M1 – Apprentissage

Michèle Sebag – Benoit Barbot LRI – LSV

13 janvier 2014

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>

Empirical assessment

Motivations

Caruana et al, 08 Crossing the chasm

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Which algorithm for which data ?
- Which hyper-parameters ?

Note on meta-learning

- Programming by optimization, http://www.prog-by-opt.net/
- Needed: descriptive features

Empirical assessment, 2

Empirical study in large scale

- 700,000 dimensions
 curse of dimensionality
- Scores: Accuracy, Area under the ROC curve, square loss

Algorithms

 SVM - adjusting C; stochastic gradient descent al. 05

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Perceptrons
- Logistic regression
- Artificial neural nets
- Naive Bayes
- k-Nearest Neighbors
- Boosted decision trees
- Random Forests

Empirical assessment, 3

Methodology

- Normalize data
- Normalize indicators (median)

Datasets

		1	1		
Problem	n Attr	Train	Valid	Test	% Pos
Sturn	761	10K	2K	9K	33.65
Calam	761	10K	2K	9K	34.32
Digits	780	48K	12K	10K	49.01
Tis	927	5.2K	1.3K	6.9K	25.13
Cryst	1344	2.2K	$1.1 \mathrm{K}$	2.2K	45.61
KDD98	3848	76.3K	19K	96.3K	5.02
R-S	20958	35K	7K	30.3K	30.82
Cite	105354	81.5K	18.4K	81.5K	0.17
Dse	195203	120K	43.2K	107K	5.46
Spam	405333	36K	9K	42.7K	44.84
Imdb	685569	84K	$18.4 \mathrm{K}$	84K	0.44

Table 1. Description of problems

Datasets

TIS1 is from the Kent Ridge Bio-medical Data Repository. The problem is to find Translation Initiation Sites (TIS) at which translation from mRNA to proteins initiates.

CRYST2 is a protein crystallography diffraction pattern analysis dataset from the X6A beamline at Brookhaven National Laboratory.

STURN and CALAM are ornithology datasets. The task is to predict the appearance of two bird species: sturnella neglecta and calamospiza melanocorys.

KDD98 is from the 1998 KDD-Cup. The task is to predict if a person donates money. This is the only dataset with missing values.

DIGITS4 is the MNIST database of handwritten digits by Cortes and

LeCun. It was converted from a 10 class problem to a hard binary problem by treating digits less than 5 as one class and the rest as the other class.

Datasets, 2

IMDB and CITE are link prediction datasets. For IMDB each attribute represents an actor, director, etc. For CITE attributes are the authors of a paper in the CiteSeer digital library. For IMDB the task is to predict if Mel Blanc was involved in the film or television program and for CITE the task is to predict if J. Lee was a coauthor of the paper. We created SPAM from the TREC 2005 Spam Public Corpora. Features take binary values showing if a word appears in the document or not. Words that appear less than three times in the whole corpus were removed.

Real-Sim (R-S) is a compilation of Usenet articles from four discussion groups: simulated auto racing, simulated aviation, real autos and real aviation. The task is to distinguish real from simulated.

DSE7 is newswire text with annotated opinion expressions. The task is to find Subjective Expressions i.e. if a particular word expresses an opinion.

