

# Tables

Event Log	Accuracy		
	Base	Modified	Enriched
−age, −gender	.916 ± .002	.915 ± .002	.915 ± .002
−age, +gender	.916 ± .002	.920 ± .004	.926 ± .003
+age, −gender	.916 ± .002	.920 ± .003	.928 ± .004
+age, +gender	.916 ± .002	.923 ± .003	.932 ± .003

Event Log	$\Delta DP$ (Age)			$\Delta DP$ (Gender)		
	Base	Modified	Enriched	Base	Modified	Enriched
−age, −gender	.014 ± .012	.014 ± .011	.014 ± .011	.009 ± .003	.021 ± .018	.021 ± .018
−age, +gender	.010 ± .005	.008 ± .001	.042 ± .028	.008 ± .005	.064 ± .041	.837 ± .081
+age, −gender	.014 ± .015	.030 ± .058	.903 ± .033	.037 ± .026	.033 ± .031	.052 ± .033
+age, +gender	.018 ± .013	.005 ± .003	.917 ± .017	.007 ± .006	.004 ± .002	.836 ± .044

**Tab. 1.1:** Evaluation of accuracy and  $\Delta DP$  for the four versions of the *Hospital Billing* log. The attributes *age* and *gender* are annotated based on whether they introduce a bias (+) or not (−). The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .

Event Log	Accuracy			$\Delta DP$		
	Base	Modified	Enriched	Base	Modified	Enriched
cs	.816 $\pm$ .005	.848 $\pm$ .003	.890 $\pm$ .001	.004 $\pm$ .004	.002 $\pm$ .003	.996 $\pm$ .002
hb (–age –gender)	.916 $\pm$ .002	.915 $\pm$ .002	.915 $\pm$ .002	.011 $\pm$ .004	.017 $\pm$ .005	.017 $\pm$ .005
hb (–age +gender)	.916 $\pm$ .002	.920 $\pm$ .004	.926 $\pm$ .003	.009 $\pm$ .001	.036 $\pm$ .039	.439 $\pm$ .563
hb (+age –gender)	.916 $\pm$ .002	.920 $\pm$ .003	.928 $\pm$ .004	.026 $\pm$ .016	.032 $\pm$ .002	.478 $\pm$ .602
hb (+age +gender)	.916 $\pm$ .002	.923 $\pm$ .003	.932 $\pm$ .003	.012 $\pm$ .008	.005 $\pm$ .001	.877 $\pm$ .057
bpi	.817 $\pm$ .003	.822 $\pm$ .001	.829 $\pm$ .004	.058 $\pm$ .004	.053 $\pm$ .006	.510 $\pm$ .094

**Tab. 1.2:** Evaluation of accuracy and  $\Delta DP$  for the *Cancer Screening (cs)* log, the *BPI Challenge 2012 (bpi)* log, and four versions of the *Hospital Billing (hb)* log, where the attributes *age* and *gender* are annotated based on whether they introduce a bias (+) or not (–). Since the *hb* event log uses two sensitive attributes, we report their average  $\Delta DP$ . The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across five validation folds, expressed as  $\mu \pm \sigma$ .

Event Log	Number of Nodes	Number of removed Nodes	Depth
cs	26.6 $\pm$ 4.7	1 $\pm$ 0	10.8 $\pm$ 0.4
hb (–age, –gender)	33.4 $\pm$ 2.6	0 $\pm$ 0	11.0 $\pm$ 0.7
hb (–age, +gender)	45.4 $\pm$ 3.2	1 $\pm$ 0	12.2 $\pm$ 0.8
hb (+age, –gender)	42.2 $\pm$ 4.3	1 $\pm$ 0	12.8 $\pm$ 1.0
hb (+age, +gender)	47.0 $\pm$ 3.4	2 $\pm$ 0	12.8 $\pm$ 0.4
bpi	38.2 $\pm$ 3.0	1 $\pm$ 0	11.2 $\pm$ 0.4

**Tab. 1.3:** Evaluation of the characteristics of the distilled decision tree for the *Cancer Screening (cs)* log, the *BPI Challenge 2012 (bpi)* log and all versions of the *Hospital Billing (hb)* log. For *hb*, the attributes *age* and *gender* are annotated based on whether they introduce bias (+) or not (–). The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .

