

Interpretable Modeling in Machine Learning

Cynthia Rudin

Associate Professor Computer Science Department Electrical and Computer Engineering Department Duke University

Cynthia's Principles of Machine Learning

1) The Rashomon Effect: There is no "best" model for a finite dataset.

2) Thou shalt not make up a model using domain expertise alone.

3) People do not like to trust models that they don't understand.

4) Thou shalt not mistake "computationally fast" for "better."



CHADS ₂				
Risk factors	Points			
CHF	1			
HTN	1			
<u>A</u> ge≥ 75	1			
DM	1			
Stroke/TIA/embolism	2			
	Max 6			



Max 6

Data

X: patient histories Y: whether patient had stroke next year



Model

CHADS ₂				
Risk factors	Points			
CHF	1			
HTN	1			
<u>A</u> ge≥ 75	1			
DM	1			
Stroke/TIA/embolism	2			
	Max 6			

TIMI Risk Score for	STEMI
Historical	
Age 65-74	2 points
≥ 75	3 points
DM/HTN or angina	1 point
Exam	
SBP < 100	3 points
HR >100	2 points
Killip II-IV	2 points
Weight < 67 kg	1 point
Presentation	
Anterior STE or LBBB	1 point
Time to $rx > 4$ hrs	1 point
Risk Score = Total	(0 -14)
(FRONT)	



Dynamic Supervision of Sexual Offenders

STATIC-99 – TALLY SHEET

Subject Name: _____

Place of Scoring: _____

Date of S	Scoring:	Name of Assessor:				
Question Number	Risk Factor Codes					
1	Young	Aged 25 or older	0			
		Aged 18 - 24.99	1			
2	Ever Lived With	Ever lived with lover for				
		at least two years?				
		Yes	0			
		No	1			
3	Index non-sexual violence -	No	0			
	Any Convictions?	Yes	1			
4	Prior non-sexual violence -	No	0			
	Any Convictions?	Yes	1			
5	Prior Sex Offences	Charges Convictions				
		None None	0			
		1-2 1	1			
		3-5 2-3	2			
		6+ 4+	3			
6	Prior sentencing dates	3 or less	0			
	(excluding index)	4 or more	1			
7	Any convictions for non-contact	No	0			
	sex offences	Yes	1			
8	Any Unrelated Victims	No	0			
	-	Yes	1			
9	Any Stranger Victims	No	0			
		Yes	1			
10	Any Male Victims	No	0			
		Yes	1			
		Add up scores from				
	Total Score	individual risk				
		factors				

	POINTS	Risk Category
	0,1	Low
Suggested Nominal Risk Categories	2,3	Moderate-Low
	4,5	Moderate-High
	6+	High

Date of Scoring: ______ Name of Assessor: ______

Question Number	Risk Factor	Codes	Score
1	Young	Aged 25 or older	0
1	Toung	A ged $18 - 24.99$	1
2	Ever Lived With	Ever lived with lover for at least two years?	
		Yes	0
		No	1
3	Index non-sexual violence -	No	0
	Any Convictions?	Yes	1
4	Prior non-sexual violence -	No	0
	Any Convictions?	Yes	1
5	Prior Sex Offences	Charges Conviction	s
		None None 1-2 1 3-5 2-3 6+ 4+	0 1 2 3
6	Prior sentencing dates	3 or less	0
	(excluding index)	4 or more	1
7	Any convictions for non-contact	No	0
	sex offences	Yes	1
8	Any Unrelated Victims	No	0
		Yes	1
9	Any Stranger Victims	No	0
		Yes	1



		1-2 1	1
		3-5 2-3	2
		6+ 4+	3
6	Prior sentencing dates	3 or less	0
	(excluding index)	4 or more	1
7	Any convictions for non-contact	No	0
	sex offences	Yes	1
8	Any Unrelated Victims	No	0
		Yes	1
9	Any Stranger Victims	No	0
		Yes	1
10	Any Male Victims	No	0
		Yes	1
		Add up scores from	
	Total Score	individual risk	
		factors	

	<u>POINTS</u>	Risk Category
	0,1	Low
Suggested Nominal Risk Categories	2,3	Moderate-Low
	4,5	Moderate-High
	6+	High









Decision Trees

• Why trees?

• CART (Breiman 1993) is arguably the most widely used predictive modeling method used in industry currently.

Decision Trees

- Example: Will the customer wait for a table at a restaurant?
 - OthOptions: Other options, True if there are restaurants nearby.
 - Weekend: This is true if it is Friday, Saturday or Sunday.
 - Area: Does it have a bar or other nice waiting area to wait in?
 - Plans: Does the customer have plans just after dinner?
 - Price: This is either \$, \$\$, \$\$\$, or \$\$\$\$
 - Precip: Is it raining or snowing?
 - Genre: French, Mexican, Thai, or Pizza
 - Wait: Wait time estimate: 0-5 min, 6-15 min, 16-30 min, or 30+
- Crowded: Whether there are other customers (no, some, or full) Credit: Adapted from Russell and Norvig



Example: Will the customer wait for a table at a restaurant?