Results

DIM	761	761	780	927	1344	3448	20958	105354	195203	405333	685569	-
ACC	STURN	CALAM	DIGITS	Tis	CRYST	KDD98	R-S	CITE	DSE	SPAM	IMDB	MEAN
MEDIAN	0.6901	0.7337	0.9681	0.9135	0.8820	0.9494	0.9599	0.9984	0.9585	0.9757	0.9980	
BSTDT	0.9962	1.0353	1.0120	0.9993	1.0178	0.9998	0.9904	1.0000	0.9987	0.9992	1.0000	1.0044
RF	0.9943	1.0103	1.0076	1.0025	1.0162	1.0000	0.9995	0.9998	1.0013	1.0044	1.0000	1.0033
SVM	1.0044	1.0018	1.0024	1.0060	1.0028	0.9999	1.0156	1.0008	1.0004	1.0008	1.0003	1.0032
BAGDT	1.0001	1.0350	0.9976	1.0017	1.0111	1.0000	0.9827	1.0000	0.9996	0.9959	1.0000	1.0021
ANN	0.9999	0.9899	1.0051	1.0007	0.9869	1.0000	1.0109	1.0001	1.0018	1.0029	1.0003	0.9999
LR	1.0012	0.9896	0.8982	1.0108	1.0080	1.0000	1.0141	1.0001	1.0014	1.0026	0.9999	0.9932
BSTST	1.0077	1.0298	0.9017	0.9815	0.9930	1.0000	0.9925	0.9999	0.9948	0.9905	0.9989	0.9900
KNN	1.0139	0.9982	1.0122	0.9557	0.9972	0.9999	0.9224	1.0000	0.9987	0.9698	0.9996	0.9880
PRC	0.9972	0.9864	0.9010	0.9735	0.9930	1.0000	1.0119	0.9999	1.0007	1.0041	1.0001	0.9880
NB	0.9695	0.9347	0.8159	0.9230	0.9724	1.0000	1.0005	1.0000	0.9878	0.9509	0.9976	0.9593
RMS	STURN	CALAM	DIGITS	Tis	Cryst	KDD98	R-S	CITE	DSE	SPAM	IMDB	MEAN
MEDIAN	0.5472	0.5800	0.8449	0.7455	0.7051	0.7813	0.8257	0.9623	0.8154	0.8645	0.9597	
RF	0.9980	1.0209	1.0186	1.0102	1.0277	1.0003	1.0011	0.9988	1.0072	1.0118	1.0006	1.0087
BSTDT	0.9993	1.0351	1.0363	0.9977	1.0323	0.9998	0.9781	1.0003	0.9983	1.0007	1.0003	1.0071
ANN	1.0042	0.9987	1.0088	1.0109	1.0014	1.0005	1.0315	1.0011	1.0068	1.0077	1.0022	1.0067
SVM	0.9979	0.9882	1.0076	1.0149	0.9972	0.9992	1.0409	1.0091	1.0067	0.9993	1.0004	1.0056
BAGDT	1.0007	1.0357	0.9924	1.0023	1.0218	0.9998	0.9587	1.0000	0.9994	0.9782	1.0012	0.9991
LR	1.0010	0.9963	0.8169	1.0232	0.9935	1.0007	1.0367	1.0009	1.0082	1.0073	0.9988	0.9894
PRC	0.9976	0.9841	0.8115	0.9537	0.9919	0.9998	1.0313	0.9979	1.0006	1.0071	0.9997	0.9796
BSTST	1.0078	1.0205	0.8202	0.9757	1.0021	1.0007	0.9861	1.0000	0.9900	0.9695	0.9952	0.9789
KNN	1.0119	1.0013	1.0365	0.9309	0.9986	1.0000	0.8468	0.9988	0.9983	0.9270	0.9941	0.9768
NB	0.9793	0.9509	0.7236	0.9031	0.9454	1.0000	0.9989	0.9981	0.9828	0.8984	0.9731	0.9412
AUC	STURN	CALAM	DIGITS	Tis	CRYST	KDD98	R-S	CITE	DSE	SPAM	IMDB	MEAN
MEDIAN	0.6700	0.7793	0.9945	0.9569	0.9490	0.5905	0.9913	0.7549	0.9008	0.9957	0.9654	
RF	0.9892	1.0297	1.0017	1.0069	1.0134	1.0140	1.0009	1.0962	1.0304	1.0022	1.0209	1.0187
KNN	1.0397	0.9992	1.0024	0.9509	1.0007	1.0165	0.9905	1.1581	1.0027	0.9902	0.9648	1.0105
LR	1.0045	0.9903	0.9424	1.0136	1.0070	1.0492	1.0041	1.0272	1.0293	0.9999	1.0084	1.0069
ANN	1.0132	1.0008	1.0001	1.0042	0.9992	1.0461	1.0031	0.9779	1.0105	1.0001	1.0021	1.0052
BSTST	1.0199	1.0304	0.9468	0.9901	0.9993	1.0512	0.9991	0.9956	0.9973	0.9989	1.0036	1.0029
SVM	0.9870	0.9645	1.0002	1.0077	0.9909	0.9324	1.0032	1.1120	1.0100	1.0011	0.9979	1.0006
BSTDT	0.9991	1.0492	1.0033	0.9958	1.0137	0.9605	0.9962	0.9646	0.9881	1.0015	1.0041	0.9978
BAGDT	1.0009	1.0551	0.9999	1.0062	1.0116	0.9768	0.9890	0.9673	0.9691	0.9925	0.9809	0.9954
PRC	0.9973	0.9630	0.9372	0.9749	0.9937	0.9724	1.0036	0.9991	0.9777	1.0006	0.9477	0.9788
NB	0.9329	0.8936	0.8574	0.9407	0.9574	0.9860	0.9990	1.0009	0.9917	0.9798	0.8787	0.9471
AVG	STURN	CALAM	DIGITS	Tis	CRYST	KDD98	R-S	Cite	DSE	SPAM	IMDB	MEAN
RF	0.9938	1.0203	1.0093	1.0065	1.0191	1.0048	1.0005	1.0316	1.0130	1.0061	1.0072	1.0102
ANN	1.0058	0.9965	1.0047	1.0053	0.9958	1.0156	1.0152	0.9930	1.0064	1.0036	1.0015	1.0039
BSTDT	0.9982	1.0399	1.0172	0.9976	1.0212	0.9867	0.9882	0.9883	0.9950	1.0004	1.0014	1.0031
SVM	0.9965	0.9848	1.0034	1.0095	0.9970	0.9772	1.0199	1.0406	1.0057	1.0004	0.9995	1.0031
BAGDT	1.0006	1.0419	0.9966	1.0034	1.0148	0.9922	0.9768	0.9891	0.9894	0.9889	0.9940	0.9989
LR	1.0022	0.9921	0.8858	1.0159	1.0028	1.0166	1.0183	1.0094	1.0129	1.0033	1.0024	0.9965
KNN	1.0219	0.9996	1.0170	0.9458	0.9988	1.0055	0.9199	1.0523	0.9999	0.9623	0.9862	0.9917
BSIST	1.0118	1.0269	0.8896	0.9824	0.9982	1.0173	0.9926	0.9985	0.9941	0.9863	0.9992	0.9906
PRC	0.9974	0.9778	0.8832	0.9674	0.9929	0.9907	1.0156	0.9990	0.9930	1.0039	0.9825	0.9821
NB	0.9606	0.9264	0.7989	-0.9223	0.9584	0.9953	0.9995	0.9997	0.9874	0.9430	0.9498	0.9492

