

# Architecture of ML Systems 11 Model Debugging & Serving

#### **Matthias Boehm**

Graz University of Technology, Austria Computer Science and Biomedical Engineering Institute of Interactive Systems and Data Science BMVIT endowed chair for Data Management





### Announcements/Org

- #1 Video Recording
  - Link in TeachCenter & TUbe (lectures will be public)
  - Live streaming through TUbe, since May 08
  - Questions: <u>https://tugraz.webex.com/meet/m.boehm</u>

#### #2 AMLS Programming Projects

- Status: all project discussions w/ 15 students (~10 PRs)
- Soft deadline: June 30
- If unable to complete: email to <u>m.boehm@tugraz.at</u>
- Doodle for oral exam slots July 2/3 (1 done + 8 scheduled)
- #3 Course Evaluation
  - Please participate; open period: June 1 July 15





TUbe



Data Science Lifecycle









### Agenda

4

- Model Debugging and Validation
- Model Deployment and Serving





# Model Debugging and Validation

Similar to Software Testing Focus on Benchmarks, Assessment, Monitoring, Trust, Finding Room for Improvements





#### 6

### Recap: Data Validation

Sanity checks on expected shape before training first model

[Neoklis Polyzotis, et al: Data Management Challenges in Production Machine Learning. Tutorial, **SIGMOD 2017**]



Check a feature's min, max, and most common value

Research)

- Ex: Latitude values must be within the range [-90, 90] or  $[-\pi/2, \pi/2]$
- The histograms of continuous or categorical values are as expected
  - Ex: There are similar numbers of positive and negative labels
- Whether a feature is present in enough examples
  - Ex: Country code must be in at least 70% of the examples
- Whether a feature has the right number of values (i.e., cardinality)
  - Ex: There cannot be more than one age of a person

(Amazon Research)

Others

S [Sebastian Schelter et al: Automating Large-Scale Data Quality Verification. **PVLDB 2018**]

According long- from Sec-	-	-	
non ton to the Major		-	
Elbarres			
auna i		78	
-			

[Mike Dreves et al: From Data to Models and Back **DEEM@SIGMOD 2020**, <u>http://deem-workshop.org/videos/</u> <u>2020/8\_dreves.mp4</u>]



From data to models and ba





### **Overview Model Debugging**

- #1 Understanding via Visualization
  - Plotting of predictions / interactions
  - Combination with dimensionality reduction into 2D:
    - Autoencoder
    - **PCA** (principal component analysis)
    - t-SNE (T-distributed Stochastic Neighbor Embedding)
  - Input, intermediate, and output layers of DNNs

#### #2 Fairness, Explainability, and Validation via Constraints

- Impose constraints like monotonicity for ensuring fairness
- Generate succinct representations (e.g., rules) as explanation
- Establish assertions and thresholds for automatic validation and alerts



[Credit: twitter.com/tim kraska]

[Andrew Crotty et al: Vizdom: Interactive Analytics through Pen and Touch. PVLDB 2015]



[Credit: nlml.github.io/in-rawnumpy/in-raw-numpy-t-sne/]



### **Basic Model-Specific Statistics**

#### Regression Statistics

8

Average response and stddev, average residuals stddev residuals

correct

label

R2 (coeff of determination) with and without bias, etc

#### Classification Statistics

- Classical: recall, precision, F1-score
- Visual: confusion matrix (correct vs predicated classes)
   Junderstand performance wrt individual classes
- Example Mnist
- Mispredictions might also be visualized via dimensionality reduction







### Understanding Other Basic Issues

- Overfitting / Imbalance
  - Compare train and test performance
  - → Algorithm-specific techniques: regularization, pruning, loss, etc
- Data Leakage

9

- Example: time-shifted external time series data (e.g., weather)
- Compare performance train/test vs production setting
- Covariance Shift (features)
  - Distribution of features in training/test data different from production data
  - Reasons: out-of-domain prediction, sample selection bias
  - Examples: NLP, speech recognition, face/age recognition
- Concept Drift (features → labels)
  - Gradual change of statistical properties / dependencies (features-labels)
  - Requires re-training, parametric approaches for deciding when to retrain





