

Explainable AI: Learning Arguments

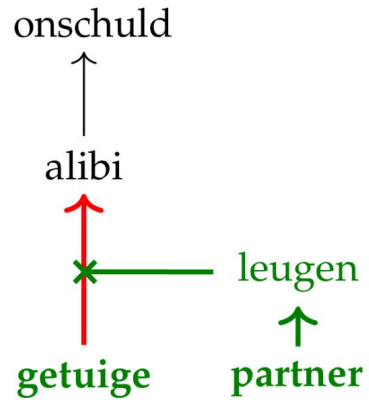
Jonas Bei, David Pomerence, Sepideh Sharbaf, Lukas Schreiner
Supervisors: Nico Roos & Pieter Collins



Agenda

- Introduction
- Learning Arguments
- Discretization
- Experiments
- Future work
- Summary

Arguments

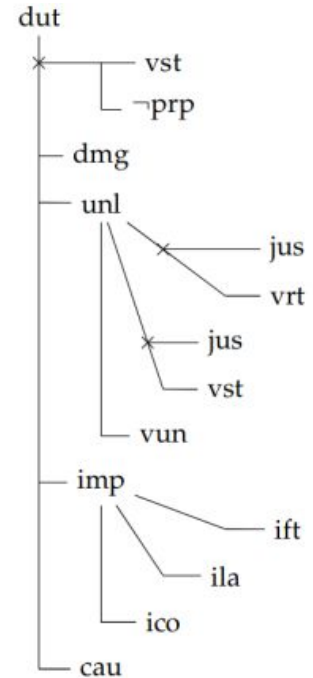


Introduction

Learning Arguments

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
¬dmg	¬dut dmg ¬unl	¬dut dmg unl ¬imp	¬dut dmg unl imp ¬cau	dut dmg unl imp cau vrt ¬vst ¬vun ift ¬ila ¬ico ¬jus prp	dut dmg unl imp cau vrt ¬vst ¬vun ift ¬ila ¬ico ¬jus prp	dut dmg unl imp cau vrt ¬vst ¬vun ift ¬ila ¬ico ¬jus prp	dut dmg unl imp cau vrt vst ¬vun ift ¬ila ¬ico ¬jus	dut dmg unl imp cau vrt vst ¬vun ift ¬ila ¬ico ¬jus	dut dmg unl imp cau vrt vst ¬vun ift ¬ila ¬ico ¬jus	dut dmg unl imp cau vrt vst vun ift ¬ila ¬ico ¬jus	dut dmg unl imp cau vrt vst vun ift ¬ila ¬ico ¬jus	dut dmg unl imp cau vrt vst vun ift ¬ila ¬ico ¬jus	¬dut dmg ¬unl vrt ¬vst jus	¬dut dmg ¬unl vrt vst jus	¬dut dmg unl imp cau vst ¬prp

1 > 2 > 3 > 4 > 5 ~ 6 ~ 7 ~ 8 ~ 9 ~ 10 ~ 11 ~ 12 ~ 13 > 14 ~ 15 ~ 16



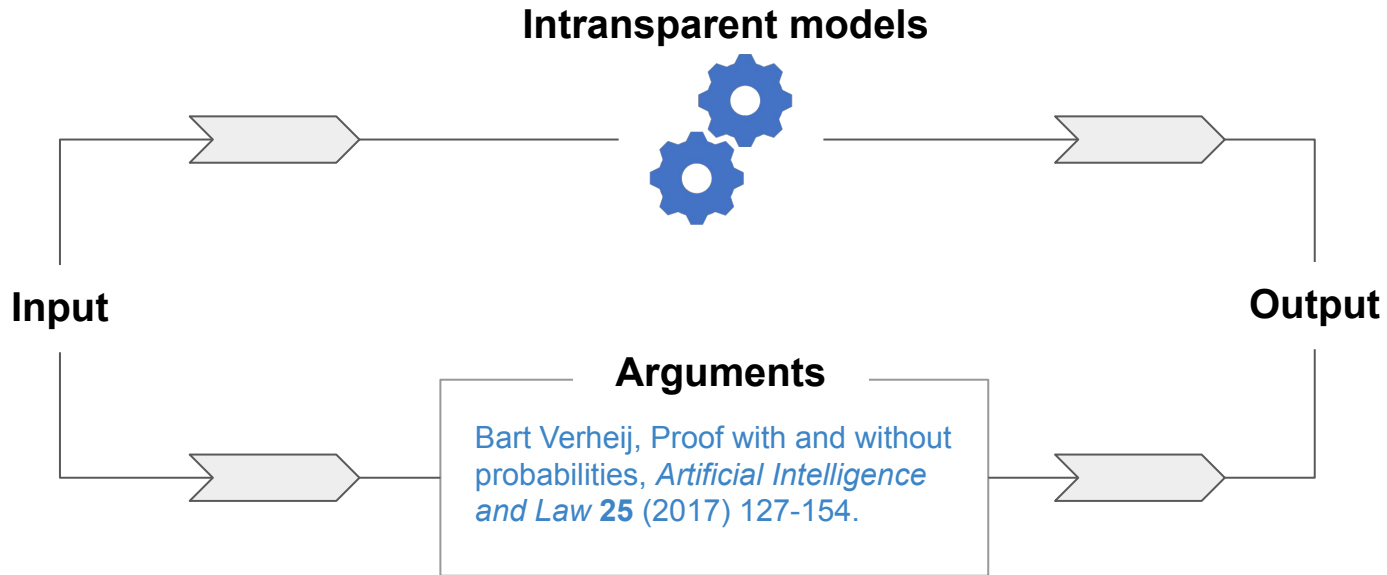
Introduction

Machine Learning setting

MEASID	TIME	LAST	MOOD	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	1000
--------	------	------	------	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	------

- Prediction
- Categorization

Explainable artificial intelligence



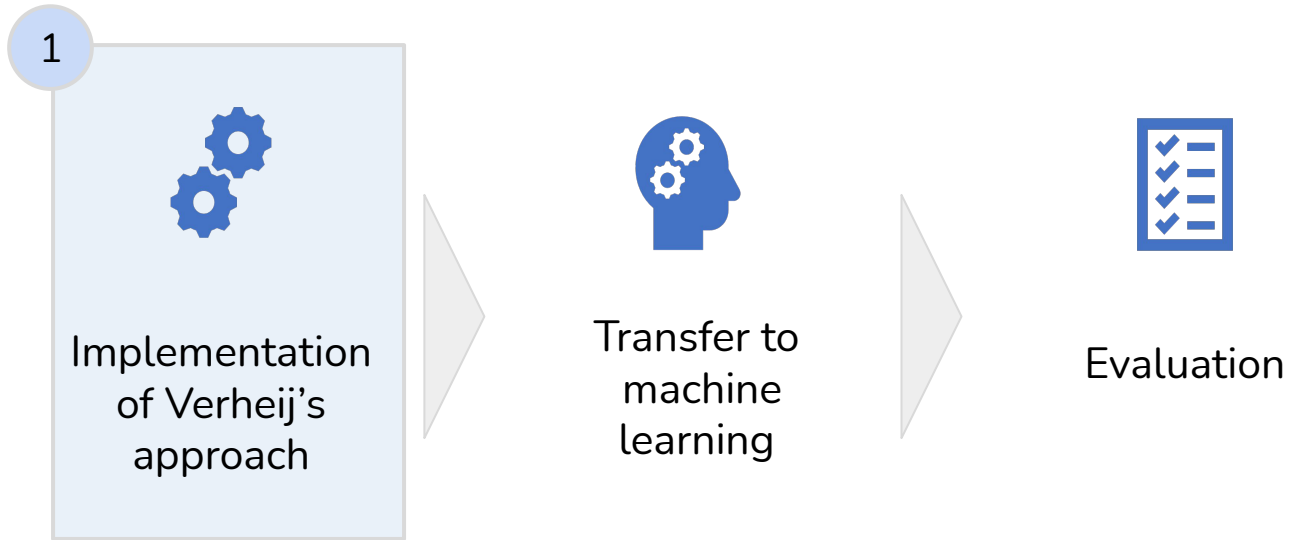
Arguments in Verheij 2017

- Coherent arguments:
Conclusion holds sometime when the premises hold.
- Presumptively valid arguments:
Conclusion holds in the most likely case where the premises hold.
- Conclusive arguments:
Conclusion always holds when the premises hold.