	OthOptions	Weekend	Area	Plans	Price	Precip	Genre	Wait	Crowded	Stay?
x ₁	Yes	No	No	Yes	\$\$\$	No	French	0-5	some	Yes
x ₂	Yes	No	No	Yes	\$	No	Thai	16-30	full	No
x ₃	No	No	Yes	No	\$	No	Pizza	0-5	some	Yes
x ₄	Yes	Yes	No	Yes	\$	No	Thai	6-15	full	Yes
x ₅	Yes	Yes	No	No	\$\$\$	No	French	30+	full	No
x ₆	No	No	Yes	Yes	\$\$	Yes	Mexican	0-5	some	Yes
X ₇	No	No	Yes	No	\$	Yes	Pizza	0-5	none	No
x ₈	No	No	No	Yes	\$\$	Yes	Thai	0-5	some	Yes
X ₉	No	Yes	Yes	No	\$	Yes	Pizza	30+	full	No
x ₁₀	Yes	Yes	Yes	Yes	\$\$\$	No	Mexican	6-15	full	No
x ₁₁	No	No	No	No	\$	No	Thai	0-5	none	No
x ₁₂	Yes	Yes	Yes	Yes	\$	No	Pizza	16-30	full	Yes

Decision Trees

• Example: Will the customer wait for a table at a restaurant? Crowded?



Standard Way to Build a Decision Tree

- Start at the top of the tree.
- Grow it by "splitting" features one by one. To split, look at how "impure" the node is.
- Assign leaf nodes the majority vote in the leaf.



• Which of these two features should we split on?





• Which of these two features should we split on?





• Which of these two features should we split on?



• Next we'll split on Plans



Standard Way to Build a Decision Tree

- Start at the top of the tree.
- Grow it by "splitting" features one by one. To split, look at how "impure" the node is.
- Assign leaf nodes the majority vote in the leaf.

• At the end, go back and prune leaves to reduce overfitting.

Why is this a good way to build a tree?

• Bottom line: Decision trees don't optimize anything

• No wonder there's a tradeoff between accuracy and interpretability in predictive modeling...









CART





Scalable Bayesian Rule Lists

- accurate
- principled
- interpretable
- scalable

Step 1 of Scalable Bayesian Rule Lists: Mine Frequent Patterns



Step 2 of Scalable Bayesian Rule Lists: Choose and Assemble Patterns into a Decision List

Example coming...



An Example of a Model for Predicting Customer Churn (IBM Watson Telco Customer Churn Data)

Input: data about each customer, and whether they churned Output: predictive model on the next slide

An Example of a Model for Predicting Customer Churn (IBM Watson Telco Customer Churn Data)

if Contract= 1 Year & StreamingMovies=Yes) -> P(churn) = 0.20else if (Contract= 1 Year) -> P(churn) = 0.05-> P (churn) = 0.70else if (Tenure<1 year & InternetService=FiberOptic) -> P(churn) = 0.03else if (Contract=2 year), else if (InternetService=FiberOptic & OnlineSecurity=No) -> P(churn) = 0.48else if (OnlineBackup=No & TechSupport=No), \rightarrow P(churn) = 0.41 else -> P(churn) = 0.22

Stroke Prediction Model (work in progress)

Input: data about each patient, and whether they had a stroke later Output: predictive model

IF past history of strokeTHEN P(stroke) = 40.5%ELSE IF patient takes warfarin reliably THEN P(stroke) = 5.6%ELSE IF age<70</td>THEN P(stroke) = 7.2%ELSE IF Blood Pressure>110THEN P(stroke) = 45%ELSE IF age<75</td>THEN P(stroke) = 6.8%ELSE IF high BMIELSEP(stroke) = 18.4\%ELSEP(stroke) = 7.2\%

Step 2 of Scalable Bayesian Rule Lists: Choose and Assemble Patterns into a Decision List

Solves a special optimization problem over decision lists.

maximize_models Posterior(model):

pile of rules size of list Posterior({c_j}_j,m,{N_{jl}}_{jl}) = $\begin{pmatrix} \prod_{l=1}^{L} \Gamma(N_{jl} + \alpha_{l}) \\ \prod_{j=0}^{m} \frac{\prod_{l=1}^{L} \Gamma(N_{jl} + \alpha_{l})}{\Gamma\left(\sum_{l=1}^{L} N_{jl} + \alpha_{l}\right)} \frac{\lambda^{m}}{\sum_{j=0}^{m} \lambda^{m}} \prod_{j=1}^{m} \frac{\eta^{c_{j}}}{\sum_{k \in R_{j-1}(c_{<A})} \eta^{k}},$

What do you want in a model anyway?