(d) Classifier, probability

of correct class

### **Occlusion-Based Explanations**

- Occlusion Explanations
  - Slide gray square over inputs
  - Measure how feature maps and classifier output changes

	Institute and Tablemanding
	Conderoad Release
16	
1 hours	
10	1000
1000	

10

[Matthew D. Zeiler, Rob Fergus: Visualizing and Understanding Convolutional Networks. **ECCV 2014**]

(1)	(,,	 -
True Label: Pomeranian	10	10
DEKRA BACI 		
True Label: Afghan Hound		

(b) Laver 5 strongest feature man

(c) Layer 5, strongest

feature man projection

- Incremental Computation of Occlusion Explanations
  - View CNN as white-box operator graph and operators as views
  - Materialize intermediate tensors and apply incremental view maintenance

[Supun Nakandala, Arun Kumar, and Yannis Papakonstantinou: Incremental and Approximate Inference forFaster Occlusion-based Deep CNN Explanations, **SIGMOD 2019**]



SIGMOD 2020 Research Highlight

706.550 Architecture of Machine Learning Systems – 11 Model Debugging and Serving Matthias Boehm, Graz University of Technology, SS 2020

(a) Input Image





## Saliency Maps

Saliency Map

11

- Given input image and specific class
- Compute saliency map of class derivatives wrt input image
- Approximated w/ a linear function (Taylor expansion)

[Karen Simonyan, Andrea Vedaldi, Andrew Zisserman: Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. ICLR Workshop 2014]



Unsupervised
 Image
 Segmentation







### **Example Model Anomalies**

#### #1 Wolf Detection based on snow cover



[Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin: Why Should I Trust You?: Explaining the Predictions of Any Classifier, **KDD 2016**]

### #2 Horse Detection based on image watermarks

Layer-wise relevance propagation

	Color and
1000	

[Sebastian Lapuschkin et al.: Analyzing Classifiers: Fisher Vectors and Deep Neural Networks, **CVPR 2016**]

### "silent but severe problems"





(a) Husky classified as wolf

(b) Explanation



#### #3 Race-biased Jail Risk Assessment

#BlackLivesMatter 🏶

[Julia Angwin et al: Machine Bias – There's software used across the country to predict future criminals. And it's biased against blacks, **2016**, <u>https://www.propublica.org/article/</u> <u>machine-bias-risk-assessments-in-criminal-sentencing</u>]







## Debugging Bias and Fairness

#### Fairness

13

- Validate and ensure fairness with regard to sensitive features (unbiased)
- Use occlusion and saliency maps to characterize and compare groups

#### Enforcing Fairness

- Use constraints to enforce certain properties (e.g., monotonicity, smoothness)
- Example: late payment → credit score





## **Explanation Tables**

Motivation

14

- Generate succinct decision rules from data
- Problem: Decision tree rules do not overlap by def
- Example athlete's exercise log:
   "Goal met" → 7 vs 7

#### Explanation Tables

- Find smallest explanation table subject to max KL divergence threshold
- Greedy and sampling algorithms

-	-	a Estentine of Second
100		er
	witter	10000
COARD.		100000
		12231
		112 211
		Filmed & Fil
		and the second second
14/62	10.75-72	

[Kareem El Gebaly, Parag Agrawal, Lukasz Golab, Flip Korn, Divesh Srivastava: Interpretable and Informative Explanations of Outcomes. **PVLDB 2014**]

id	day	time	meal	goal met?
1	Fri	Dawn	Banana	Yes
2	Fri	Night	Green salad	Yes
3	Sun	Dusk	Oatmeal	Yes
4	Sun	Morning	Banana	Yes
5	Mon	Afternoon	Oatmeal	Yes
6	Mon	Midday	Banana	Yes
7	Tue	Morning	Green salad	No
8	Wed	Night	Burgers	No
9	Thu	Dawn	Oatmeal	Yes
10	Sat	Afternoon	Nuts	No
11	Sat	Dawn	Banana	No
12	Sat	Dawn	Oatmeal	No
13	Sat	Dusk	Rice	No
14	Sat	Midday	Toast	No



day	time	meal	goal met=Yes?	count
*	*	*	.5	14
Sat	*	*	0	5
*	*	Banana	.75	4
*	*	Oatmeal	.75	4