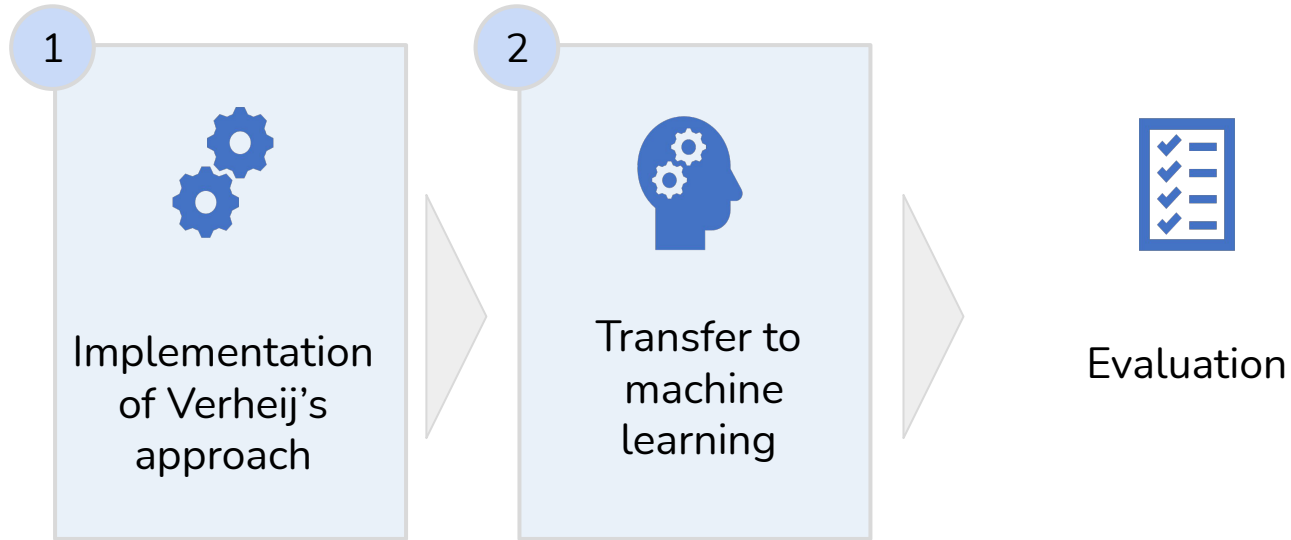
Research questions

1. Can we **reproduce** the examples from Verheij 2017 and Verheij 2020?
2. Can we **find an (efficient) algorithm** for learning arguments with this approach? How do we decide which arguments are relevant and which ones can be discarded?
3. Can we transfer the approach to a **general attribute-value classification** machine learning setting?
4. What **existing techniques** are there for learning arguments, and how do they relate to each other?
What insights can we transfer to the implementation of the approach by Verheij 2017?
5. Can we show the (in)**applicability of the approach on a real-world dataset**?
How does the approach compare with similar rule-based approaches in terms of
(a) accuracy and (b) runtime on real world datasets?
What can we infer about explainability by looking at the theories generated by the algorithms?

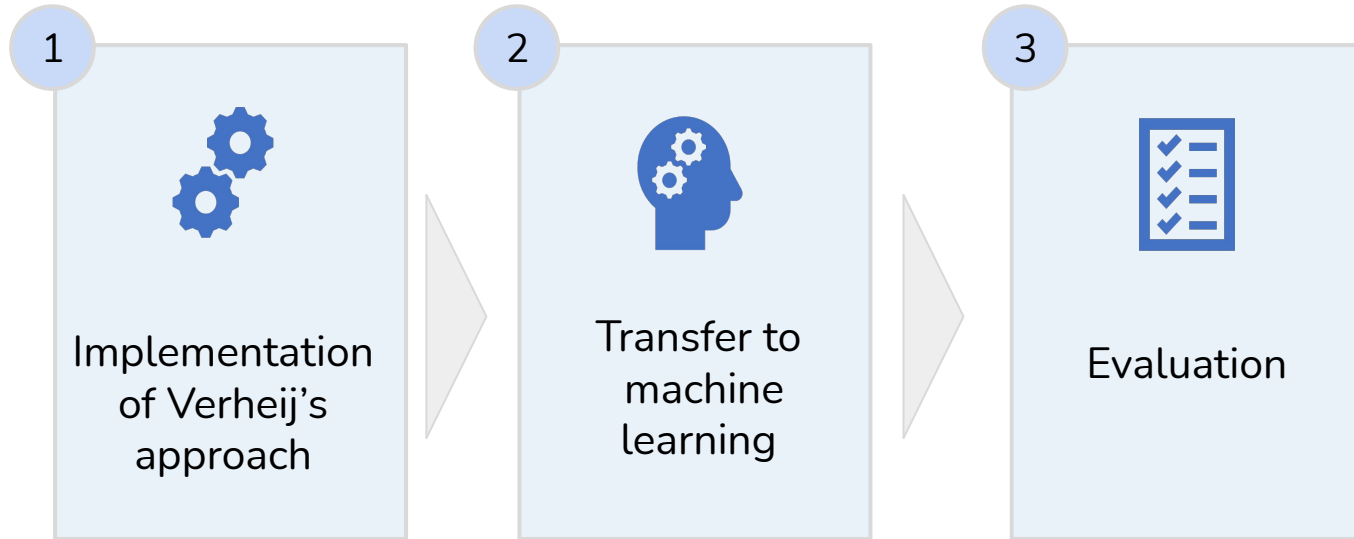
Review of progress



Review of progress



Review of progress





Learning Arguments



Learning Arguments

1. Naive search
2. Pruned search
3. HeRO algorithm
4. Decision trees

Learning Arguments

- 1. Naive search
 - 2. Pruned search
 - 3. HeRO algorithm
 - 4. Decision trees
- } Require categorical data

Learning Arguments

- 1. Naive search
 - 2. Pruned search
 - 3. HeRO algorithm
 - 4. Decision trees
- } Require categorical data

Hyperparameter optimization

Learning arguments

(1) Naive search

1. Enumerate all possible arguments:

$\rightarrow c$

$a \rightarrow c$

$\neg a \rightarrow c$

$b \rightarrow c$

$\neg b \rightarrow c$

$a, b \rightarrow c$

...

Learning arguments

(1) Naive search

1. Enumerate all possible arguments:

$\rightarrow c$

$a \rightarrow c$

$\neg a \rightarrow c$

$b \rightarrow c$

$\neg b \rightarrow c$

$a, b \rightarrow c$

...

2. Filtering irrelevant arguments

Learning arguments

(1) Naive search

1. Enumerate all possible arguments:

$$\left. \begin{array}{l} \rightarrow c \\ a \rightarrow c \\ \neg a \rightarrow c \\ b \rightarrow c \\ \neg b \rightarrow c \\ a, b \rightarrow c \\ \dots \end{array} \right\} O(2^k)$$

2. Filtering irrelevant arguments

Learning arguments

(1) Naive search

Filtering irrelevant arguments

- Discard overly specific arguments
- Keep them if they are an exception
- Merge arguments with identical premises
- Eliminate arguments that do not affect the closure of arguments

Learning arguments

(1) Naive search

Filtering irrelevant arguments

- Discard overly specific arguments
- Keep them if they are an exception
- Merge arguments with identical premises
- Eliminate arguments that do not affect the closure of arguments

$a \rightarrow b$

$b \rightarrow c$

~~$a \rightarrow c$~~

Learning arguments

(1) Naive search

Filtering irrelevant arguments

- Discard overly specific arguments
- Keep them if they are an exception
- Merge arguments with identical premises
- Eliminate arguments that do not affect the closure of arguments

$$a \rightarrow b$$

$$b \rightarrow c$$

~~$$a \rightarrow c$$~~

$$a \rightsquigarrow b$$

$$b \rightsquigarrow c$$

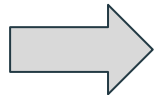
$$a \rightsquigarrow c$$

$$a \wedge b \rightsquigarrow \neg c$$

Learning arguments

(1) Naive search

```
case_model = CaseModel.fromStr([
    (3, 'inn, ¬gui, ¬evi'),
    (2, "¬inn, gui, evi, ¬evi'"),
    (1, "inn, ¬gui, evi, ¬evi'"),
    (0, "¬inn, gui, evi, evi'"),
])
```



```
evi ← ¬evi'
evi ∧ gui ← ¬inn
evi ∧ gui ∧ ¬inn ← evi'
evi ∧ ¬inn ← gui
inn ← ¬gui
inn ∧ ¬gui ← ¬evi
¬evi' ← evi ∧ inn
¬evi' ← evi ∧ ¬gui
¬gui ← inn
gui ∧ ¬evi' ∧ ¬inn ← evi
gui ∧ ¬inn ← evi ∧ ¬evi'
gui ∧ ¬inn ← ¬evi'
inn ∧ ¬evi ∧ ¬gui ←
¬evi' ← gui
¬evi' ← ¬inn
```

Learning arguments

(2) Pruned search

Observations:

(A, B) coherent

\Rightarrow

(S, B) coherent for all $S \subseteq A$

(A, B) conclusive

\Rightarrow

(S, B) conclusive for all $S \supseteq A$

where (S, B) is coherent

Learning arguments

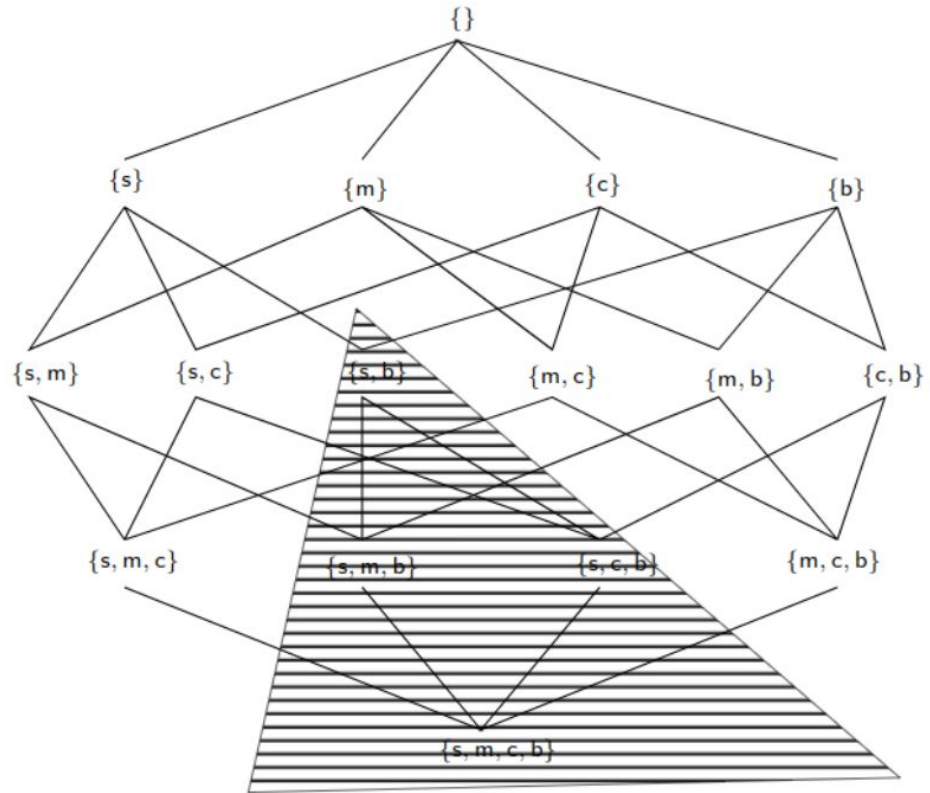
(2) Pruned search

Observations:

(A, B) coherent

\Rightarrow

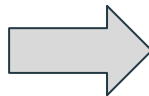
(S, B) coherent for all $S \subseteq A$



Learning arguments

(2) Pruned search

nox	,rm	,age	,dis	,rad	,tax	,ptratio	,b	,lstat	,medv
0.538	,6.575	,65.2	,4.09	,1	,296	,15.3	,396.9	,4.98	,24
0.469	,6.421	,78.9	,4.9671	,2	,242	,17.8	,396.9	,9.14	,21.6
0.469	,7.185	,61.1	,4.9671	,2	,242	,17.8	,392.83	,4.03	,34.7
0.458	,6.998	,45.8	,6.0622	,3	,222	,18.7	,394.63	,2.94	,33.4
0.458	,7.147	,54.2	,6.0622	,3	,222	,18.7	,396.9	,5.33	,36.2
0.458	,6.43	,58.7	,6.0622	,3	,222	,18.7	,394.12	,5.21	,28.7
0.524	,6.012	,66.6	,5.5605	,5	,311	,15.2	,395.6	,12.43	,22.9
0.524	,6.172	,96.1	,5.9505	,5	,311	,15.2	,396.9	,19.15	,27.1
0.524	,5.631	,100	,6.0821	,5	,311	,15.2	,386.63	,29.93	,16.5
0.524	,6.004	,85.9	,6.5921	,5	,311	,15.2	,386.71	,17.1	,18.9
0.524	,6.377	,94.3	,6.3467	,5	,311	,15.2	,392.52	,20.45	,15
0.524	,6.009	,82.9	,6.2267	,5	,311	,15.2	,396.9	,13.27	,18.9
0.524	,5.889	,39	,5.4509	,5	,311	,15.2	,390.5	,15.71	,21.7
0.538	,5.949	,61.8	,4.7075	,4	,307	,21	,396.9	,8.26	,20.4
0.538	,6.096	,84.5	,4.4619	,4	,307	,21	,380.02	,10.26	,18.2
0.538	,5.834	,56.5	,4.4986	,4	,307	,21	,395.62	,8.47	,19.9
0.538	,5.935	,29.3	,4.4986	,4	,307	,21	,386.85	,6.58	,23.1
0.538	,5.99	,81.7	,4.2579	,4	,307	,21	,386.75	,14.67	,17.5
0.538	,5.456	,36.6	,3.7965	,4	,307	,21	,288.99	,11.69	,20.2
0.538	,5.727	,69.5	,3.7965	,4	,307	,21	,390.95	,11.28	,18.2
0.538	,5.57	,98.1	,3.7979	,4	,307	,21	,376.57	,21.02	,13.6



```

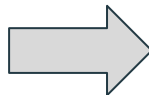
~high_rad <- ~high_lstat & ~high_medv & ~high_ptratio & ~high_tax
~high_tax <- high_b & high_ptratio & high_rad & ~high_lstat & ~high_medv
~high_tax <- high_b & high_ptratio & ~high_lstat & ~high_medv
~high_tax <- high_b & high_ptratio & ~high_lstat & ~high_medv & ~high_rad
~high_tax <- high_lstat & high_medv & high_ptratio & ~high_b & ~high_rad
high_b <- high_lstat & high_rad & high_tax & ~high_ptratio
high_b <- high_lstat & ~high_ptratio & ~high_rad & ~high_tax
high_b <- high_lstat & ~high_ptratio & ~high_tax
high_b <- high_medv & high_ptratio & high_rad & ~high_tax
high_b <- high_medv & high_ptratio & high_tax & ~high_lstat & ~high_rad

```

Learning arguments

(2) Pruned search

nox	,rm	,age	,dis	,rad	,tax	,ptratio	,b	,lstat	,medv
0.538	,6.575	,65.2	,4.09	,1	,296	,15.3	,396.9	,4.98	,24
0.469	,6.421	,78.9	,4.9671	,2	,242	,17.8	,396.9	,9.14	,21.6
0.469	,7.185	,61.1	,4.9671	,2	,242	,17.8	,392.83	,4.03	,34.7
0.458	,6.998	,45.8	,6.0622	,3	,222	,18.7	,394.63	,2.94	,33.4
0.458	,7.147	,54.2	,6.0622	,3	,222	,18.7	,396.9	,5.33	,36.2
0.458	,6.43	,58.7	,6.0622	,3	,222	,18.7	,394.12	,5.21	,28.7
0.524	,6.012	,66.6	,5.5605	,5	,311	,15.2	,395.6	,12.43	,22.9
0.524	,6.172	,96.1	,5.9505	,5	,311	,15.2	,396.9	,19.15	,27.1
0.524	,5.631	,100	,6.0821	,5	,311	,15.2	,386.63	,29.93	,16.5
0.524	,6.004	,85.9	,6.5921	,5	,311	,15.2	,386.71	,17.1	,18.9
0.524	,6.377	,94.3	,6.3467	,5	,311	,15.2	,392.52	,20.45	,15
0.524	,6.009	,82.9	,6.2267	,5	,311	,15.2	,396.9	,13.27	,18.9
0.524	,5.889	,39	,5.4509	,5	,311	,15.2	,390.5	,15.71	,21.7
0.538	,5.949	,61.8	,4.7075	,4	,307	,21	,396.9	,8.26	,20.4
0.538	,6.096	,84.5	,4.4619	,4	,307	,21	,380.02	,10.26	,18.2
0.538	,5.834	,56.5	,4.4986	,4	,307	,21	,395.62	,8.47	,19.9
0.538	,5.935	,29.3	,4.4986	,4	,307	,21	,386.85	,6.58	,23.1
0.538	,5.99	,81.7	,4.2579	,4	,307	,21	,386.75	,14.67	,17.5
0.538	,5.456	,36.6	,3.7965	,4	,307	,21	,288.99	,11.69	,20.2
0.538	,5.727	,69.5	,3.7965	,4	,307	,21	,390.95	,11.28	,18.2
0.538	,5.57	,98.1	,3.7979	,4	,307	,21	,376.57	,21.02	,13.6



```

~high_rad <- ~high_lstat ^ ~high_medv ^ ~high_ptratio ^ ~high_tax
~high_tax <- high_b ^ high_ptratio ^ high_rad ^ ~high_lstat ^ ~high_medv
~high_tax <- high_b ^ high_ptratio ^ ~high_lstat ^ ~high_medv
~high_tax <- high_b ^ high_ptratio ^ ~high_lstat ^ ~high_medv ^ ~high_rad
~high_tax <- high_lstat ^ high_medv ^ high_ptratio ^ ~high_b ^ ~high_rad
high_b <- high_lstat ^ high_rad ^ high_tax ^ ~high_ptratio
high_b <- high_lstat ^ ~high_ptratio ^ ~high_rad ^ ~high_tax
high_b <- high_lstat ^ ~high_ptratio ^ ~high_tax
high_b <- high_medv ^ high_ptratio ^ high_rad ^ ~high_tax
high_b <- high_medv ^ high_ptratio ^ high_tax ^ ~high_lstat ^ ~high_rad

```

Parameters:

- maximum size of premises
- maximum depth of exceptions

HeRO algorithm



- Incremental approach
- Information gain
- Maximum information gain

An algorithm for the induction of defeasible logic theories from databases

Benjamin Johnston

Guido Governatori

School of Information Technology and Electrical Engineering
The University of Queensland
Brisbane, Queensland, Australia

Email: superhero@benjaminjohnston.com.au, guido@itee.uq.edu.au

Abstract

Defeasible logic is a non-monotonic logic with applications in rule-based domains such as law. To ease the development and improve the accuracy of expert systems based on defeasible logic, it is desirable to automatically induce a theory of the logic from a training set of precedent data. Empirical evidence suggests that minimal theories that describe the training set tend to be more faithful representations of reality. We show via transformation from the hitting set problem that this global minimization problem is intractable, belonging to the class of NP optimisation problems. Given the inherent difficulty of finding the optimal solution, we instead use heuristics and demonstrate that a best-first, greedy, branch and bound algorithm can be used to find good theories in short time. This approach displays significant improvements in both accuracy and theory size as compared to recent work in the area that post-processed the output of an Apriori association rule-mining algorithm, with comparable execution times.

