

# Hindsight Bias: How to deal with label leakage at scale

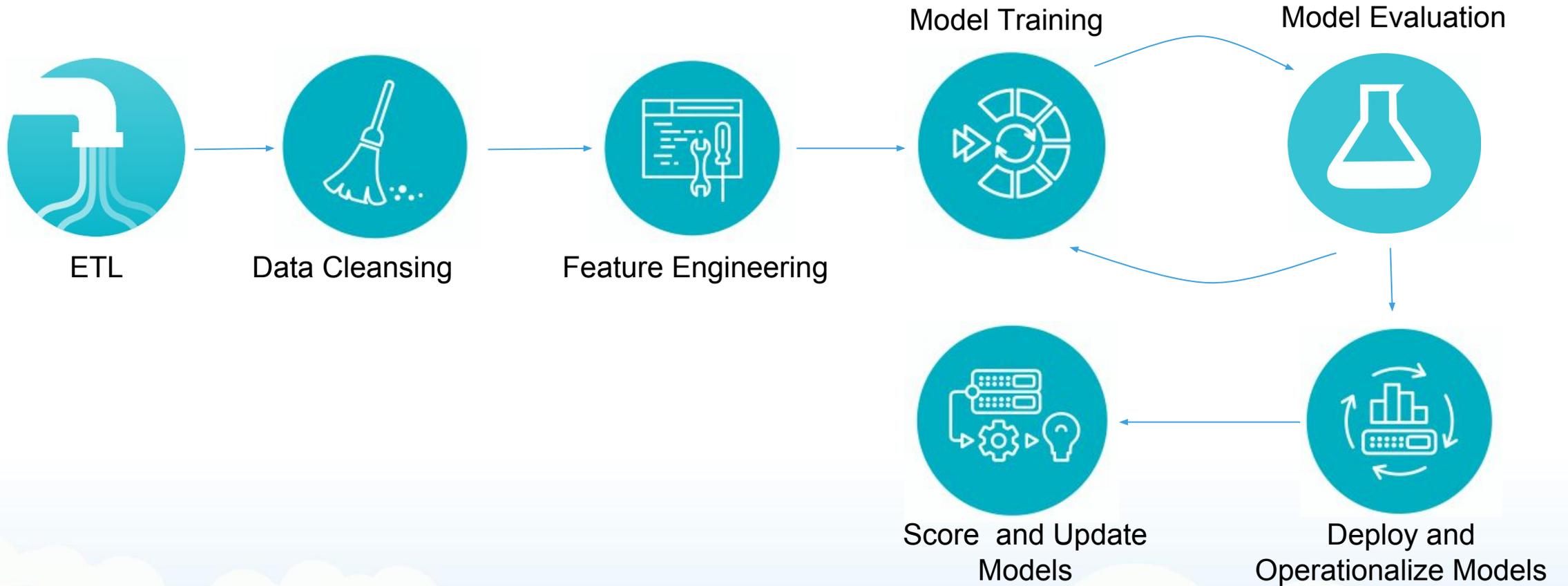
Till Bergmann (PhD), Senior Data Scientist  
tbergmann@salesforce.com