Table 2. Standardized scores of each learning algorithm

200

Results w.r.t. dimension



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

3

Moving average standardized scores of each learning algorithm as a function of the dimension.

The rank for the algorithms to perform consistently well: (1) random forest (2) neural nets (3) boosted tree (4) SVMs

Overview

Introduction

Boosting

PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Example (AT&T)

Schapire 09

- Categorizing customer's queries (Collect, CallingCard, PersonToPerson,..)
- Example queries:
 - yes I'd like to place a collect call long distance please Collect
 - operator I need to make a call but I need to bill it to my office ThirdNumber
 - yes I'd like to place a call on my master card please

CallingCard

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

 I just called a number in sioux city and I musta rang the wrong number because I got the wrong party and I would like to have that taken off of my bill
 BillingCredit

Example (AT&T)

Schapire 09

- Categorizing customer's queries (Collect, CallingCard, PersonToPerson,..)
- Example queries:
 - yes I'd like to place a collect call long distance please Collect
 - operator I need to make a call but I need to bill it to my office ThirdNumber
 - yes I'd like to place a call on my master card please

CallingCard

 I just called a number in sioux city and I musta rang the wrong number because I got the wrong party and I would like to have that taken off of my bill
 BillingCredit

Remark

- Easy to find rules of thumb with good accuracy IF 'card', THEN CallingCard
- Hard to find a single good rule

Procedure

- ▶ A learner (fast, reasonable accuracy, i.e. accuracy > random)
- Learn from (a subset of) training set
- Find a hypothesis
- Do this a zillion times (T rounds)
- Aggregate the hypotheses

Critical issues

- Enforce the diversity of hypotheses
- How to aggregate hypotheses

The Wisdom of Crowds: Why the Many Are Smarter Than the Few and How Collective Wisdom Shapes Business, Economies, Societies and Nations. J. Surowiecki, 2004.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

$$\mathcal{E} = \{(x_i, y_i)\}, \ x_i \in X, \ y_i \in \{-1, 1\}\} \quad (x_i, y_i) \sim D(x, y)$$

