	WQ	WR	WS	WT	WU	WV	WW	WX	WY	WZ	XA	XB	XC	XD	XE	XF	XG	XH	XI	XJ	XK	XL	XM	XN	XO	XP	XQ	XR	XS	XT	XU	XV	XW	XX	XY	XZ	YA	YB	YC	YD	YE	YF	YG	YH	YI	YJ	YK	YL	YM	YN	YO	YP	YQ	YR	YS	YT	YU	YV	YW	YX	YZ	ZA	ZB	ZC	ZD	ZE	ZF	ZG	ZH	ZI	ZJ	ZK	ZL	ZM	ZN	ZO	ZP	ZQ	ZR	ZS	ZT	ZU	ZV	ZW	ZX	ZY	ZZ
NCP29	2-24-08	Gaithers Property	1876023	11:29	PAVING LOT	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0																																																																																																		

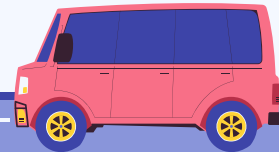


Goal

- Predict the Driver at Fault class
- Settle disagreements in a vehicle collision

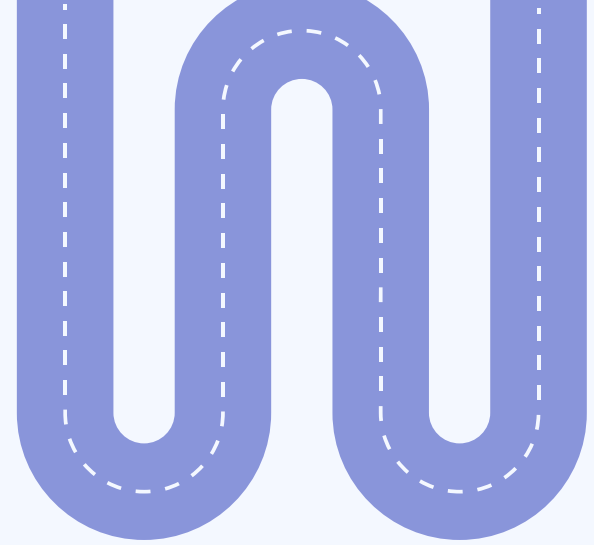
Class Distribution:

- 52% Yes
- 45% No



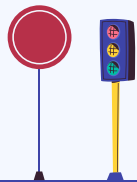
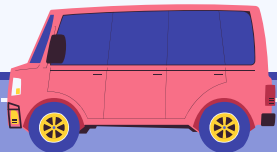
Initial Problems

- Issues opening in Weka
 - Apostrophes
 - New line in values
 - Double quotes



Preprocessing

- Remove Instances Missing the Class value
 - 184,897 instances → 180,211 instances
- Remove Attributes Missing > 70% of Their Values
- Remove Derived Attributes
 - Location vs Longitude/Latitude
- Remove Unnecessary Attributes
 - Driverless Vehicle, Person ID, Vehicle ID
 - Dataset now has 29 attributes



Preprocessing

- Stratified Random Sampling
 - 180,211 instances → 10,012 instances
- Unify Values
 - “CLEAR” vs. “Clear” and “Montgomery” vs. “Montgomery Police Department”
- Fill in Missing Values
 - Mean & Mode



Train-Validation-Test Split



80%

Train

8,008 instances



10%

Validation

1,002 instances



10%

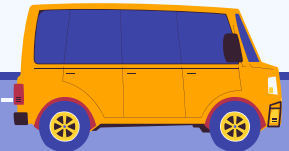
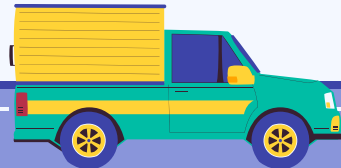
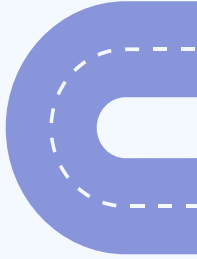
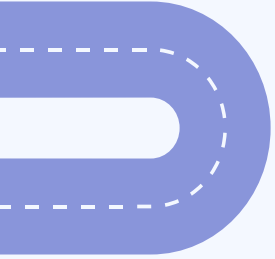
Test

1,002 instances



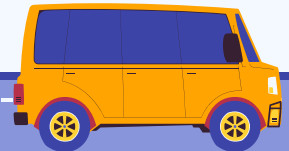
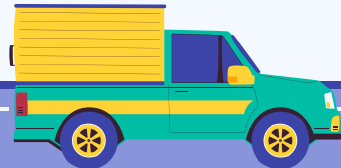
Attribute Selection