- Accuracy
 - -Data should look like it could have been generated by the model
- Gorgeousness
 - -Sparsity, Your own beliefs about the truth



P(data | model)Likelihood of data to come from model





$P(\text{model} | \text{data}) \propto P(\text{model}) \times P(\text{data} | \text{model})$ $P(\text{model} | \text{data}) \propto P(\text{model}) \times P(\text{data} | \text{model})$ $P(\text{model} | \text{model}) \times P(\text{data} | \text{model})$ $P(\text{model} | \text{model}) \times P(\text{data} | \text{model})$ $P(\text{model} | \text{model}) \times P(\text{data} | \text{model})$

Step 2 of Scalable Bayesian Rule Lists: Choose and Assemble Patterns into a Decision List

Solves a special optimization problem over decision lists.

 $P(\text{model} | \text{data}) \propto P(\text{model}) \times P(\text{data} | \text{model})$ $P(\text{model}) \propto P(\text{model}) \times P(\text{data} | \text{model})$ $P(\text{model} | \text{model}) \times P(\text{model} | \text{model}) \times P(\text{model} | \text{model})$ $P(\text{model} | \text{model}) \times P(\text{model} | \text{model}) \times P(\text{model} | \text{model})$

Step 2 of Scalable Bayesian Rule Lists: Choose and Assemble Patterns into a Decision List

Solves a special optimization problem over decision lists.

maximize_{models} Posterior(model): pile of rules size of list Posterior({c_j}_j,m,{N_{jl}}_j) = $\left(\prod_{j=0}^{m} \frac{\prod_{l=1}^{L} \Gamma(N_{jl} + \alpha_{l})}{\Gamma\left(\sum_{l=1}^{L} N_{jl} + \alpha_{l}\right)} \frac{\lambda^{m}}{\sum_{j=0}^{m} \lambda^{j}} \prod_{j=1}^{m} \frac{\eta^{c_{j}}}{\sum_{k \in R_{j-1}(c_{q,k})} \eta^{k}},$ counts for label I, rule j

Other SBRL ingredients

- Very fast bit-vector manipulation. Computational reuse. Pre-processing expensive computation.
- Theoretical bounds that are used as optimization cuts.





An Example of a Decision List for Predicting Customer Churn (IBM Watson Telco Customer Churn Data) Model from Fold 1

if Contract= 1 Year & StreamingMovies=Yes) -> P(churn) = 0.20 else if (Contract=1 Year) -> P(churn) = 0.05-> P (churn) = 0.70else if (Tenure<1 year & InternetService=FiberOptic) -> P(churn) = 0.03else if (Contract=2 year), else if (InternetService=FiberOptic & OnlineSecurity=No) \rightarrow P(churn) = 0.48 else if (OnlineBackup=No & TechSupport=No), \rightarrow P(churn) = 0.41 else -> P(churn) = 0.22



Model from Fold 2

Model from Fold 3

if (Contract=One_year & StreamingMovies=Yes) -> P(churn) = 0.25else if (Contract=One Year)-> P(churn) = 0.03else if (Contract=Two Year)-> P(churn) = 0.05else if (Tenure<1year & InternetService=Fiber_optic)</td>-> P(churn) = 0.69else if (TechSupport=No&OnlineSecurity=No)-> P(churn) = 0.45else-> P(churn) = 0.25





Runtime in seconds

	BRL	LR	SVM	CART	C4.5	RF	ADA
fold 1	0.81	0.30	2.45	0.12	0.42	2.70	6.20
fold 2	0.85	0.19	2.28	0.12	0.22	2.56	6.07
fold 3	0.87	0.18	2.34	0.11	0.23	2.60	6.05

Limits: 1 million data points, hundreds of rules: 2700 sec 50K data points, 50K rules: 9000 sec

Cynthia's Principles of Machine Learning

- 1) The Rashomon Effect: There is no "best" model for a finite dataset.
- 2) Thou shalt not make up a model using domain expertise alone.
- 3) People do not like to trust models that they don't understand.
- 4) Thou shalt not mistake "computationally fast" for "better."
- 5) If you want both accuracy and interpretability....

optimize for them

Implications in

- Energy (equipment failure prediction)
- Healthcare (scoring systems, diagnose/predict)
- Criminology (who will commit a crime?)
- Marketing (understanding customer preferences)

*Code for SBRL is publicly available on CRAN and on my website. Creative Commons License.



Thanks

*code for SBRL is publicly available on CRAN and on my website