Loop

For $t = 1 \dots T$, learn h_t from \mathcal{E}_t

Result: $H = sign(\sum_t \alpha_t h_t)$

Requisite: Classifiers h_t must be diverse

Enforcing diversity through:

- Using different training sets \mathcal{E}_t
- Using diverse feature sets
- Enforce h_t decorrelation

bagging, boosting bagging boosting

Diversity of h_t , Stability of H

Stability: slight changes of \mathcal{E} hardly modifies h

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Stable learners

k-nearest neighbors Linear discriminant analysis (LDA)

Unstable learners

Neural nets Decision trees

Instability: Why is it useful? Bias of \mathcal{H} : Err(h*)

 $h^* = \arg \min\{Err(h), h \in \mathcal{H}\}$

Decreases with size/complexity of \mathcal{H} LDA poor...

Variance:
$$h_1, \dots h_T$$

Variance $(H) = \frac{1}{T} \sum_t ||H - h_t||^2$

Instable learners

- Large variance
- Small bias (e.g. decision trees and NN are universal approximators)
- Variance(Ensemble) decreases with size of ensemble if ensemble elements are not correlated.

Variance
$$(H) \approx \frac{1}{T} Variance(h_t)$$

Why does it work, basics

- Suppose there are 25 base classifiers
- Each one with error rate $\epsilon = .35$
- Assume classifiers are independent

Then, probability that the ensemble classifier makes a wrong prediction:

Pr (ensemble makes error)
$$=\sum_{i=13}^2 5C_{25}^i\epsilon^i(1-\epsilon)^{25-i}=.06$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Overview

Introduction

Boosting PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

PAC Learning: Probably Approximately Correct



Valiant 84

(日)、

Turing award, 2011.

PAC Learning: Probably Approximately Correct Setting

iid samples drawn after distribution $D(\mathbf{x}, y)$

 $\mathcal{E} = \{(x_i, y_i)\}, \ x_i \in \ X, \ y_i \ \in \ \{-1, 1\}\} \ (x_i, y_i) \sim D(x, y)$

Strong learnability

Language (= set of target concepts) C is PAC learnable if there exists algorithm A s.t.

- For all D(x, y)
- For all $0 < \delta < 1$, with probability 1δ probably
- For all error rate $\varepsilon > 0$, approximately correct
- There exists a number *n* of samples, $n = Polynom(\frac{1}{\delta}, \frac{1}{\varepsilon})$

• s.t.
$$A(\mathcal{E}_n) = \hat{h}_n$$
 with

$$Pr(Err(\hat{h}_n) < \epsilon) > 1 - \delta$$

C is polynomially PAC-learnable if Computational cost learning $(\hat{h}_n) = \text{Pol}(\frac{1}{\delta}, \frac{1}{\varepsilon})_{a = b + \varepsilon}$

PAC Learning: Probably Approximately Correct, 2 Weak learnability

- Idem strong learnability
- Except that one only requires error to be < 1/2 (just better than random guessing)

$$\varepsilon = \frac{1}{2} - \gamma$$

Question

Kearns & Valiant 88

- Strong learnability \Rightarrow weak learnability
- ► Weak learnability ⇒ some stronger learnability ??

PAC Learning: Probably Approximately Correct, 2 Weak learnability

- Idem strong learnability
- Except that one only requires error to be < 1/2 (just better than random guessing)

$$\varepsilon = \frac{1}{2} - \gamma$$

Question

Kearns & Valiant 88

- Strong learnability \Rightarrow weak learnability
- ► Weak learnability ⇒ some stronger learnability ??

Yes !

Strong learnability \Leftrightarrow Weak learnability

- PhD Rob. Schapire 89
- ▶ Yoav Freund, MLJ 1990: The strength of weak learnability
- Adaboost: Freund & Schapire 95

Overview

Introduction

Boosting PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Weak learnability \Rightarrow Strong learnability

Thm

Schapire MLJ 1990

Given algorithm able to learn with error $\eta = 1/2 - \gamma$ with complexity c

under any distribution D(x, y)

then there exists algorithm with complexity Pol(c), able to learn with error ε .