Bias Strength	Accuracy			$\Delta DP$		
	Base	Modified	Enriched	Base	Modified	Enriched
0.50	.797 $\pm$ .004	.851 $\pm$ .003	.851 $\pm$ .003	.005 $\pm$ .008	.001 $\pm$ .001	.001 $\pm$ .001
0.55	.805 $\pm$ .003	.848 $\pm$ .003	.861 $\pm$ .006	.001 $\pm$ .002	.002 $\pm$ .002	.800 $\pm$ .444
0.60	.808 $\pm$ .003	.849 $\pm$ .003	.869 $\pm$ .001	.000 $\pm$ .001	.003 $\pm$ .003	.996 $\pm$ .002
0.65	.813 $\pm$ .003	.851 $\pm$ .003	.882 $\pm$ .002	.002 $\pm$ .002	.001 $\pm$ .001	.973 $\pm$ .055
0.70	.817 $\pm$ .003	.850 $\pm$ .003	.892 $\pm$ .003	.002 $\pm$ .002	.000 $\pm$ .000	.999 $\pm$ .001
0.75	.825 $\pm$ .003	.851 $\pm$ .004	.904 $\pm$ .002	.002 $\pm$ .002	.000 $\pm$ .001	.999 $\pm$ .001
0.80	.825 $\pm$ .005	.848 $\pm$ .007	.912 $\pm$ .004	.001 $\pm$ .001	.001 $\pm$ .001	.998 $\pm$ .001
0.85	.838 $\pm$ .002	.853 $\pm$ .002	.925 $\pm$ .003	.001 $\pm$ .001	.001 $\pm$ .001	.999 $\pm$ .002
0.90	.841 $\pm$ .002	.852 $\pm$ .004	.936 $\pm$ .002	.001 $\pm$ .002	.000 $\pm$ .001	.999 $\pm$ .001
0.95	.844 $\pm$ .004	.850 $\pm$ .005	.945 $\pm$ .003	.004 $\pm$ .004	.000 $\pm$ .001	.000 $\pm$ .000
1.00	.854 $\pm$ .003	.854 $\pm$ .003	.958 $\pm$ .001	.004 $\pm$ .002	.001 $\pm$ .001	.000 $\pm$ .001

Bias Strength	Number of Nodes	Number of removed Nodes	Depth
0.50	27.0 $\pm$ 4.7	0.0 $\pm$ 0.0	10.8 $\pm$ 0.4
0.55	25.0 $\pm$ 1.4	0.8 $\pm$ 0.4	10.6 $\pm$ 0.5
0.60	24.6 $\pm$ 1.7	1.0 $\pm$ 0.0	10.4 $\pm$ 0.5
0.65	27.8 $\pm$ 4.4	1.0 $\pm$ 0.0	10.2 $\pm$ 0.4
0.70	27.0 $\pm$ 4.7	1.0 $\pm$ 0.0	10.2 $\pm$ 0.4
0.75	26.2 $\pm$ 5.0	1.0 $\pm$ 0.0	10.6 $\pm$ 0.5
0.80	29.4 $\pm$ 5.9	1.0 $\pm$ 0.0	11.0 $\pm$ 0.0
0.85	25.4 $\pm$ 1.7	1.0 $\pm$ 0.0	10.8 $\pm$ 0.4
0.90	25.8 $\pm$ 1.1	1.0 $\pm$ 0.0	10.8 $\pm$ 0.4
0.95	27.8 $\pm$ 5.4	1.0 $\pm$ 0.0	10.8 $\pm$ 0.4
1.00	23.0 $\pm$ 0.0	1.0 $\pm$ 0.0	10.0 $\pm$ 0.0

**Tab. 1.4:** Evaluation of the accuracy,  $\Delta DP$  and characteristics of the distilled decision tree for varying bias strengths. The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .

Num. Attributes	Accuracy			$\Delta DP$		
	Base	Modified	Enriched	Base	Modified	Enriched
1	.807 $\pm$ .003	.820 $\pm$ .003	.862 $\pm$ .003	.000 $\pm$ .000	.000 $\pm$ .000	.999 $\pm$ .001
2	.796 $\pm$ .002	.807 $\pm$ .003	.830 $\pm$ .002	.001 $\pm$ .001	.045 $\pm$ .099	.501 $\pm$ .009
4	.799 $\pm$ .001	.807 $\pm$ .002	.822 $\pm$ .004	.001 $\pm$ .001	.005 $\pm$ .010	.374 $\pm$ .005
6	.799 $\pm$ .005	.805 $\pm$ .005	.816 $\pm$ .005	.002 $\pm$ .002	.028 $\pm$ .014	.267 $\pm$ .077
8	.799 $\pm$ .004	.801 $\pm$ .001	.812 $\pm$ .002	.001 $\pm$ .001	.020 $\pm$ .013	.238 $\pm$ .029
10	.801 $\pm$ .004	.803 $\pm$ .003	.811 $\pm$ .005	.004 $\pm$ .007	.033 $\pm$ .021	.168 $\pm$ .094