- CorrelationAttributeEval (11 attributes)
 - Evaluates the correlation between an attribute and the class variable
- CfsSubsetEval (3 attributes)
 - Selects attributes that have a high predictive ability and low intercorrelation



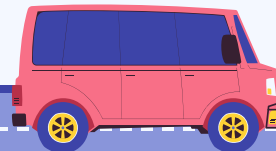
Attribute Selection

- OneRAttributeEval (13 attributes)
 - Evaluates attributes using OneR classifier
- InfoGainAttributeEval (9 attributes)
 - Selects attributes based on how much information is gained from that attribute with respect to the class
- Intuition (18 attributes)



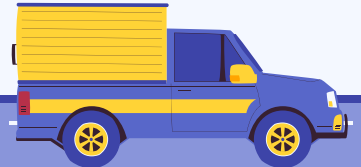
Classifiers

- NaïveBayes
 - Assume attributes are independent
 - Makes predictions by combining the likelihood of individual attributes
- KStar
 - Memorizing training instances
 - Comparing testing instance to previously existing training instances



Classifiers

- DecisionTable
 - Build and use a simple decision table using majority combinations
- OneR
 - Creates a ruleset based on one attribute



CorrelationAttributeEval

	Accuracy	True Positive	False Positive	ROC Area
NaïveBayes	76.3473	0.700	0.163	0.876
KStar	80.8383	0.828	0.215	0.883
DecisionTable	80.4391	0.834	0.230	0.886
OneR	69.5609	0.879	0.515	0.682



CfsSubsetEval

	Accuracy	True Positive	False Positive	ROC Area
NaïveBayes	74.1517	0.616	0.114	0.823
KStar	74.1517	0.601	0.097	0.824
DecisionTable	74.5509	0.612	0.101	0.829
OneR	45.01	0.000	0.032	0.484



OneRAttributeEval

	Accuracy	True Positive	False Positive	ROC Area
NaïveBayes	78.6427	0.881	0.322	0.881
KStar	69.5609	0.879	0.515	0.682
DecisionTable	80.1397	0.802	0.200	0.882
OneR	69.5609	0.879	0.515	0.682



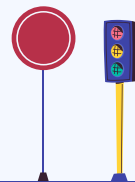
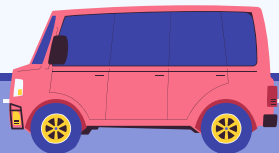
InfoGainAttributeEval

	Accuracy	True Positive	False Positive	ROC Area
NaïveBayes	73.5529	0.759	0.292	0.836
KStar	74.3513	0.769	0.285	0.797
DecisionTable	79.5409	0.866	0.285	0.881
OneR	45.01	0.000	0.032	0.484



Intuition

	Accuracy	True Positive	False Positive	ROC Area
NaïveBayes	77.5449	0.817	0.273	0.867
KStar	75.1497	0.748	0.245	0.824
DecisionTable	80.1397	0.802	0.200	0.886
OneR	69.5609	0.879	0.515	0.682



Best Model

DecisionTable classifier with CorrelationAttributeEval

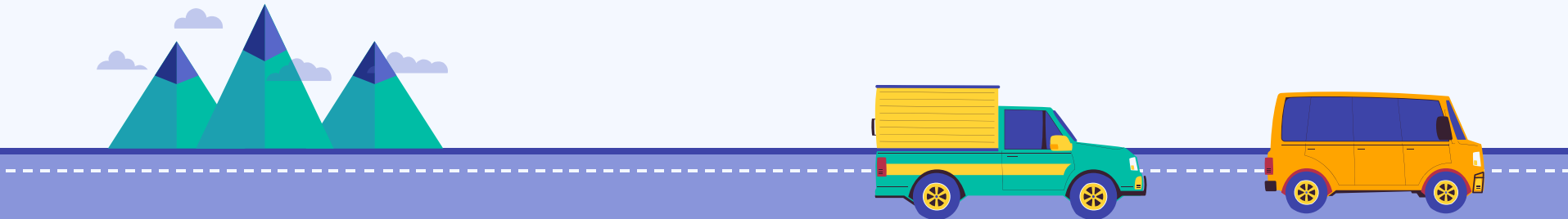
Accuracy	True Positive	False Positive	ROC Area
80.4391	0.834	0.230	0.886

- Second Highest Accuracy
- One of the highest True Positive Rates
- Highest ROC Area



Potential Improvements

- DecisionTable classifier with CorrelationAttributeEval could be improved by optimizing the cutoff value for correlation analysis



Thank you!