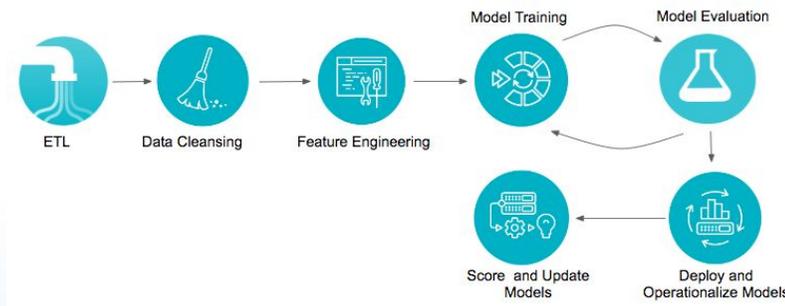
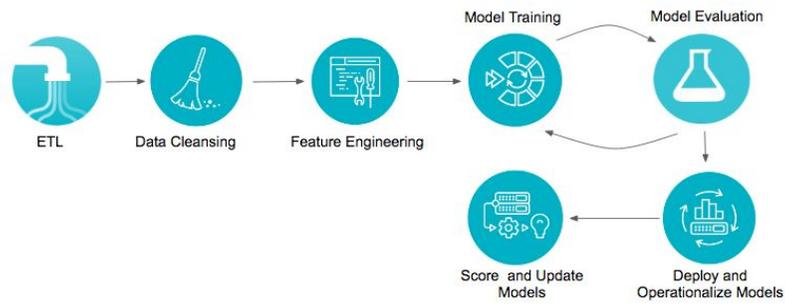
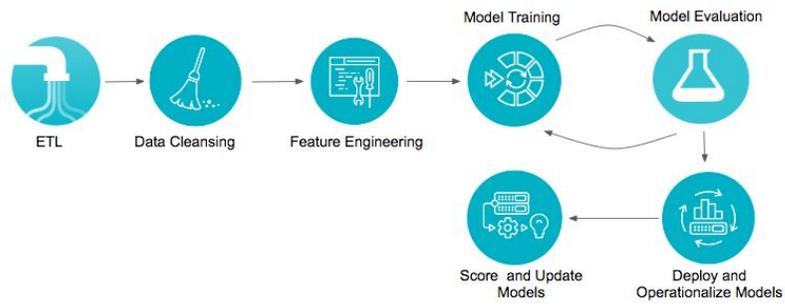
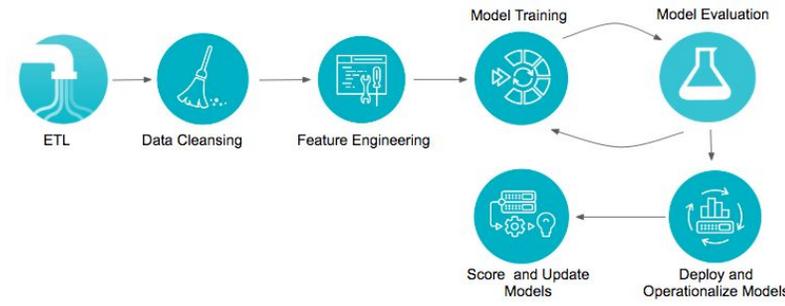
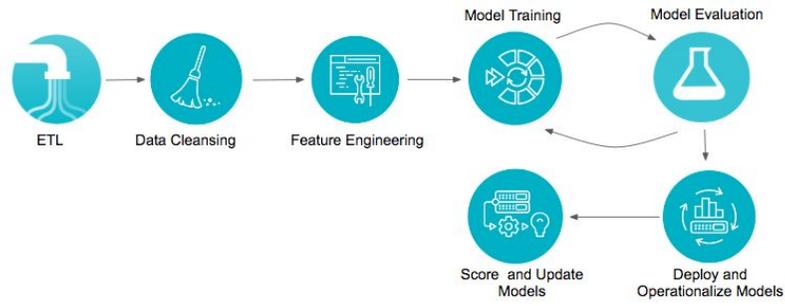
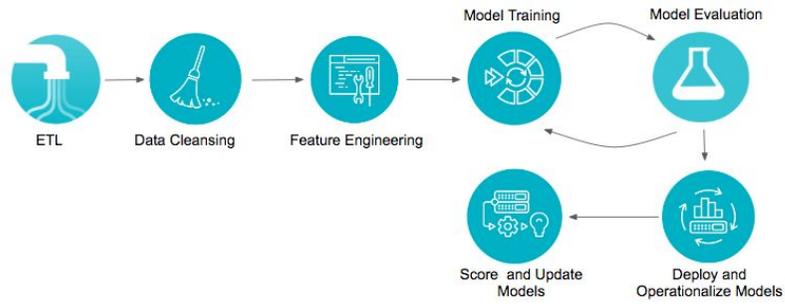
# 24 Hours in the Life of Salesforce