Keywords: Defeasible Logic, Machine Learning, Association Rules

1 Introduction

Expert and decision support systems are slowly making inroads into the legal community, but unfortunately they currently appear to be limited in terms of either the difficulty of their construction or their inability to justify their reasoning processes to the user. Existing systems could be roughly classified into two broad categories: expert systems that are constructed by manual encoding of knowledge (Zeleznirow & Hunter 1994), and classification tools that are automatically trained from precedent data using machine learning or data mining techniques (Zeleznirow & Stranieri 1997, Brüninghaus & Ashley 1999). Unsurprisingly, the expense involved in employing human experts and the difficulty that experts have in expressing the reasoning behind their “intuition” can eliminate the option of building expert systems in spite

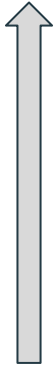
expert systems, but facilitating construction via automatic induction.

Recent work in the field of non-monotonic logics suggests the suitability of the formalism as an underlying model for such reasoning, that turns out (as we will show) to be conducive to automatic induction. Non-monotonic logics, such as defeasible logic, were originally developed to simplify reasoning with incomplete information (Ginsberg 1993). In contrast to monotonic logics whereby a conclusion of a theory remains valid irrespective of how many assertions are added to the theory, non-monotonic logics can reach tentative conclusions that may be overridden (and replaced with a contrary conclusion) in light of additional information. Defeasible logic is one of many non-monotonic logics in use, but is particularly desirable for use in information systems because it matches the non-monotonicity of legal reasoning and is computationally efficient without sacrificing too much expressiveness. The extension of a defeasible logic theory has been shown to be computable in linear time (Antonioni, Billington, Governatori & Maher 2001, Maher 2001), as opposed to the NP-hardness or even undecidability of most non-monotonic and monotonic logics (Prakken 1997, Ginsberg 1993). While some expressiveness is sacrificed in using defeasible logic over first order logic, it still remains quite suited to legal domains. Antonioni et al (1999) have demonstrated the extremely high correspondence between regulatory documents and their equivalent expression as defeasible logic theories, in some cases the correspondence is almost 1-1 between sentences in legal documents and logical rules.

2 Defeasible Logic

In this section we will present a formal explanation of defeasible logic. Because we are focussing our attention to a specific application of defeasible logic, for simplicity our terminology slightly deviates from that used in other papers. A more com-

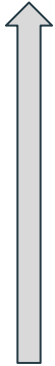
HeRO algorithm



rule 2



HeRO algorithm



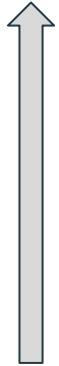
rule 2



rule 3



HeRO algorithm



rule 1



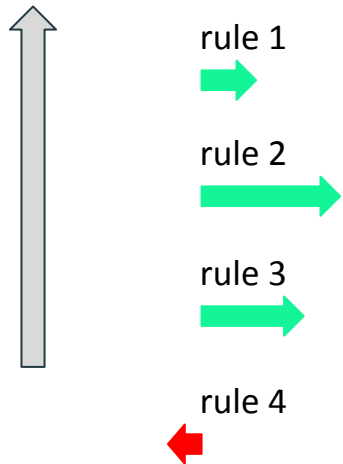
rule 2



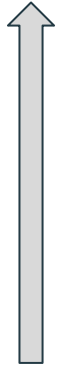
rule 3



HeRO algorithm



HeRO algorithm



rule 1



rule 2



rule 3



Learning arguments

(4) Decision trees

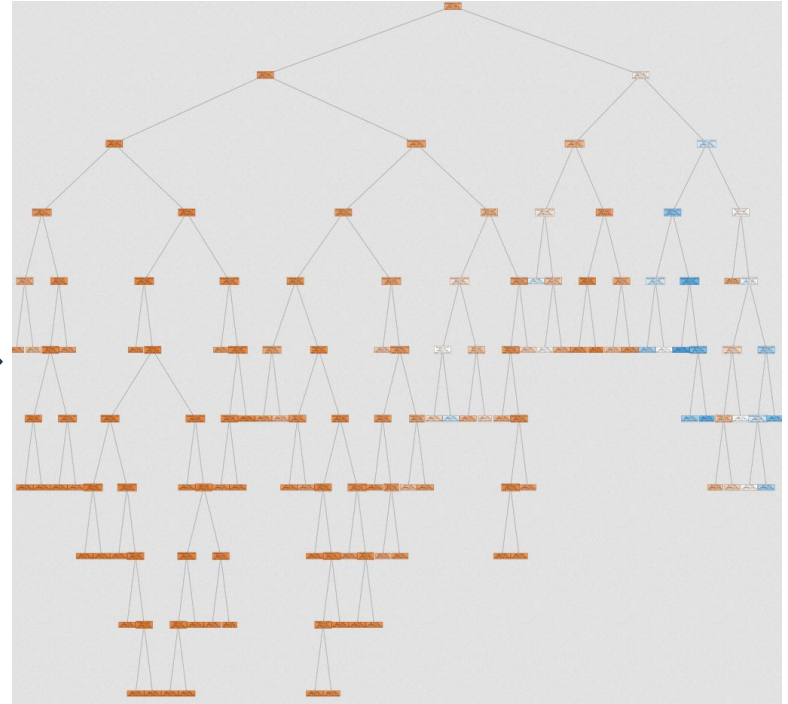
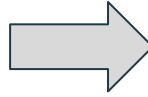
- Efficient on large data sets
- Entropy-based discretization
- Small number of arguments thanks to pruning
- No exceptions

Learning arguments

(4) Decision trees

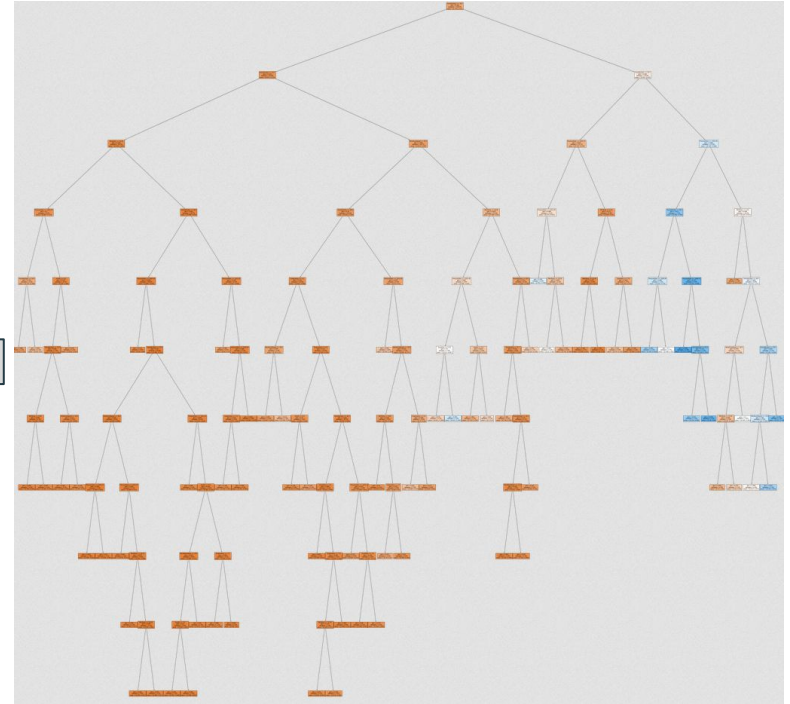
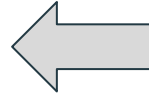
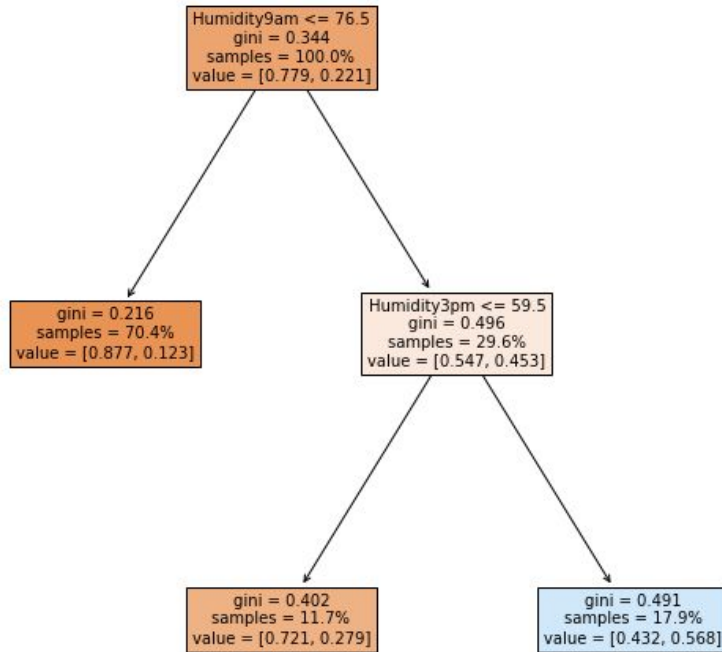
	MinTemp	MaxTemp	Evaporation	Sunshine	WindGustSpeed	WindSpeed9am	WindSpeed3pm	Humidity
6049	17.9	35.2	12.0	12.3	48.0	6.0	20.0	2
6050	18.4	28.9	14.8	13.0	37.0	19.0	19.0	3
6052	19.4	37.6	10.8	10.6	46.0	30.0	15.0	4
6053	21.9	38.4	11.4	12.2	31.0	6.0	6.0	3
6054	24.2	41.0	11.2	8.4	35.0	17.0	13.0	1
...
142298	19.3	33.4	6.0	11.0	35.0	9.0	20.0	6
142299	21.2	32.6	7.6	8.6	37.0	13.0	11.0	5
142300	20.7	32.8	5.6	11.0	33.0	17.0	11.0	4
142301	19.5	31.8	6.2	10.6	26.0	9.0	17.0	6
142302	20.2	31.7	5.6	10.7	30.0	15.0	7.0	7