Proof (sketch)

- Learn h under D(x, y)
 Define D'(x, y): D(x, y) ∧ (Pr(h(x) ≠ y) = 1/2)
 Learn h' under D'(x, y)
 Define D''(x, y): D(x, y) ∧ (h(x) ≠ h'(x))
- ► Learn h" under D''(x, y)
- ► Use Vote(h, h', h'')

Proof (sketch)

Vote(h, h', h'') true if h and h' true, or ((h or h' wrong), and h'' true)

$$Pr(Vote(h, h', h'') OK) = Pr(h OK and h' OK) + Pr(h or h' \neg OK).Pr(h'' OK) \geq 1 - (3\eta^2 - 2\eta^3)$$

 $\textit{Err}(h) < \eta \Rightarrow \textit{Err}(\textit{Vote}(h, h', h'')) < 3\eta^2 - 2\eta^3$



Overview

Introduction

Boosting PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Adaboost

Freund Schapire 95

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

http://videolectures.net/mlss09us_schapire_tab/

Given

algorithm A, weak learner

$$\mathcal{E} = \{(x_i, y_i)\}, x_i \in \mathcal{X}, y_i \in \{-1, 1\}, i = 1 \dots n\}$$

Iterate

• For
$$t = 1, \ldots T$$

• Define distribution D_t on $\{1, \ldots n\}$

focussing on examples misclassified by h_{t-1}

- Draw \mathcal{E}_t after D_t or use example weights
- Learn h_t

$$Pr_{x_i \sim D_t}(h_t(x_i) \neq y_i) = \varepsilon$$

• Return: weighted vote of h_t

Adaboost

Init: D_1 uniform distribution

$$D_1(i)=\frac{1}{n}$$

Define D_{t+1} as follows

$$D_{t+1}(i) = \frac{1}{Z_t} D_t(i) \times \begin{cases} exp(-\alpha_t) & \text{if } h_t(x_i) = y_i \\ exp(\alpha_t) & \text{if } h_t(x_i) \neq y_i \end{cases}$$
$$= \frac{1}{Z_t} D_t(i) \times exp(-\alpha_t h_t(x_i) y_i)$$

With

►
$$Z_t$$
: normalisation term
► $\alpha_t = \frac{1}{2} ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right) > 0$
► $\varepsilon_t = Pr_{x_i \sim D_t} (h_t(x_i) \neq y_i)$

Adaboost

Init: D_1 uniform distribution

$$D_1(i)=\frac{1}{n}$$

Define D_{t+1} as follows

$$D_{t+1}(i) = \frac{1}{Z_t} D_t(i) \times \begin{cases} \exp(-\alpha_t) & \text{if } h_t(x_i) = y_i \\ \exp(\alpha_t) & \text{if } h_t(x_i) \neq y_i \end{cases}$$
$$= \frac{1}{Z_t} D_t(i) \times \exp(-\alpha_t h_t(x_i) y_i)$$

With

►
$$Z_t$$
: normalisation term
► $\alpha_t = \frac{1}{2} ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right) > 0$
► $\varepsilon_t = Pr_{x_i \sim D_t} (h_t(x_i) \neq y_i)$

Final hypothesis

$$H(x) = sign\left(\sum_{t} \alpha_t h_t(x)\right)$$





weak classifiers = vertical or horizontal half-planes

Round 1



α₁=0.42

Round 2



Round 3



 $\epsilon_{3=0.14}$ $\alpha_{3}=0.92$

Final Classifier



Bounding training error

Thm
Let
$$\varepsilon_t = \frac{1}{2} - \gamma_t$$

Then
 $Err_{train}(H) \leq \prod_t \left[2\sqrt{\varepsilon_t(1 - \varepsilon_t)} \right]$
 $= \prod_t \sqrt{1 - 4\gamma_t^2}$
 $\leq exp \left(-2\sum_t \gamma_t^2 \right)$

Analysis

• If A weak learner, $\exists \gamma \ s.t. \ \forall t, \ \gamma_t > \gamma > 0$

$$Err_{train}(H) < exp\left(-2\gamma^2 T\right)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

• Does not require γ or T to be known a priori

Proof

Note

$$F(x) = \sum_t \alpha_t h_t(x)$$

Step 1: final distribution

$$D_{\mathcal{T}}(i) = D_{1}(i) \prod_{t} \left[\frac{1}{Z_{t}} \exp\left(-y_{i} \sum_{t} \alpha_{t} h_{t}(x_{i})\right) \right]$$
$$= \frac{1}{n} \prod_{t} \left[\frac{1}{Z_{t}} \right] \exp\left(-y_{i} F(x_{i})\right)$$