# The typical Machine Learning pipeline



# Multiply it by M\*N (M = customers; N = use cases)



# Problems with enterprise data

## Not enough data scientists to hand tune each model

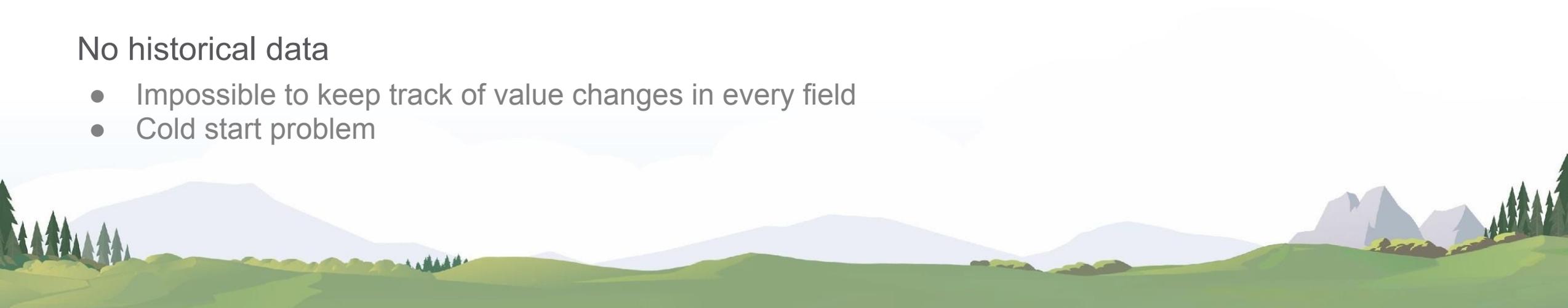
- We don't know the specific business use case and data
- Each step in the pipeline needs to be automated

## Messy data

- Nobody likes data entry - missing fields, typos
- Automated business practices can lead to patterns in the data
- Custom fields get added, removed or deprecated at any time

## No historical data

- Impossible to keep track of value changes in every field
- Cold start problem



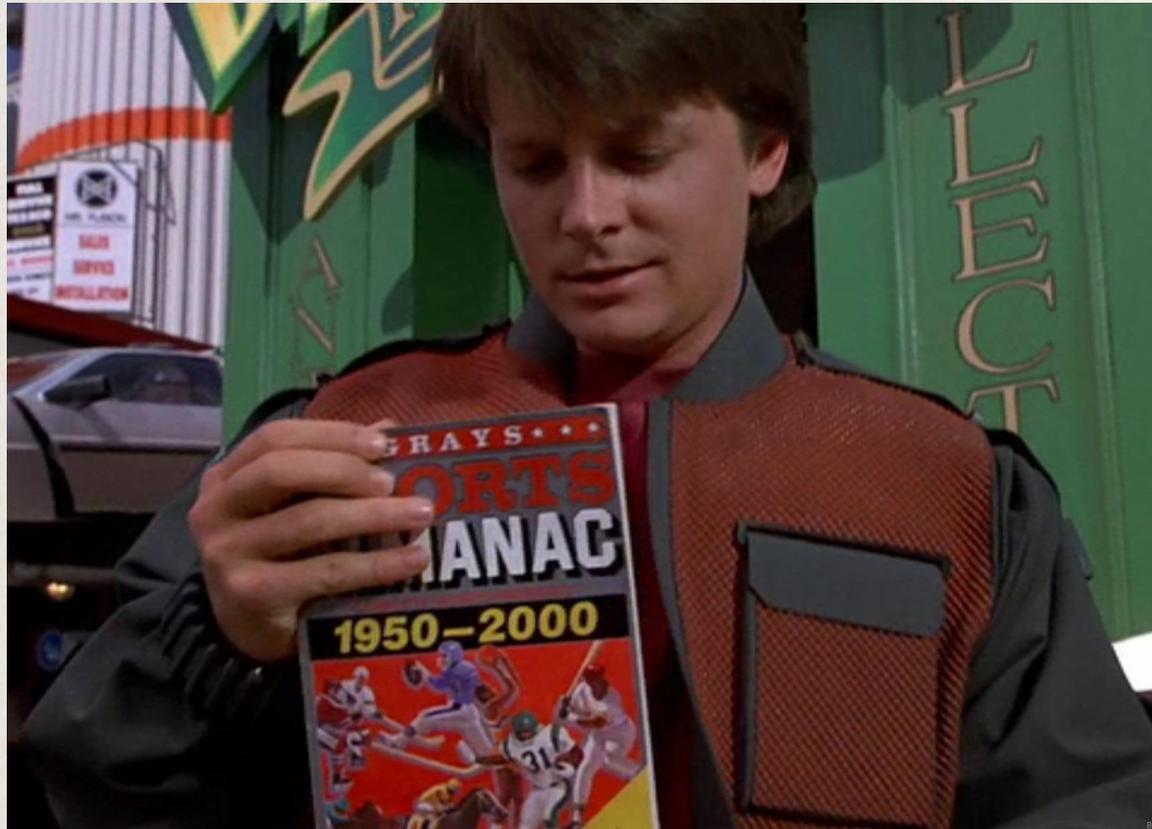
# What is hindsight bias?

