How Developers Iterate on Machine Learning Workflows

-- A Survey of the Applied Machine Learning Literature

Doris Xin¹, Litian Ma¹, Shuchen Song¹, Rong Ma², Aditya Parameswaran¹ University of Illinois at Urbana-Champaign² Peking University





Highlights

Data prep.

Data joins

GLMs

SVMs

Model class

Human def. Features

(long live feat. eng.!)

SocS	NS	WWW	NLP	CV		
Join (31.0%)	Feature Def. (40.6%)	Feature Def. (36.1%)	Feature Def. (32.1%)	Feature Def. (37.5%)		
eature D ef. (27.6%)	Univar. FS (18.8%)	Join (22.2%)	BOW (17.9%)	B OW (25.0%)		
Normalize (17.2%)	Normalize (12.5%)	Normalize (13.9%)	Join (14.3%)	Interaction (25.0%)		
Impute (6.9%)	PCA (9.4%)	Discretize (8.3%)	Normalize (10.7%)	Join (12.5%)		

_imitations of data

- Incomplete picture of iterations
- Results not presented in iteration order
- Small per domain corpus can lead to spurious results

Remedies

- Multiple surveyors to reduce
- chance of spurious results
- Design iteration estimators
 - that do not rely on order

Table 1: Common DPR operations ordered top to bottom by popularity. Join = joining multiple data sources; Feat. def. = custom logic for fine-grained feature extraction; Univer. FS = univariate feature selection, using criteria such as support and correlation per feature; BOW = bag of words; PCA = principal component analysis, a common dimensionality reduction technique.

SocS	NS	WWW	NLP	CV
GLM (36.0%)	SVM (32.7%)	G LM (37.0%)	RNN (32.4%)	CNN (38.2%)
SVM (28.0%)	GLM (15.4%)	R F (11.1%)	GLM (14.7%)	S VM (17.6%)
R F (20.0%)	R F (13.5%)	S VM (11.1%)	S VM (11.8%)	R NN (17.6%)
Decision Tree (12.0%)	DNN (13.5%)	Matrix Factoriz. (11.1%)	CNN (8.8%)	R F (5.9%)

Table 2: Common model classes ordered top to bottom by popularity per domain. GLM = generalized linear models (e.g., logistic regression); RF = random forest; SVM = support vector machine; R/CNN = recursive/convolutional neural networks.

SocS	NS	WWW	NLP	CV
Regularize (40.0%)	CV (31.8%)	Regularize (41.2%)	LR (39.4%)	LR (46.2%)
CV (30.0%)	LR (22.7%)	LR (23.5%)	Batch Size (24.2%)	Batch Size (30.8%)
LR (10.0%)	DNN Arch. (18.2%)	B atch S ize (11.8%)	DNN Arch. (18.2%)	DNN Arch. (11.5%)
Batch Size (10.0%)	Kernel (9.1%)	CV (11.8%)	Kernel (6.1%)	Regularize (11.5%)

Table 3: Most popular model tuning operations by domain. CV = cross validation; LR = learning rate; DNN arch. = DNN architecture modification; Kernel specifically applies to SVM.

SocS	NS	WWW	NLP	CV		
P /R (25.7%)	Accuracy (28.6%)	Accuracy (20.8%)	P /R (29.2%)	Visualization (33.3%)		
Accuracy (20.0%)	P /R (18.6%)	P /R (20.8%)	Accuracy (27.1%)	Accuracy (29.8%)		
Feat. Contrib. (17.1%)	Visualization (15.7%)	Case S tudies (13.2%)	Case Studies (14.6%)	P /R (17.5%)		
Visualization (14.3%)	Correlation (11.4%)	DCG (9.4%)	Human Eval. (8.3%)	Case Studies (12.3%)		

Table 4: Most popular evaluation methods by domain. P/R = precision/recall; Feat. Contrib. = feature contribution to model performance; NCG = discounted cumulative gain, popular in ranking tasks; Case studies = case studies of individual results.

DNN only in NLP/CV

Model Tuning Learning rate fast Batch size ^Jtraining

D Evaluation

Coarse grained:

P/R, Accuracy

Fine-grained: Case studies, Vis.

Prioritize system support for common operations



(2.2) Estimating Iterations

		Data Prep.			ML	Model	Class	s ML Tuning Evaluation Metrics						
		norm.	impute	•••	DT	SVM		Reg.	λ		AUC		# tables	# figs
		v	v			v					v		5	2
Rebart Wolfer Fak Prediction Rystem: A March March March March March March March March March March March March March March March <td< th=""><th></th><th>~</th><th></th><th></th><th></th><th>~</th><th></th><th>~</th><th></th><th></th><th>~</th><th></th><th>5</th><th>2</th></td<>		~				~		~			~		5	2
		v	v			~					~		5	2
Α	ggregate	~	✓			~					~		5	2
			$n_{\mathcal{D}} = 2$		י <u>ן</u>	$n_{\mathcal{M}} =$	1	$n_{\mathcal{F}}$	γ = γ	 1	$n_{\mathcal{E}} =$	1	$n_{table} = 5$	n_{figur}
				_		_					_		4.	= 2

- Three surveyors independently identify operations mentioned in a paper.
- Take majority vote to aggregate the three results. Collect statistics on the aggregate annotation.

Estimator for # data prep. iterations $\hat{t}_{DPR} = n_{\mathcal{D}}$ Estimator for # ML iterations $\hat{t}_{LI} = (n_{\mathcal{M}} - 1) + (n_{\mathcal{P}} - 1)$ Estimator for # post proc. iter. $\hat{t}_{PPR} = \min(n_{\mathcal{E}}, n_{table} + n_{figure})$

Open source dataset at https://github.com/helix-ml/AppliedMLSurvey