56420 rows x 92 columns



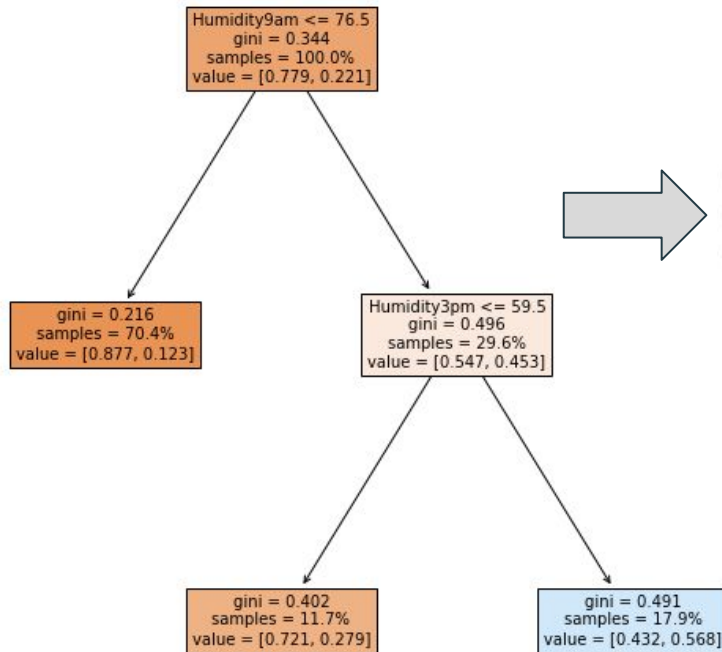
Learning arguments

(4) Decision trees



Learning arguments

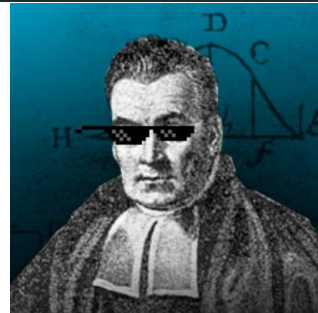
(4) Decision trees



→

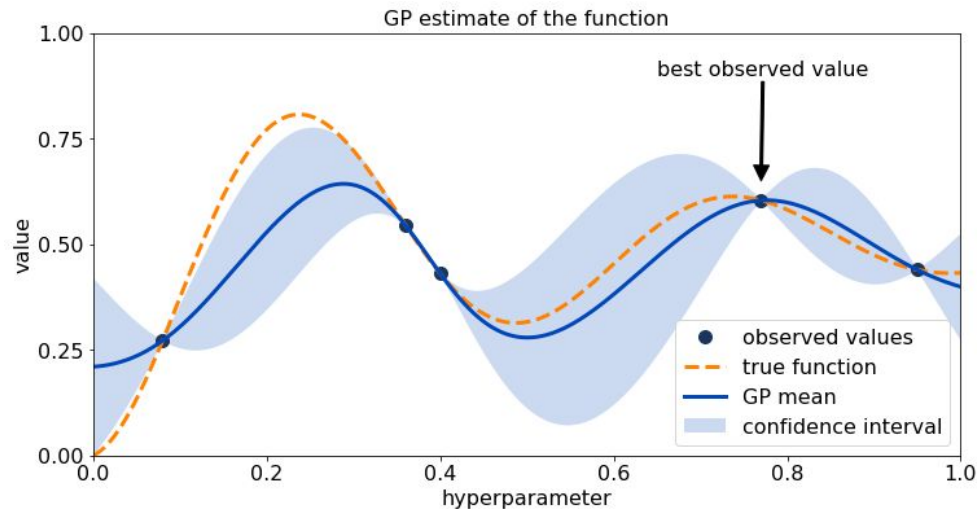
(Humidity9am <= 76.5) -> NoRain
(Humidity9am > 76.5), (Humidity3pm > 59.5) -> Rain
(Humidity9am > 76.5), (Humidity3pm <= 59.5) -> NoRain

Bayesian Optimization

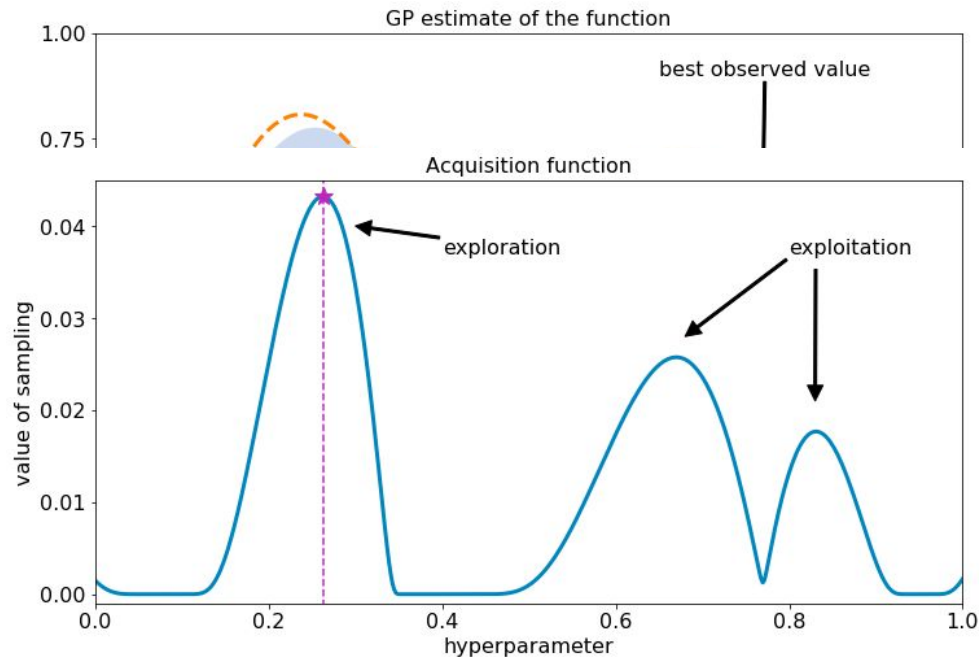


- Sampled points are used to estimate the objective function (prior)
- Points are sampled using an acquisition function and the prior is updated
 - Acquisition function balances exploration & exploitation via uncertainty in the prior
- Prior is updated
- After a given number of iterations, use a numerical method to find the estimated optimum