Label/data leakage



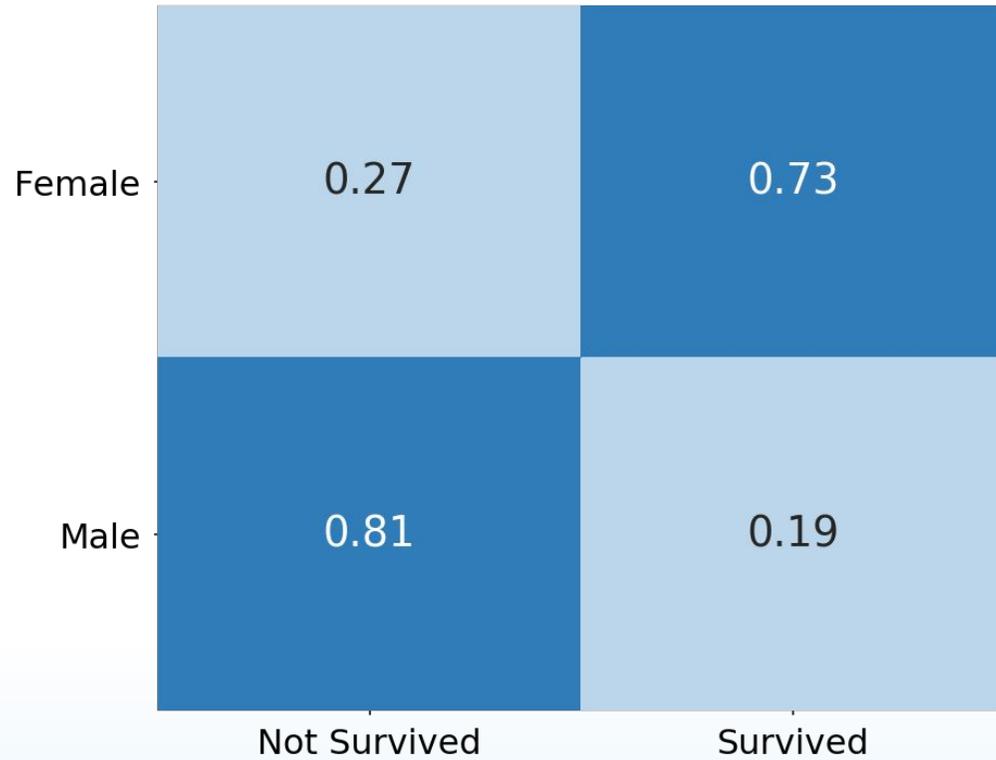
# Back to the future

Knowing things you shouldn't know