### Data Slices and Automatic Slice Finding

Problem

15

 Find top-k data slices where model performs worse than average [Yeounoh Chung et al.: Slice Finder: Automated Data Slicing for Model Validation. CoRR 2018/ICDE2019]



- Data slice: S<sup>DG</sup> := D=PhD AND G=female (subsets of features)
- Objective: w \* err(S<sup>DG</sup>)/err(S<sup>\*</sup>) + (1-w) \* |S<sup>DG</sup>|

#### Existing Algorithms

- Preparation: Binning +
   One-Hot Encoding of features
- Lattice search with heuristic, level-wise termination

#### Extensions

- #1 Lower/upper bounds → pruning & termination
- #2 Large-scale task-/data-parallel computation

SystemDS/scripts/
staging/slicing



ISDS

706.550 Architecture of Machine Learning Systems – 11 Model Debugging and Serving Matthias Boehm, Graz University of Technology, SS 2020

"find largest error vs find large slices"





### Model Assertions

Motivation

16

[Daniel Kang, Deepti Raghavan, Peter Bailis, Matei Zaharia: Model Assertions for Debugging Machine Learning, **NIPS Workshop ML Systems, 2018**]

17,5

- ML models might fail in complex ways that are not captured in loss function
- Inspired by assertions in SW dev  $\rightarrow$  Model assertions via Python rules

**Example:** Flickering of object detection



(a) Frame 1, base SSD





(b) Frame 2, base SSD

(c) Frame 3, base SSD

#### Assertion Use Cases

- #1 Runtime monitoring (collect statistics on incorrect behavior)
- #2 Corrective Action (trigger corrections at runtime) -> but how in retrospect?
- #3 Active Learning (decide which difficult data points to give to user)
- #4 Weak supervision (propose alternative labels and use for retraining)





#### 17

### **Continuous Integration**

 System Architecture ease.ml/ci [Cedric Renggli, Bojan Karlaš, Bolin Ding, Feng Liu, Kevin Schawinski, Wentao Wu, Ce Zhang: Continuous Integration of Machine Learning Models with ease.ml/ci: Towards a Rigorous Yet Practical Treatment, **SysML 2019**]

**Test Condition and Reliability Guarantees** 



#### ml: Github Repository - script : ./test\_model.pv : n - o > 0.02 + / - 0.01- condition Define test script - reliability: 0.9999 ./.travis.yml - mode : fp-free - adaptivity : full ./.testset steps 32 **2** Provide N test examples ./ml\_codes or **Technical Contribution Example Test** Provide guidelines on Condition 3 Commit how large N is in a New model has at a new ML ml/ci passed ml/ci failed declarative, rigorous, least 2% higher model/code but still practical way, accuracy, estimated enabled by novel within 1% error, system optimization with probability or techniques. 0.9999. 4 Get pass/fail signal





# Model Deployment and Serving

How to Exchange, Deploy, and Update Trained Models How to Efficiently Serve Prediction Tasks





### <sup>19</sup> Model Exchange Formats

- Definition Deployed Model
  - #1 Trained ML model (weight/parameter matrix)
  - #2 Trained weights AND operator graph / entire ML pipeline
    - → especially for DNN (many weight/bias tensors, hyper parameters, etc)

#### Recap: Data Exchange Formats (model + meta data)

- General-purpose formats: CSV, JSON, XML, Protobuf
- Sparse matrix formats: matrix market, libsvm
- Scientific formats: NetCDF, HDF5
- ML-system-specific binary formats (e.g., SystemML binary block)
- Problem ML System Landscape
  - Different languages and frameworks, including versions



%
% 0 or more comment lines
%
5 5 8
1 1 1.000e+00
2 2 1.050e+01
3 3 1.500e-02
1 4 6.000e+00
4 2 2.505e+02
4 4 -2.800e+02
4 5 3.332e+01
5 5 1.200e+01

%MatrixMarket matrix coordinate real general

20



### Model Exchange Formats, cont.

- Why Open Standards?
  - Open source allows inspection but no control
  - Open governance necessary for open standard
  - Cons: needs adoption, moves slowly
- #1 Predictive Model Markup Language (PMML)
  - Model exchange format in XML, created by Data Mining Group 1997
  - Package model weights, hyper parameters, and limited set of algorithms