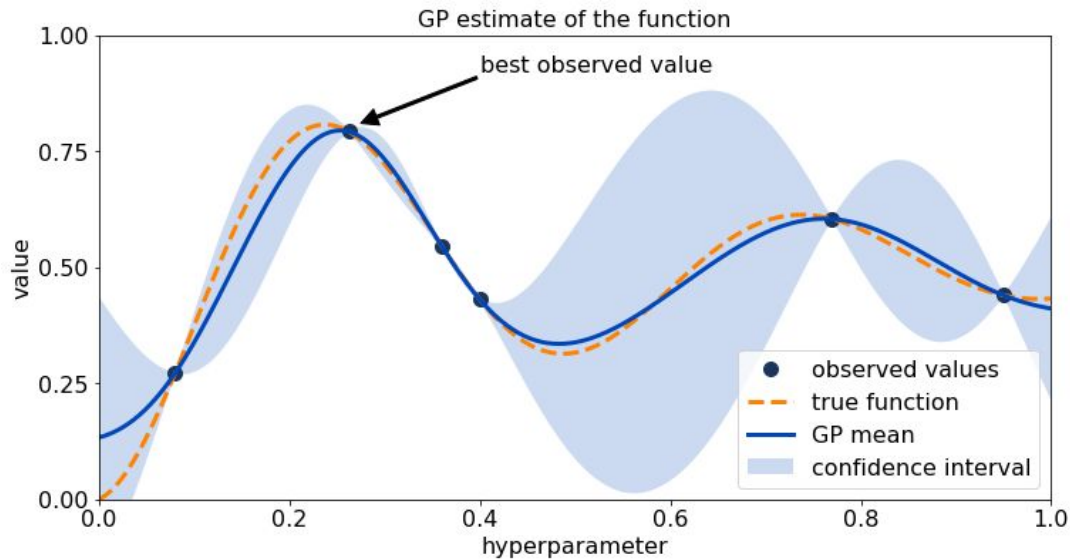
Bayesian Optimization



Bayesian Optimization



Bayesian Optimization



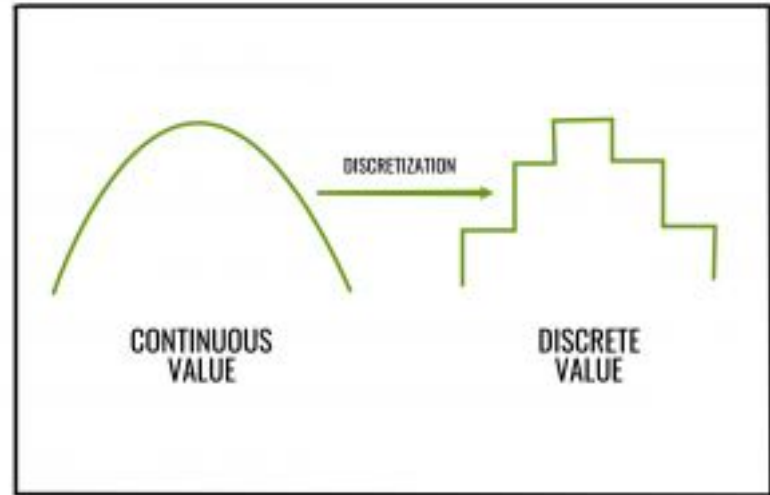


Discretization Algorithms

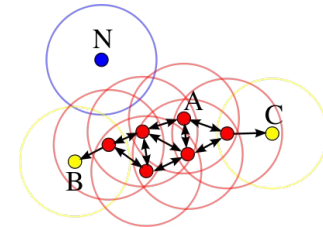
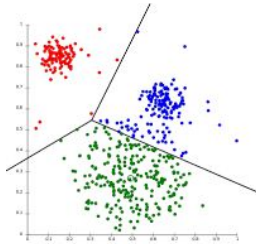
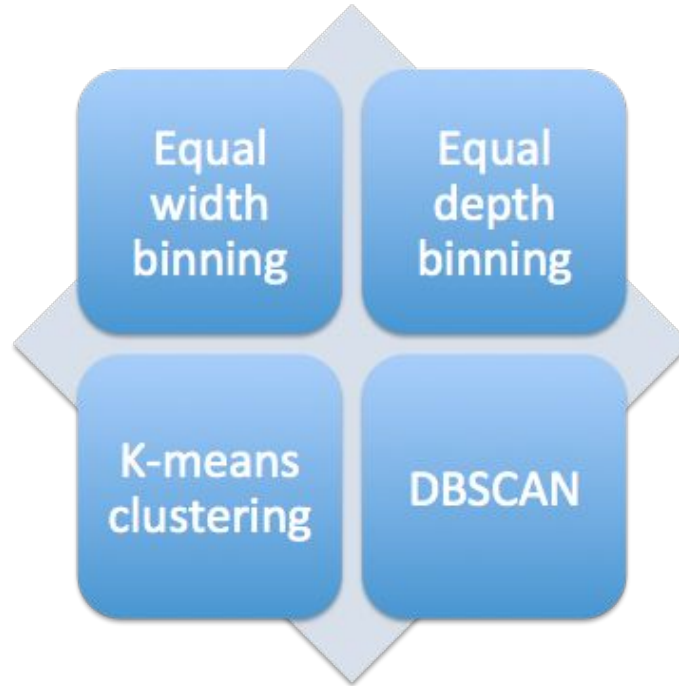
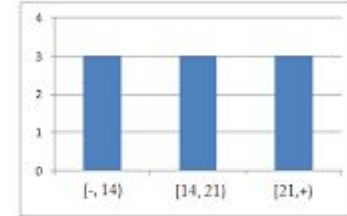
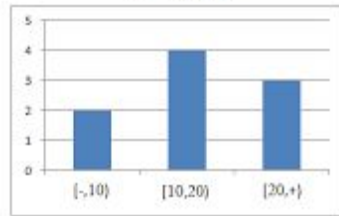


Techniques used

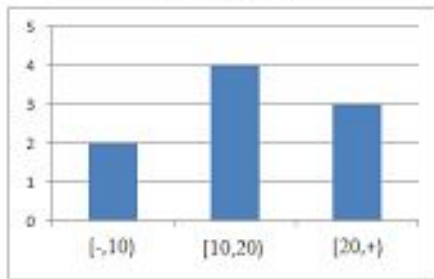
With the exception of decision trees, the rule-mining algorithms cannot be trained on continuous data. Therefore, in order to apply the rule-mining algorithms to datasets, we must rely on data discretization techniques to preprocess the data before mining the rules.



Discretization Techniques

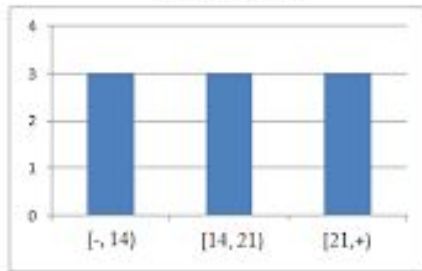


Equal Width Binning



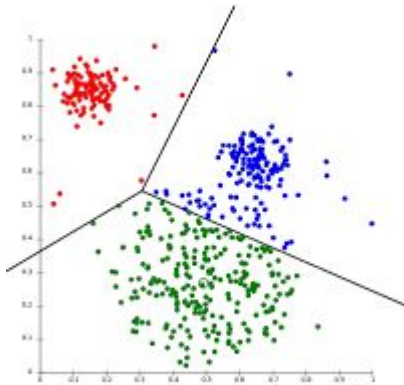
- This algorithm is a comparatively simple binning technique.
- Assuming each cluster, same diameter, each of bins have size $\text{max-min}/K$.
- To discretize, values are assigned to the respective bin they fall into.

Equal Depth Binning



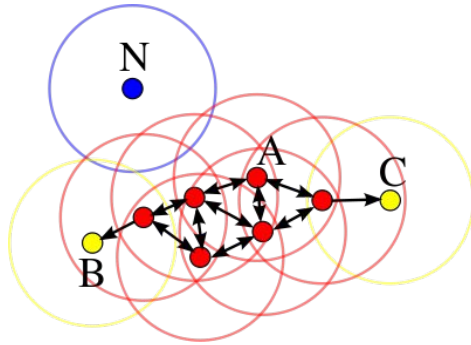
- Equal-depth or equal-frequency binning is another simple discretization approach.
- Each bin approximately holds the same number of instances.
- This is done by sorting the values of the feature.

k-means



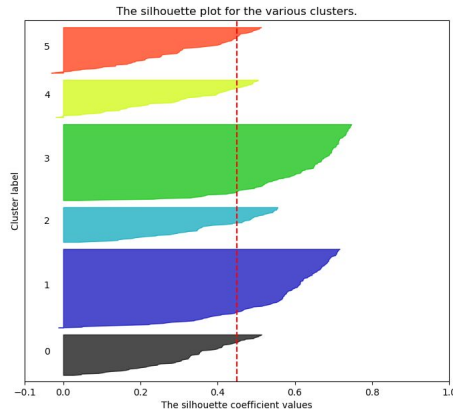
- K-Means is based on the idea of centroids, which are points in the centre of the cluster.
- K centroids are initialized randomly, and the instances are assigned to the cluster whose centroid is closest.
- The algorithm converges when the movement of centroids is below a certain threshold.
- Quite fast, sensitive to outliers

DBSCAN



- DBSCAN considers clusters to be regions of high density.
- For each instance, the algorithm counts the number of instances within a distance ϵ , also called the instance's ϵ -neighbourhood.
- The neighbours of this core instance are considered to be in the same cluster, where some neighbours may also be core instances themselves.
- A cluster consists of a multitude of core instances.

Cluster Optimization: Silhouette score



- The silhouette score has been utilized to provide a metric for accuracy of clusters in this project.
- This score computes the mean silhouette coefficient of all samples.

$$\text{silhouette_score} = \frac{b-a}{\max(a,b)}$$

- Clusters are optimized by exhaustive search this project, i.e., every combination of parameters is tested using the silhouette score, before returning the parameters resulting in the highest score.