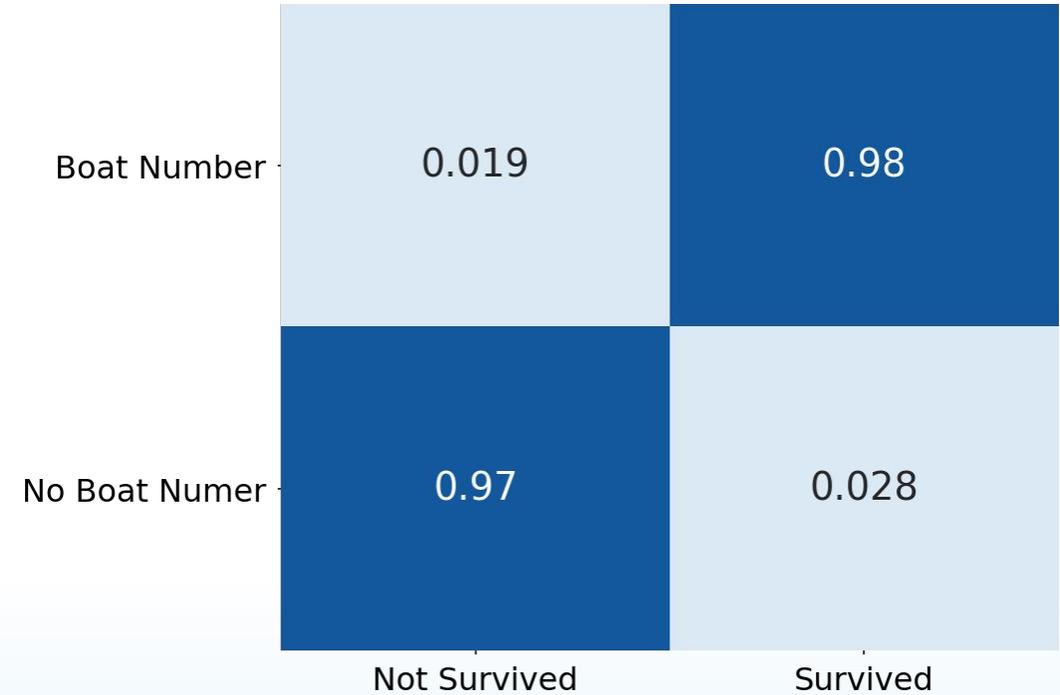
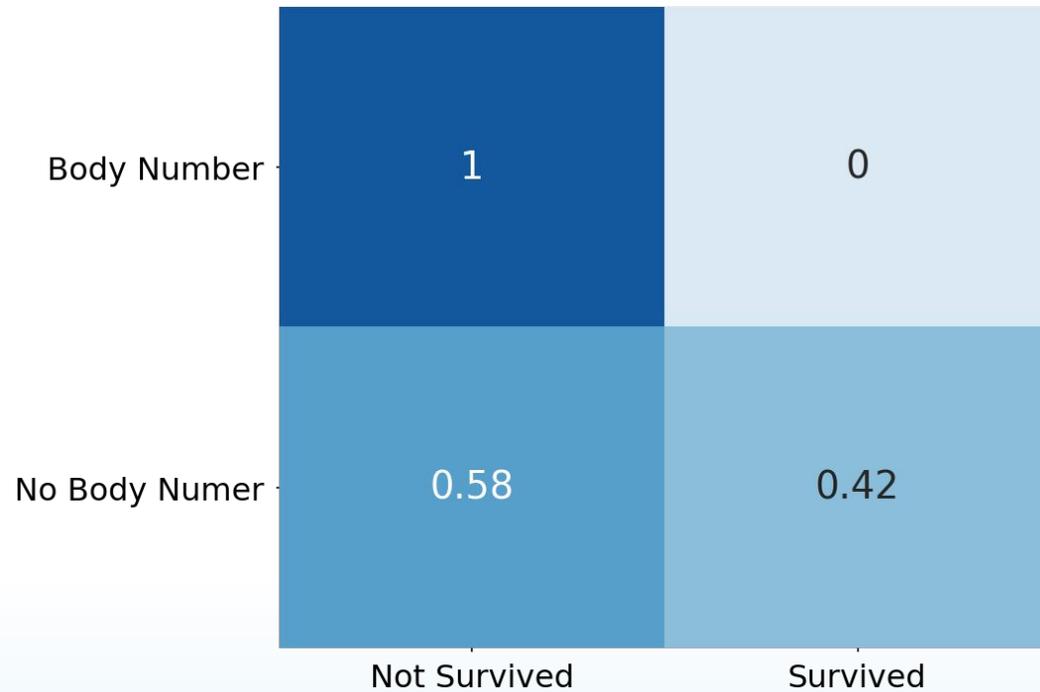
# A classic example