Proof

Step 2 $Err_{train}(H) \leq \prod_{t} Z_t$

as

$Err_{train}(H) = \frac{1}{n}\sum_{i} \begin{cases} 1 & \text{if } H(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases}$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <
Step 2

$$Err_{train}(H) \leq \prod_{t} Z_{t}$$

as

$$\begin{aligned} & \textit{Err}_{train}(H) &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } H(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } y_i F(x_i) < 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

Step 2

$$Err_{train}(H) \leq \prod_{t} Z_{t}$$

as

$$\begin{aligned} \mathsf{Err}_{train}(H) &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } H(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } y_i F(x_i) < 0 \\ 0 & \text{otherwise} \end{cases} \\ &\leq \frac{1}{n} \exp\left(-y_i F(x_i)\right) \end{aligned}$$

Step 2

$$Err_{train}(H) \leq \prod_{t} Z_{t}$$

as

$$\begin{aligned} \mathsf{Err}_{train}(H) &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } H(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } y_i F(x_i) < 0 \\ 0 & \text{otherwise} \end{cases} \\ &\leq \frac{1}{n} \exp\left(-y_i F(x_i)\right) \\ &= \sum_{i} D_T(i) \prod_t Z_t \end{aligned}$$

Step 2

$$Err_{train}(H) \leq \prod_{t} Z_{t}$$

as

$$Err_{train}(H) = \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } H(x_i) \neq y_i \\ 0 & \text{otherwise} \end{cases}$$
$$= \frac{1}{n} \sum_{i} \begin{cases} 1 & \text{if } y_i F(x_i) < 0 \\ 0 & \text{otherwise} \end{cases}$$
$$\leq \frac{1}{n} exp(-y_i F(x_i))$$
$$= \sum_{i} D_T(i) \prod_t Z_t$$
$$= \prod_t Z_t$$

Step 3

$$Z_t = 2\sqrt{\varepsilon_t(1-\varepsilon_t)}$$

Because

$$Z_t = \sum_i D_t(i) \exp(-\alpha_t y_i h_t(x_i))$$

= $\sum_i / h_t(x_i) \neq y_i D_t(i) e^{\alpha_t} + \sum_i / h_t(x_i) = y_i D_t(i) e^{-\alpha_t}$
= $\varepsilon_t e^{\alpha_t} + (1 - \varepsilon_t) e^{-\alpha_t}$
= $2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$

Training error \neq test error ! (overfitting ?) Observed



Why?

Explanation based on the margin



The margin



	# rounds										
	5	100	1000								
train error	0.0	0.0	0.0								
test error	8.4	3.3	3.1								
% margins ≤ 0.5	7.7	0.0	0.0								
minimum margin	0.14	0.52	0.55								

Analysis

- 1. Boosting \Rightarrow larger margin
- 2. Larger margin \Rightarrow lower generalization error Why: if margin is large, hypothesis can be approximated by a simple one.

Intuition about Margin

Infant





Elderly



Man





Woman



Partial conclusion

Adaboost is:

- a way of boosting a weak learner
- a margin optimizer
- (other interpretations related to the minimization of an exponential loss)

However

In if Adaboost minimizes a criterion, it would make sense to directly minimize this criterion...

but direct minimization *degrades* performances....

Main weakness: sample noise

Noisy examples are rewarded.

Application: Boosting for Text Categorization [with Singer]

- weak classifiers: very simple weak classifiers that test on simple patterns, namely, (sparse) *n*-grams
 - find parameter α_t and rule h_t of given form which minimize Z_t
 - use efficiently implemented exhaustive search
- "How may I help you" data:
 - 7844 training examples
 - 1000 test examples
 - categories: AreaCode, AttService, BillingCredit, CallingCard, Collect, Competitor, DialForMe, Directory, HowToDial, PersonToPerson, Rate, ThirdNumber, Time, TimeCharge, Other.