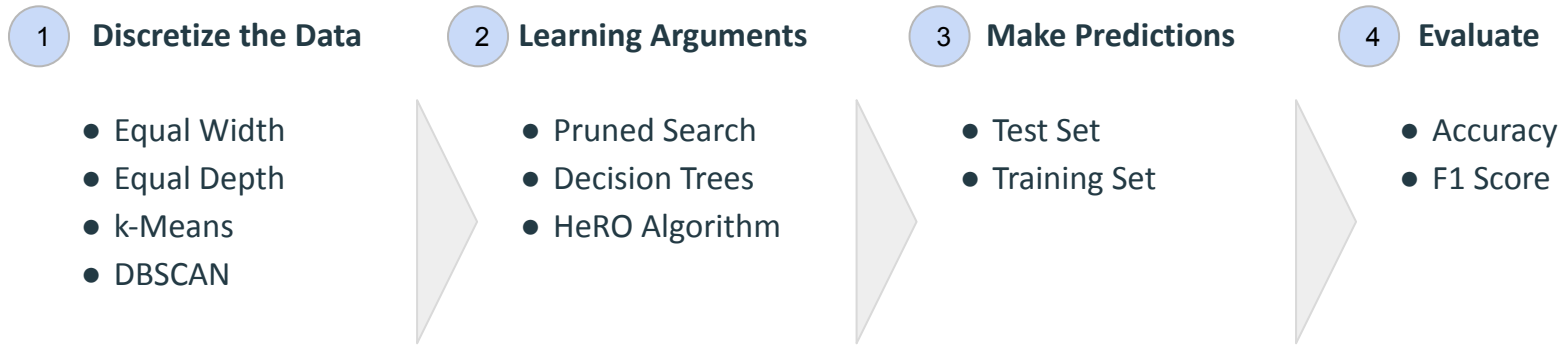


Experimental Results



Experimental Results

Set Up



Hyperparameter

- Algorithm Type
- No. Bins (optional)
- Search Depth
- Max Premises

Experimental Results

Evaluation Pruned Search

<i>n</i> =198	Acc	F1	No. bins	Depth	Runtime	Max premises
Acc	1,000					
F1	0,940	1,000				
No. bins	-0,008	0,057	1,000			
Depth	0,000	0,000	0,000	1,000		
Runtime	-0,173	-0,001	0,170	0,035	1,000	
Max premises	0,000	0,000			0,207	1,000

- Depth and Premise constraint do not affect precision

- Algorithms with higher accuracy have lower the runtime

Evaluation Decision Trees and HeRO

Decision Trees

- Very high accuracy and F1 score for all discretization techniques
- Reasonable runtime in comparison
- Work as well with undiscretized input data

HeRO

- Performance strongly depends on the discretization technique. Accuracy varies from 0.5 to 0.97
- Very high variations in runtime. Some configurations were 30x times slower
- Still has potential for improvements

Discussion

Evaluation

- Accuracy does not draw a complete picture because of explainability also matters
- Number of bins has a tremendous impact on the difficulty of the problem (e.g. no. bins=1)



Future Work: Measure explainability of the algorithm

Runtime

- Runtime increases exponentially in search and discretization algorithms
- Experiments were run on reduced data sets
- Needs improvement before becoming applicable



Future Work: Measure explainability of the algorithm



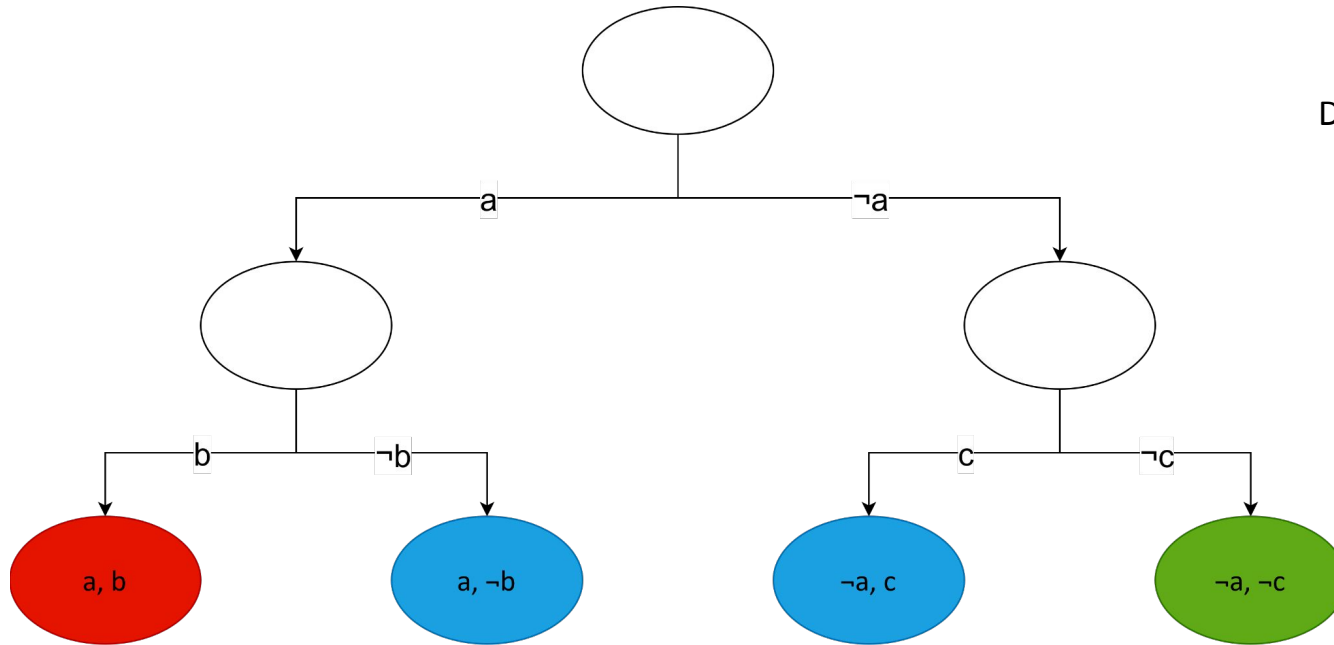
Future Work



Future Work

- Arguments with Exceptions for Decision Tree Rule-Mining
 - Derive arguments with exceptions from decision trees
 - Form arguments with exceptions, which allows prediction on incomplete data
- Optimizing Discretization
 - Choosing the columns to be discretized by selecting columns containing numerical values

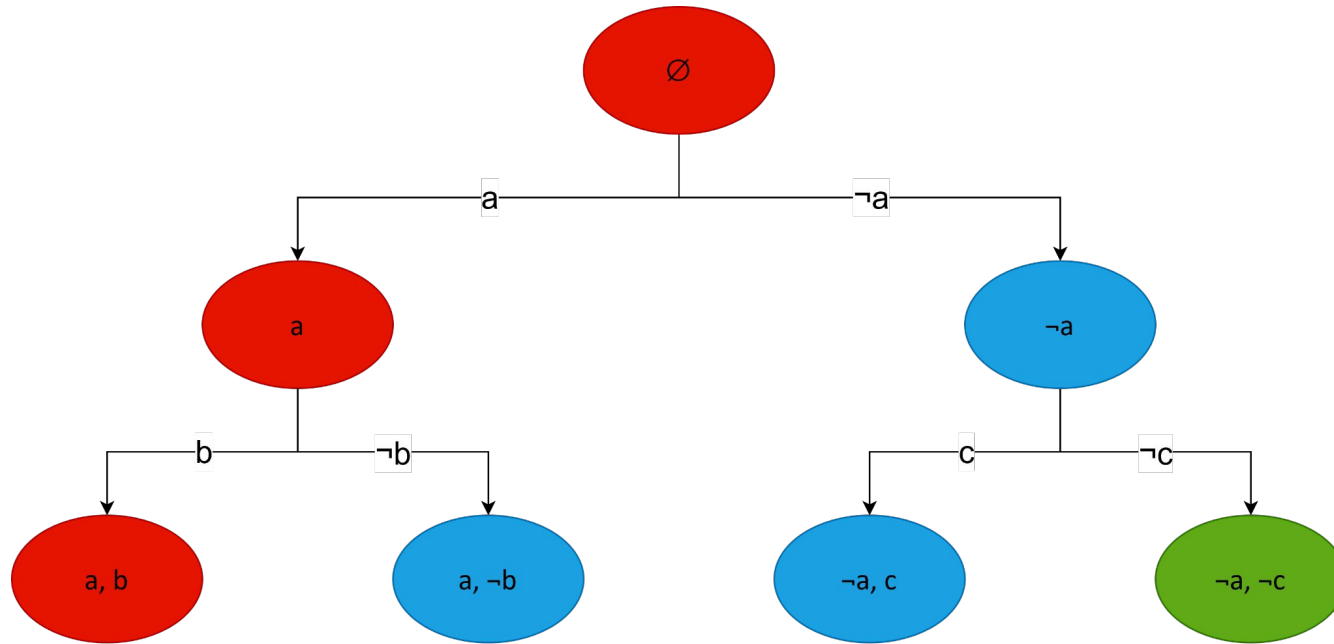
Arguments with exceptions from decision trees



Derived rules:

- $a, b \rightarrow$ red
- $a, \neg b \rightarrow$ blue
- $\neg a, c \rightarrow$ blue
- $\neg a, \neg c \rightarrow$ green

Arguments with exceptions from decision trees



Derived rules:

- $a, b \rightarrow$ red
- $a, \neg b \rightarrow$ blue
- $\neg a, c \rightarrow$ blue
- $\neg a, \neg c \rightarrow$ green

New rules:

- $\emptyset \rightarrow$ red
- $a \rightarrow$ red
- $\neg a \rightarrow$ blue

Arguments with exceptions from decision trees: Some rules imply others!