Predicting survival on the Titanic



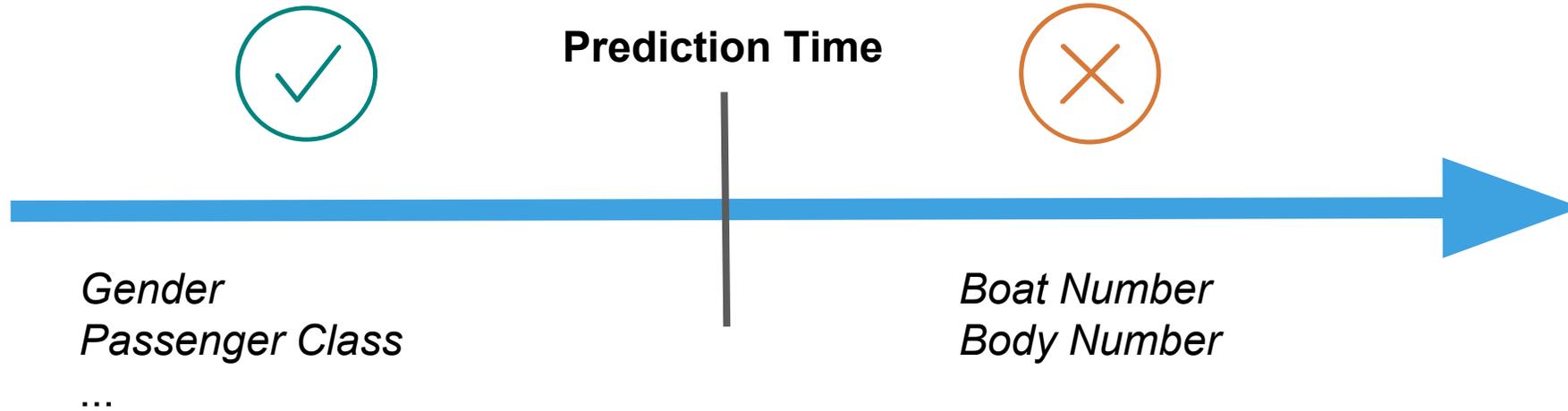
# A classic example

Predicting survival on the Titanic



# A classic example

Predicting survival on the Titanic



# A modern example

## Predicting lead conversion in Salesforce



### Before Conversion

The screenshot shows a Salesforce lead record for Ms. Lana Miller. The record is categorized as a 'Lead' and includes the following fields:

Name	Ms. Lana Miller	Email	lmiller@example.com
Company	Creativenet	Lead Status	New
Title	President	Days Since Last Activity	
Phone	(650) 455-3029	Lead Owner	Ely East
Segmentation			
Lead Source	Website	Industry	Consulting
Region	West	No. of Employees	1,800
Annual Revenue	\$60,000,000	Deal Value	
Address			
Address	101 Market Street San Francisco, CA 94105 United States		

A map below the address shows the location in the Financial District of San Francisco, near Rincon Park and Dragon's Gate.