#### #2 Portable Format for Analytics (PFA)

- Attempt to fix limitations of PMML, created by Data Mining Group
- JSON and AVRO exchange format
- Minimal functional math language → arbitrary custom models
- Scoring in JVM, Python, R



[Nick Pentreath: Open Standards for Machine Learning Deployment, bbuzz 2019]



21



### Model Exchange Formats, cont.

- #3 Open Neural Network Exchange (ONNX)
  - Model exchange format (data and operator graph) via Protobuf
  - First Facebook and Microsoft, then IBM, Amazon → PyTorch, MXNet
  - Focused on deep learning and tensor operations
  - ONNX-ML: support for traditional ML algorithms
  - Scoring engine: <u>https://github.com/Microsoft/onnxruntime</u>
  - Cons: low level (e.g., fused ops), DNN-centric → ONNX-ML

Lukas Timpl python/systemds/ onnx\_systemds

#### TensorFlow Saved Models

- TensorFlow-specific exchange format for model and operator graph
- Freezes input weights and literals, for additional optimizations (e.g., constant folding, quantization, etc)
- Cloud providers may not be interested in open exchange standards





Apache

SvstemML<sup>\*\*</sup>

### **ML Systems for Serving**

- #1 Embedded ML Serving
  - TensorFlow Lite and new language bindings (small footprint, dedicated HW acceleration, APIs, and models: MobileNet, SqueezeNet)
  - SystemML JMLC (Java ML Connector)

#### #2 ML Serving Services

- Motivation: Complex DNN models, ran on dedicated HW
- RPC/REST interface for applications
- TensorFlow Serving: configurable serving w/ batching
- Clipper: Decoupled multi-framework scoring, w/ batching and result caching
- Pretzel: Batching and multi-model optimizations in ML.NET
- Rafiki: Optimization for accuracy under latency constraints, and batching and multi-model optimizations

_	hand	a bran	_	
			- 10	
- 2	83	1.3		
		-		
8			23	
100			121/21	
膜			10	
100	euno	2.7.00	9173	
2103				

22

[Christopher Olston et al: TensorFlow-Serving: Flexible, High-Performance ML Serving. NIPS **ML Systems 2017**]



[Daniel Crankshaw et al: Clipper: A Low-Latency Online Prediction Serving System. **NSDI 2017**]



[Yunseong Lee et al.: PRETZEL: Opening the Black Box of Machine Learning Prediction Serving Systems. **OSDI 2018**]



[Wei Wang et al: Rafiki: Machine Learning as an Analytics Service System. **PVLDB 2018**]

Example: Google Translate 140B words/day → 82K GPUs in 2016

**TensorFlowLite** 

23



### Example SystemDS JMLC







### Example SystemDS JMLC, cont.

Background: Frame

24

- Abstract data type with schema (boolean, int, double, string)
- Column-wise block layout
- Local/distributed operations:
   e.g., indexing, append, transform



**Distributed representation:** ? x ncol(F) blocks

(shuffle-free conversion of csv / datasets)







### Example SystemML JMLC, cont.

Motivation

25

- ➔ Embedded scoring
- → Latency ⇒ Throughput
- $\rightarrow$  Minimize overhead per  $\Delta X$

Example

#### Typical compiler/runtime overheads:

Script parsing and config:	~100ms
Validation, compile, IPA:	~10ms
HOP DAG (re-)compile:	~1ms
Instruction execute:	<0.1µs

// single-node, no evictions,

```
1: Connection conn = new Connection(); // no recompile, no multithread.
2: PreparedScript pscript = conn.prepareScript(
    getScriptAsString("glm-predict-extended.dml"),
    new String[]{"FX", "MX", "MY", "B"}, new String[]{"FY"});
3: // ... Setup constant inputs
4: for( Document d : documents ) {
5: FrameBlock FX = ...; //Input pipeline
6: pscript.setFrame("FX", FX);
7: FrameBlock FY = pscript.executeScript().getFrame("FY");
8: // ... Remaining pipeline
9: } // execute precompiled script
// many times
```