Num. Attributes	Number of Nodes	Number of removed Nodes	Depth
1	24.6 $\pm$ 2.6	1.0 $\pm$ 0.0	10.4 $\pm$ 0.9
2	31.0 $\pm$ 4.2	1.8 $\pm$ 0.4	10.4 $\pm$ 0.9
4	47.8 $\pm$ 4.1	7.4 $\pm$ 0.9	10.8 $\pm$ 0.4
6	81.8 $\pm$ 19.7	21.0 $\pm$ 6.0	11.6 $\pm$ 0.9
8	77.8 $\pm$ 15.7	22.4 $\pm$ 5.2	10.8 $\pm$ 0.4
10	60.2 $\pm$ 16.0	13.2 $\pm$ 8.8	10.8 $\pm$ 0.4

**Tab. 1.5:** Evaluation of the accuracy,  $\Delta DP$  and characteristics of the distilled decision tree for a varying amount of sensitive attributes. The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .

Num. Decisions	Accuracy			$\Delta DP$		
	Base	Modified	Enriched	Base	Modified	Enriched
2	.765 $\pm$ .001	.793 $\pm$ .002	.843 $\pm$ .002	.003 $\pm$ .003	.001 $\pm$ .001	.998 $\pm$ .003
4	.684 $\pm$ .006	.710 $\pm$ .019	.791 $\pm$ .003	.000 $\pm$ .001	.007 $\pm$ .002	.997 $\pm$ .002
8	.627 $\pm$ .004	.648 $\pm$ .009	.749 $\pm$ .003	.003 $\pm$ .003	.009 $\pm$ .009	.983 $\pm$ .021
12	.595 $\pm$ .003	.648 $\pm$ .001	.729 $\pm$ .003	.003 $\pm$ .002	.005 $\pm$ .004	.968 $\pm$ .019
16	.582 $\pm$ .004	.638 $\pm$ .004	.716 $\pm$ .002	.004 $\pm$ .003	.002 $\pm$ .001	.939 $\pm$ .028
20	.570 $\pm$ .002	.625 $\pm$ .005	.702 $\pm$ .006	.003 $\pm$ .002	.004 $\pm$ .001	.921 $\pm$ .035

Num. Decisions	Number of Nodes	Number of removed Nodes	Depth
2	19.8 $\pm$ 2.3	1.0 $\pm$ 0.0	6.8 $\pm$ 0.4
4	33.0 $\pm$ 3.5	2.0 $\pm$ 0.0	10.4 $\pm$ 0.5
8	61.8 $\pm$ 3.6	5.8 $\pm$ 1.1	16.4 $\pm$ 0.9
12	90.2 $\pm$ 1.8	9.6 $\pm$ 0.9	25.0 $\pm$ 0.0
16	124.6 $\pm$ 0.9	14.6 $\pm$ 0.9	33.0 $\pm$ 0.0
20	154.2 $\pm$ 1.8	17.6 $\pm$ 0.9	41.0 $\pm$ 0.0

**Tab. 1.6:** Evaluation of the accuracy,  $\Delta DP$  and characteristics of the distilled decision tree for a varying amount of biased decisions. The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .

Event Log	Accuracy (NN)	Accuracy (DT)
cs	.848 $\pm$ .003	.847 $\pm$ .003
hb (−age, −gender)	.915 $\pm$ .002	.913 $\pm$ .002
hb (−age, +gender)	.920 $\pm$ .004	.918 $\pm$ .004
hb (+age, −gender)	.920 $\pm$ .003	.917 $\pm$ .003
hb (+age, +gender)	.923 $\pm$ .003	.919 $\pm$ .002
bpi	.822 $\pm$ .001	.823 $\pm$ .002

**Tab. 1.7:** Comparison of the accuracy of the modified neural network (NN) against the accuracy of the modified decision tree (DT) for the *Cancer Screening* (cs) log, the *BPI Challenge 2011* (bpi) log and all versions of the *Hospital Billing* (hb) log. For hb, the attributes age and gender are annotated based on whether they introduce bias (+) or not (−). The reported values represent the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) across validation folds, expressed as  $\mu \pm \sigma$ .