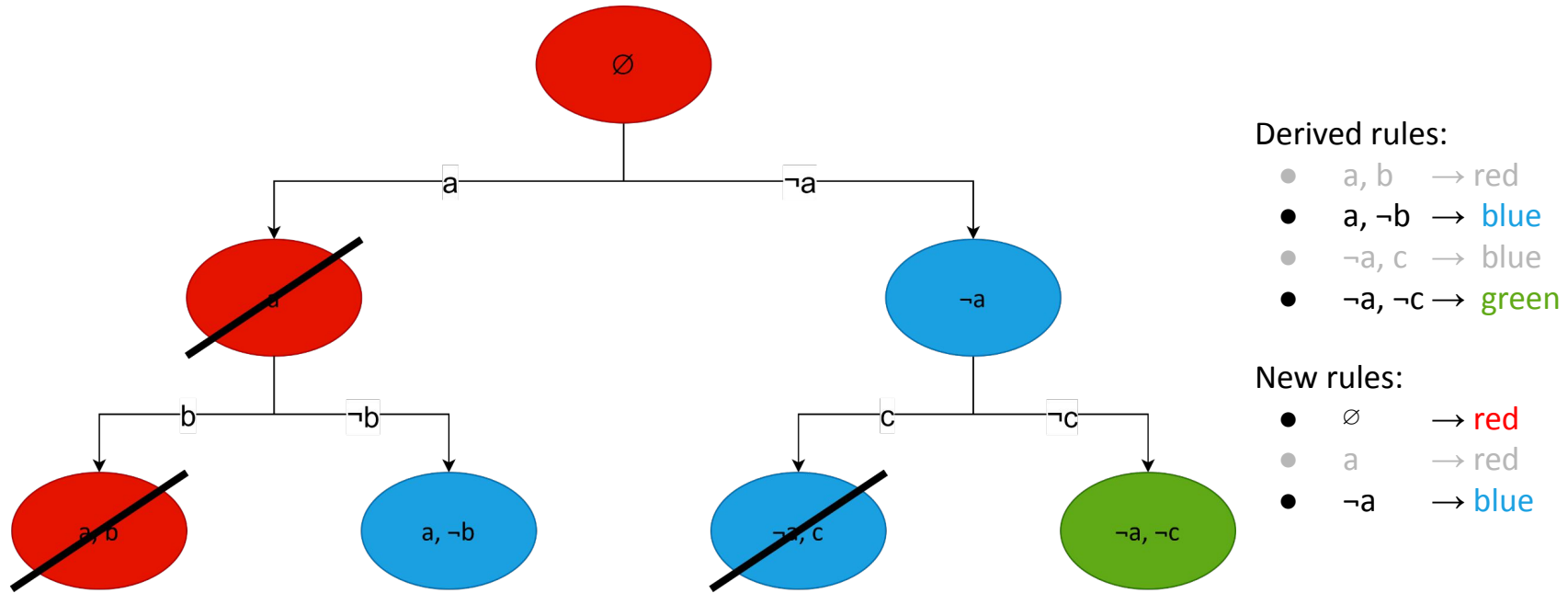
$\emptyset \rightarrow \text{red}$	implies	$a, b \rightarrow \text{red}$
$\neg a \rightarrow \text{blue}$	implies	$\neg a, c \rightarrow \text{blue}$

If rule a implies rule b:

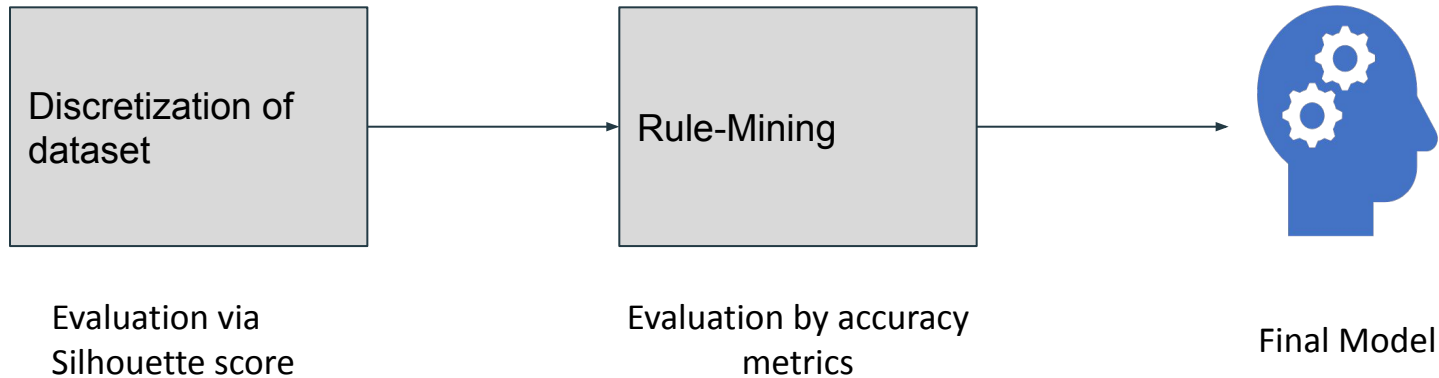
- The premises of rules a are a subset of the premises of rule b
- The conclusions of rule a and b are the same

Since these cases are covered by the less specific rules,
we can prune rules that are implied.

Arguments with exceptions from decision trees



Merging discretization and rule-mining into one classifier



Merging discretization and rule-mining into one classifier



Optimization of all
parameters



Final Model

- Parameters are combined from both techniques
- Can entail techniques used

Possible problem with the merge

- Target feature is not taken into account for discretization
- Instead, discretization will aim to maximize clustering metrics by design
- Improvement may be negligible



Conclusion



Project Conclusion

1. Examples from Verheij are implemented and replicated
2. Four different algorithms for learning arguments have been investigated:
 - Naive Search
 - Pruned Search
 - HerO
 - Decision Tree
3. The approach has been transferred to an attribute-value dataset:
 - Equal-Width Binning
 - Equal-Depth Binning
 - K-Means Clustering
 - DBSCAN Clustering

Project Conclusion

4. We have given a survey over existing techniques for mining rules and arguments. We have applied ideas from logical learning and from the Apriori algorithm. We have implemented the HeRO algorithm. We also sketched an algorithm that allows learning arguments from decision trees
5. Applicability of the algorithms has been shown by experimentation



Thank you!
Questions?





Bart Verheij. 'Proof with and without Probabilities'. *Artificial Intelligence and Law* 25, no. 1 (2017): 127–54.

Bart Verheij. 'Arguments for Good Artificial Intelligence'. Inaugural lecture. Groningen: University of Groningen, 2018.

Benjamin Johnston and Guido Governatori. 'An Algorithm for the Induction of Defeasible Logic Theories from Databases'. In *Proceedings of the 14th Australasian Database Conference-Volume 17*, 75–83, 2003.



Appendix



What specific algorithms are used in Bayesian Optimization?

- To create the prior: Kriging/Gaussian process regression
- For the acquisition function: Probabilistically, one of the following functions are chosen:
 - Lower confidence bound
 - Negative expected improvement
 - Negative probability of improvement
- An unknown numerical approach is used to find the optimum of the prior (not in documentation)

		<i>Verheij 2017</i>	<i>Naive search</i>	<i>Pruned search</i>	<i>HeRO</i>
1	inn, \neg gui	inn \Leftarrow \neg gui \leftarrow inn	inn $\wedge \neg$ gui \Leftarrow	inn $\wedge \neg$ gui \Leftarrow \neg gui \leftarrow inn gui $\leftarrow \neg$ inn evi \wedge gui $\leftarrow \neg$ inn evi $\wedge \neg$ inn \leftarrow gui gui $\wedge \neg$ inn \leftarrow evi	inn $\wedge \neg$ gui \wedge evi \Leftarrow
0	\neg inn, gui, evi				
(a) Case model		gui \leftarrow evi			gui $\wedge \neg$ inn \Leftarrow evi
		(b) Learned arguments			

Figure 5.1: Learning arguments in case model 1 from Verheij (2017): *Presumption of innocence*.

2	a	b	c	y
1	a	b	$\neg c$	$\neg y$
0	a	$\neg b$	$\neg c$	y

(a) Case model

<i>Manual</i>	<i>Naive search</i>	<i>Pruned search</i>	<i>HeRO</i>
$y \Leftarrow$	$y \Leftarrow$	$y \Leftarrow$	$y \Leftarrow$
$\neg y \Leftarrow \neg c$		$y \leftarrow c$ $\neg y \Leftarrow \neg c$	
		$\neg y \leftarrow b \wedge \neg c$	$\neg y \Leftarrow b \wedge \neg c$
		$\neg y \Leftarrow a \wedge \neg c$	
$y \leftarrow \neg b$		$y \leftarrow \neg b$	

(b) Learned arguments

Figure 5.2: Learning arguments in case model 1 from Verheij (2017): *Presumption of innocence*.