26









## Serving Optimizations – MQO

#### Result Caching

27

 Establish a function cache for X → Y (memoization of deterministic function evaluation)

#### Multi Model Optimizations

- Same input fed into multiple partially redundant model evaluations
- Common subexpression elimination between prediction programs
- Done during compilation or runtime
- In PRETZEL, programs compiled into physical stages and registered with the runtime + caching for stages (decided based on hashing the inputs)



FrontEnd

(1)

**(**5)



[Yunseong Lee et al.: PRETZEL: Opening the Black Box of Machine Learning Prediction Serving Systems. **OSDI 2018**]

Runtime

706.550 Architecture of Machine Learning Systems – 11 Model Debugging and Serving Matthias Boehm, Graz University of Technology, SS 2020



Predict(m: ModelId, x: X) -> y: Y



## Serving Optimizations – Quantization

Quantization

28

- Lossy compression via ultra-low precision / fixed-point
- Ex.: 62.7% energy spent on data movement

08 Data Access Methods

[Amirali Boroumand et al.: Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks. **ASPLOS 2018**]



#### Quantization for Model Scoring

- Usually much smaller data types (e.g., UINT8)
- Quantization of model weights, and sometimes also activations
   → reduced memory requirements and better latency / throughput (SIMD)

```
import tensorflow as tf
converter = tf.lite.TFLiteConverter.from_saved_model(saved_model_dir)
converter.optimizations = [tf.lite.Optimize.OPTIMIZE_FOR_SIZE]
tflite_quant_model = converter.convert()
```

[Credit: https://www.tensorflow.org/lite/performance/post\_training\_quantization ]





### Serving Optimizations – Compilation



- Compile entire TF graph into binary function w/ low footprint
- Input: Graph, config (feeds+fetches w/ fixes shape sizes)
- Output: x86 binary and C++ header (e.g., inference)
- Specialization for frozen model and sizes

#### PyTorch Compile

- Compile Python functions into ScriptModule/ScriptFunction
- Lazily collect operations, optimize, and JIT compile
- Explicit jit.script call or@torch.jit.script

```
a = torch.rand(5)
def func(x):
    for i in range(10):
        x = x * x # unrolled into graph
    return x
jitfunc = torch.jit.script(func) # JIT
```

```
New PySoch Optimizes
Deep Learning Computations
```

[Vincent Quenneville-Bélair: How PyTorch Optimizes Deep Learning Computations, Guest Lecture Stanford 2020]

```
jitfunc = torch.jit.script(func) # J
jitfunc.save("func.pt")
```

706.550 Architecture of Machine Learning Systems – 11 Model Debugging and Serving Matthias Boehm, Graz University of Technology, SS 2020



04 Adaptation,

**Fusion**, and **JIT** 

[Chris Leary, Todd Wang: XLA – TensorFlow, Compiled!, **TF Dev Summit 2017**]







### Serving Optimizations – Specialization

- NoScope Architecture
  - Baseline: YOLOv2 on 1 GPU per video camera @30fps
  - Optimizer to find filters

Helicipe, IS	ownering Marrier Methods (Section For Thins of Freid)
1000100-000	Name of the Association of the Association
MUTCH	and Manageria

[Daniel Kang et al: NoScope: Optimizing Deep CNN-Based Queries over Video Streams at Scale. **PVLDB 2017**]

- #1 Model Specialization
  - Given query and baseline model
  - Trained shallow NN (based on AlexNet) on output of baseline model
  - Short-circuit if prediction with high confidence

#### #2 Difference Detection

- Compute difference to ref-image/earlier-frame
- Short-circuit w/ ref label if no significant difference



ISDS





### Summary and Conclusions

- Summary 11 Model Deployment
  - Model Debugging and Validation
  - Model Deployment and Serving
- #1 Finalize Programming Projects by ~June 30

#### #2 Oral Exam

31

- Doodle for July 2/3, 45min each (done via skype/webex)
- Part 1: Describe you programming project, ram-up questions
- Part 2: Questions on 2-3 topics of 11 lectures (basic understanding of the discussed topics / techniques)

cisco Webex