Weak Classifiers

rnd	term	AC	AS	BC	CC	CO	CM	DM	DI	HO	PP	RA	3N	ΤI	тс	OT
1	collect	I	I	T	T	1	I	-	I	T	T	I	T	I	I	T
		T	T	T	-	T	T	-	T	T	T	T		T	T	•
2	card	T	-	-	L	-	-	-	•	-	-	•	T	T	■	•
2	mu hama	-	_			_	_	_	_	_	_	_	-			-
3	my nome	I -	I	-	-	-	-	-	I	_	-	-	-	I	I	-
4	person? person	I	I	•	•	-	I	-	I	-	I	I	-	I	-	T
5	code	•	-	-	-	-	-	-	-	-	_	-	_	-	T	-
			_	-	_	_	_	_	_	_	_		_	_	_	_
6	1	_	_	-	_	_	_	_	_	-	_	_	_	-	_	_
		_	-	T	-	_	-	_	_	-	_	_	_	_	_	_

More Weak Classifiers

rnd	term	AC	AS	BC	CC	CO	CM	DM	DI	HO	PP	RA	3N	ΤI	тс	ОТ
7	time	-	-	-	-	-	-	-	-	-	T	-	-			-
		_	_	_	_	_	_	_	_	_	_	_	_	T	-	
8	wrong number	I	T	L	•	-	T	■	T	•	T	T	•	T	T	T
		-	—	-	-	-	-	-	-	-	-	-	-		-	-
9	how	-	-	-	•	-	-	-	-	•	T	•	-	-	-	-
		-	_	-	_	-	-	_	—	-	_	-	_		-	_
10	call	-	-	—	-	—	-	-	-	_	-	-	-	-	-	-
		-	-	—	-	-	-	-	-	—	-	—	-	—	-	-
11	seven	∎	-	_	-	-	-	-	-	-	-	_	-	-	T	-
		-		-	-	-	-	-	-	-	-	-	-	-	-	-
12	trying to	-	-	-	-	-	-	-	-	-	-	-	-	•	T	-
		-	-	-	-	-	-	-	-	-	-	-	-	-	-	_
13	and	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
		-	-	-	_	—	-	-	_	-	_	_	-	_	-	-

More Weak Classifiers

rnd	term	AC	AS	BC	CC	CO	CM	DM	DI	НΟ	PP	RA	3N	тι	тс	от
14	third		T	-	T	T	T	-	T	T	-	T	L	T	T	_
		-	_	_	_	_	_	_	_	_	_	_	-	—	_	_
15	to	1	-	-	-	-	-	-	-	—	-	-	-	-	-	-
		-	_	-	_	_	_	_	_	_	-	-		-	-	-
16	for	-	-	-	-	-	-	-	-	-	-	-	-	T	-	-
		_	_	_	_	_	_	_	_	_	_	_	_	_	_	_
17	charges		-	-	-	L	-	-	-	-	-	-	-	T		-
		_	_	_	_	_	_	_	_	_	_	_	_	_	-	_
18	dial	-	-	-	-	-	-	-	-	-	T	_	-	T	T	-
		-	-	-	—	-	_	_	_	-	_	_	_	-	_	_
19	just	-	-	-	-	-	-	-	-	_	-	-	-	-	_	-
		-	_	_	_	_	_	_	_	_	_	_	_	_	_	_

Finding Outliers

examples with most weight are often outliers (mislabeled and/or ambiguous)

- I'm trying to make a credit card call (Collect)
- hello (Rate)
- yes I'd like to make a long distance collect call please (CallingCard)
- calling card please (Collect)
- yeah I'd like to use my calling card number (Collect)
- can I get a collect call (CallingCard)
- yes I would like to make a long distant telephone call and have the charges billed to another number (CallingCard DialForMe)
- yeah I can not stand it this morning I did oversea call is so bad (BillingCredit)
- yeah special offers going on for long distance (AttService Rate)
- mister allen please william allen (PersonToPerson)
- yes ma'am I I'm trying to make a long distance call to a non dialable point in san miguel philippines (AttService Other)

Application: Human-computer Spoken Dialogue [with Rahim, Di Fabbrizio, Dutton, Gupta, Hollister & Riccardi]

- application: automatic "store front" or "help desk" for AT&T Labs' Natural Voices business
- caller can request demo, pricing information, technical support, sales agent, etc.
- interactive dialogue



- NLU's job: classify caller utterances into 24 categories (demo, sales rep, pricing info, yes, no, etc.)
- weak classifiers: test for presence of word or phrase

Overview

Introduction

Boosting

PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - のへで

Bagging

Breiman96

Enforcing diversity through bootstrap: iterate

Draw

 $\mathcal{E}_t = n$ examples uniformly drawn with replacement from \mathcal{E} Learn h_t from \mathcal{E}_t Finally:

$$H(x) = Vote(\{h_t(x)\})$$

Bagging vs Boosting

Bagging: h_t independent

parallelisation is possible

Boosting: h_t depends from the previous hypotheses $(h_t \text{ covers up } h_1 \dots h_{t-1} \text{ mistakes}).$

Visualization In the 2d plane: distance, error Boosting Dietterich Margineantu 97

Bagging



Figure 1: Kappa-Error diagrams for ADABOOST (left) and bagging (right) on the Expf domain.

