

SmartCoding

Design and Implementation of ML project

Contents



Introduction

Why SmartCoding?

What was done?

Conclusion

Time

Resources

The SmartCoding Project

Data

People

Introduction





Dasware Now it all just happens

Who am I?

- Olli-Pekka Valtonen
- Education: Software Engineer
- Various roles at Basware:
 - Software Developer
 - System Architect
 - Project / Program Manager
 - Team Lead
 - Technical Product Owner







What is Basware?

- Provides leading world-class Networked Accounts Payable Saas service.
- Basware invented the Invoice
 Automation space
- Founded in Finland (came to Tampere through buying Analyste)
- Previously listed in HEX,
- Delisted, owned by AKKR
- ~1500 employees.



Why SmartCoding?





Why SmartCoding?

\checkmark

Invoices without purchase order are time consuming & painful

- Examples: Team Dinner, Services, Utilities, Coffee for the office
 - (These also could be with purchase order, Basware offers tools to automate most of these)
- PO:less invoicess need to be coded manually.
- This is work that nobody wants to do.

Dasware Now it all just happens^{**}

Oppoturnity

Olli-Pekka Valtonen • You

Basware | Code Spelunker

2mo • Edited • S Imagine that your organization receives around 500k invoices annually and let's assume each invoice takes around five minutes to process.

Even with 'traditional' automation in place, it would take over 20 AP Clerks to process these invoices. This is a significant investment in workforce. As the business grows, headcount needs to be increased to keep up with the volume of processed invoices.

...

For a big chunk of these invoices, there'll be a purchase order. By linking the order to the invoice we can automate the manual work of creating coding. For PO:less invoices, we can use **#Basware #SmartCoding** to create the coding rows with machine learning.

Unfortunately, not everything can be immediately automated. Getting high quality digitized invoices also takes significant effort. Luckily, these are solved problems.

Let us help you to direct investment where it is needed. We can help you to start the journey towards fully scalable Touchless Invoice Processing. #nowitalljusthappens #machinelearning #automation #digitalization #accountspayable #accountspayableautomation

Increased automation

• Fully automating processing of 50k invoices saves around 2 FTE.



What was done?





SmartCoding





- Accounts Payable: Feed succesfully coded invoices into data platform.
- Clean, sanitize and dedup the data in Data Platform.
- Feed the data from Data Platform to SmartCoding microservice.
- In the microservice train machine learning model(s) that predict the best coding for given invoice.
- As PO:less invoice arrives to invoice process, ask the microservice to infer the best coding using ML and display it to the user.

Accounts Payable



< REST cXML Ltd Invoice Valid						183.00 EUR GROSS	150.0	DO EUR NET Send to process	•
Images 1		Discussions 🖓	Header data 🕑	Related documents 🕑	Workflow	Attachments 🖓 🛛 I	nvoice line	es 🔄 History 🕾	
Viva.png Open image De	Organization		Involce type		Supplier code *		Supplier bank name		
Q 100% - Q	o c	Basware Oyj	-	en-US:Invoice/en-US	***	100010, REST cXML Ltd		Bank	
	Supplier bank IBAN		Supplier bank BBAN		Invoice number *		Reference person		
100				111111-1111111		SMC_1			
	lens 1	: Invoice date *		Base date *		Exchange rate base date *		Payment term code *	
C 88.5		4/28/2023	m	4/28/2023	8	4/28/2023		0014, Within 14 days Due net	•
and the second s		Currency code *		Currency code (company)		Currency code (organization)		Reference number	
		EUR	-	EUR		EUR			
The second se		Exchange rate (company) *		Exchange rate (organization) *		Gross total *		Gross total (company)	
		1		1.000000		183.00		183.00	
		firess total (organiz	ation)	Net total *		Net total (company) *		Net total (organization) *	
Coding 1	1041X								
Coding rows generated by SmartCoding. View off	er proposals								
Iter co Q Gross total: 183.00 Coding gros	s total: 186.00 🗌 C	oding net total: 150.00	Coding tax total:	36.00 Coding tax total 2:	0.00 Net	sum difference: 0.00 Gross	sum diffe	rrence: -3.00 Tax s Add cod	ding
# Next reviewer * Account	Code Cos	t Center Code	Project Code	Profit Center Code	Profit Co	ntër Name Latest Cor	nment	Net Total (Organization)	Action

Dasware Now it all just happens™

ML Microservice



Data Platform

We didn't do this alone

- Data Lake / Data Warehouse solution
- Responsible for ingesting and storing the raw data
- Responsible for computing and transforming the data to required format consumed by the ML microservice.

Built by another team - blackbox for the SmartCoding team

The SmartCoding Project





Step 1: Plans & POCs

Plans

- Project plan
- Business plan
- Architecture plan
- E2E testing plan
- Delivery plan
- Syncing deliveries with Data Platform

POCs & Research

- Training & Inference POC
- ML Model Research

Most plans were iterated along the way.

Prereqs for success



Well-formed data Correct data Data practitioners Cloud native developers & QA. Sales that can speak ML Computation environment such as Azure or AWS Operations

resources

Whew

Data





Let's talk about data

".. bbbut data is the new oil..."

- Basware has well-formed data from human created and operated systems.
- Basware has typically data that is valid and often correct as it is used for accounting purposes.

Is there machine learning to be done?

- First task for a data scientist is to figure out the nature of the problem (classification vs regression)
- Second task is to uncovern the patterns and relationships to determine whether to use statistical methods or machine learning.
 - Multiple models? A single model?
 - Which method?
 - Etc



"Data has no value when sitting idly in a database. This simply is not so. Just as inventory sitting on a shelf in a warehouse has discernible value, so do idle information assets. This is the difference between realized value and the true definition of an asset, which takes into consideration its probable future economic benefit. All information has a probable future economic benefit."

Douglas B. Laney



People





Let's talk about roles

			\mathbf{V}		
Data Scientist	Architect	Dev & Test	MLOPS		
Solves the data problem	Solves the software problem	Develops the service	Configures the models		
Research & Model Design	Fitting inference into the system	Implementation of data pipelines	Monitors the system Deploys the models Retrains.		
Has requirements for the data architecture	Microservice architecture Has opinion on the data architecture	Test automation Test tooling E2E testing			

Missing skills

ML & MLOPS

- Data Scientist
- MLOPS engineer

Architecture, Dev & QA

- Cloud native development
- Serverless microservice
- Serverless QA

We didn't really find tasks for LinkedIn Influencers with MANIC ENERGY 🛞



Data Practitioner Rampup

Data Scientist

• Critical for project success

Rest of the organization

- Very difficult domain for multi-taskers
- Requires basic statistics / math skills
 - (High-school level)
- Both data and ML understanding needs to be developed in sales, marketing and consulting.

MLOPS

- Training existing personnel
- Basic Math / Data Skills
- Ops skills
- Curious mindset

Team Composition

Normal software development tasks

- Architect roles are vital to tie together the 'data problem', the 'ML problem' and the 'software problem'.
- Wide variety of tasks: data pipelines, writing ML code, writing tests, writing user interfaces, cost estimations, security and writing code.





It takes around three months of 100% full-time-equivalent on the job learning to bring competent Senior Developer familiar with AWS to get fluent with Serverless Microservice Architecture on AWS



Resources





No cloud computing? No ML (for us)

Mature production stack

- Access to cloud computing resources
- AWS SageMaker for the ML
- Python based microservice stack to run everything.
 - Python seems to be the language of ML
- Something like Splunk will make your life much easier.
 - You want to be able to trace the 'an action' thoughout the system.

Operability

Right skill for the correct task.

- Using R&D Data Scientist for MLOPS tasks is overkill.
 - Results in unhappy people.
- Overlap with production tasks
 - Monitoring, configuration, problem solving
- Navigating data practitioners can be hard: who is the right person?
 - Production Data Scientist, Data Engineer, ML Engineer...

Growing new muscles is very hard...

• and should not be underestimated.

Delivery

Think very hard before you make anything configurable...

- Especially if you intend to take it into use for 1000+ customers.
- Could it still be automated?

Whether you like it or not, somebody needs to support the delivery...

• Which is particularly annoying when you can't rely on on person due to timezone differences.

Training an organization to do something for the very first time...

• is very frustrating.



Machine learning is not a deterministic system.

Biggest challenge is to tell an end user that the system is predicting a particular case with 65% confidence and that it might be wrong.

Yet looking at the statistics, the system is predicting with 85% accuracy and performing well.



Data Scientists require access to prod data

In some orgs R&D doesn't have production access...

- This changes with the introduction of data scientist. Production data is important resouce for any ML project.
- Access should be lightweight.
- Building suitable test data is difficult for complex systems
 - Map, pertubate, mask but retain statistical distribution of prod data and bring it to lesser environments.

Time





Let's talk about time



Conclusion





Key learnings 1/2

Planning & Pocs

• Time spent on building POCs & training was pretty small.

Still mostly a software problem (for project manager)

• 85% software problem, max 15% ML problem

Machine learning is not deterministic

 People's math skills suck and conceptualizing confidence vs accuracy is difficult.

Building & operating a distributed system is expensive

- Requires relatively rare talent and and advanced skills through the organization.
 Training competent people is easy, but...
- Finding the right talent is hard.

Key learnings 2/2

Leave room for business iteration

- We redid many things as we learned: introduced multiple models per tenant, redid billing
 Product Management is important
- Users are impacted by the non-determinism.

Working with R&D people is easy

• So easy it distorts your reality. With non R&D people not so easy.

Do not underestimate developing new skills and talent

• Skill rampup and recruitment take time.

You don't really need big data solution in the beginning...

• but it is useful later on.

Thank you

Olli-Pekka Valtonen