Analysis

Assume

- \mathcal{E} drawn after distribution P
- \mathcal{E}_t uniformly sampled from \mathcal{E}_t , h_t learned from \mathcal{E}_t
- Error of H, average of the h_t :

$$H(x) = \mathbb{E}_{\mathcal{E}_t}[h_t(x)]$$

Error

Direct error:

$$e = \mathbb{E}_{\mathcal{E}} \mathbb{E}_{X,Y}[(Y - h(X))^2]$$

Bagging error

$$e_B = \mathbb{E}_{X,Y}[(Y - H(X))^2]$$

Rewriting e:

$$e = \mathbb{E}_{X,Y}[Y^2] - 2\mathbb{E}_{X,Y}[YH] + \mathbb{E}_{X,Y}\mathbb{E}_{\mathcal{E}}[h^2]$$

and with Jensen inequality, ${\rm I\!E}[Z^2] \geq {\rm I\!E}[Z]^2$

 $e \ge e_B$

Overview

Introduction

Boosting

PAC Learning Boosting Adaboost

Bagging

General bagging Random Forests

▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - のへで

Random Forests



http://videolectures.net/sip08_biau_corfao/

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Classification / Regression



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─ のへで



E 990



E nar







E 990



三 のへで









▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで



◆□ → ◆□ → ◆三 → ◆三 → ◆○ ◆



◆□▶ ◆□▶ ◆目▶ ◆目▶ ● ● ● ●



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで
Build a tree



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

Build a tree



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 _ のへで

Build a tree



Random forests

Breiman00 Principle

stat.berkeley.edu/users/breiman/RandomForests

- Randomized trees
- Average

Properties

- Fast, easy to implement
- Excellent accuracy
- No overfitting % number of features

Theoretical analysis difficult

KDD 2009 - Orange

Targets

- 1. Churn
- 2. Appetency
- 3. Up-selling

Core Techniques

- 1. Feature Selection
- 2. Bounded Resources
- 3. Parameterless methods



Random tree

- In each node, uniformly select a subset of attribute
- Compute the best one
- Until reaching max depth



(日) (同) (日) (日)

Tree aggregation

- *h_t*: random tree
- Fast and straightforward parallelization



(日) (四) (王) (日) (日) (日)

Analysis

Biau et al. 10

$$\mathcal{E} = \{(x_i, y_i)\}, x_i \in [0, 1]^d, y_i \in \mathbb{R}, i = 1..n\} (x_i, y_i) \sim P(x, y)$$

Goal: estimate

$$r(x) = \mathbb{E}[Y|X = x]$$

Criterion

consistency

$$\mathbb{E}[(r_n(X) - r(X))^2]$$

Set $k_n \ge 2$, iterate $\log_2 k_n$ fois:

Select in each node a feature s.t.

 $Pr(feature. j selected) = p_{n,j}$

Split: feature j < its median</p>





▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで







Analysis

- Each tree has $2^{\log_2 k_n} = k_n$ leaves
- Each leaf covers $2^{-\lfloor \log_2 k_n \rfloor} = 1/k_n$ volume
- ► Assuming x_i uniformly drawn in [0, 1]^d, number of examples per leaf is ≈ n/k_n

► If
$$k_n = n$$
, very few examples per leaf
 $r_n(x) = \mathbb{E}\left[\frac{\sum_i y_i \mathbf{I}_{x_i,x} \text{ in same leaf}}{\sum_i \mathbf{I}_{x_i,x} \text{ in same leaf}}\right]$

Consistency

Thm

 r_n is consistent if $p_{nj} \log k_n \to \infty$ forall j and $k_n/n \to 0$ when $n \to \infty$.

◆□▶ ◆□▶ ◆∃▶ ◆∃▶ = のへで