### After Conversion

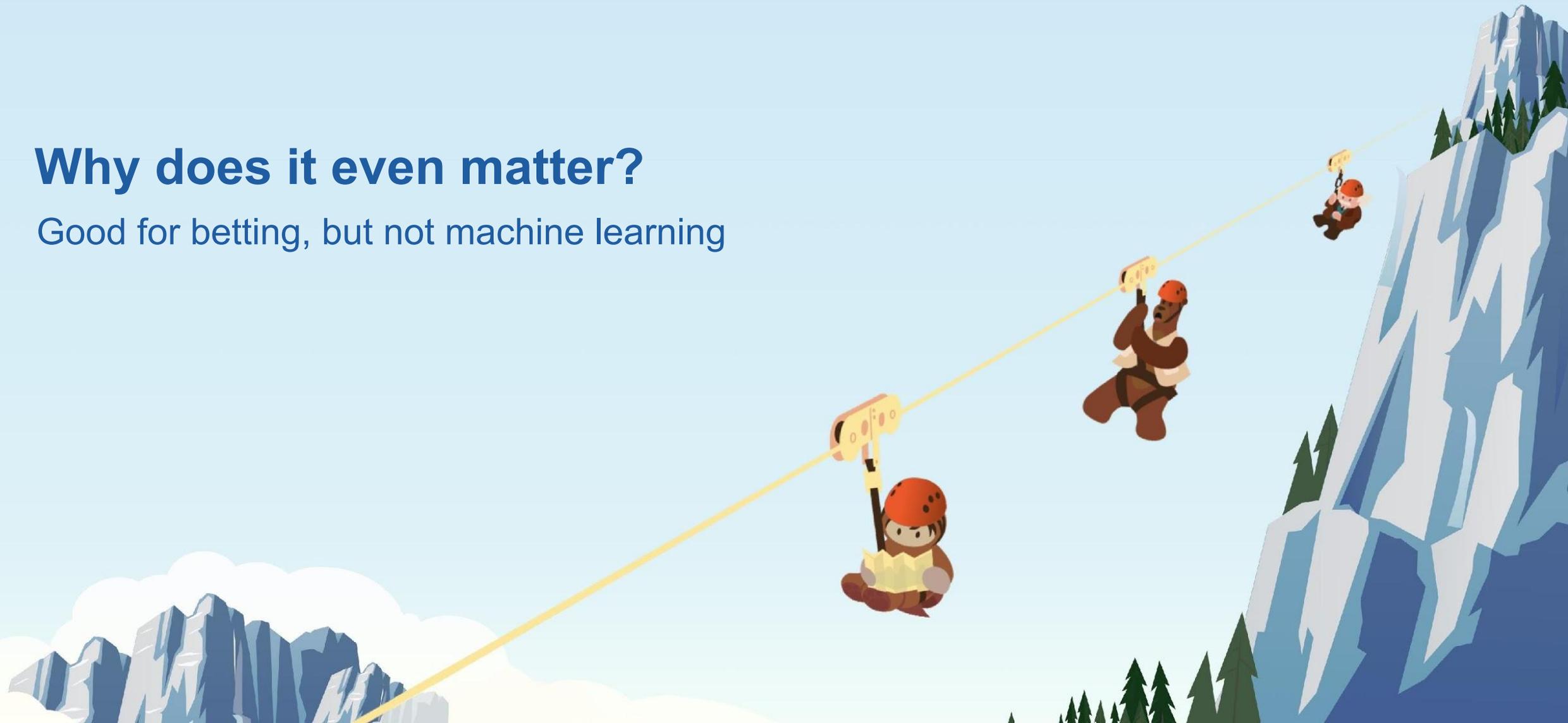
The screenshot shows the same Salesforce record for Ms. Lana Miller, but now it is categorized as an 'Account'. The record includes the following fields:

Name	Ms. Lana Miller	Email	lmiller@example.com
Company	Creativenet	Lead Status	New
Title	President	Days Since Last Activity	
Phone	(650) 455-3029	Lead Owner	Ely East
Segmentation			
Lead Source	Website	Industry	Consulting
Region	West	No. of Employees	1,800
Annual Revenue	\$60,000,000	Deal Value	\$1,000
Address			
Address	101 Market Street San Francisco, CA 94105 United States		

A red circle highlights the 'Deal Value' field, which now contains the value '\$1,000'. The map below the address is identical to the 'Before Conversion' screenshot.

# Why does it even matter?

Good for betting, but not machine learning

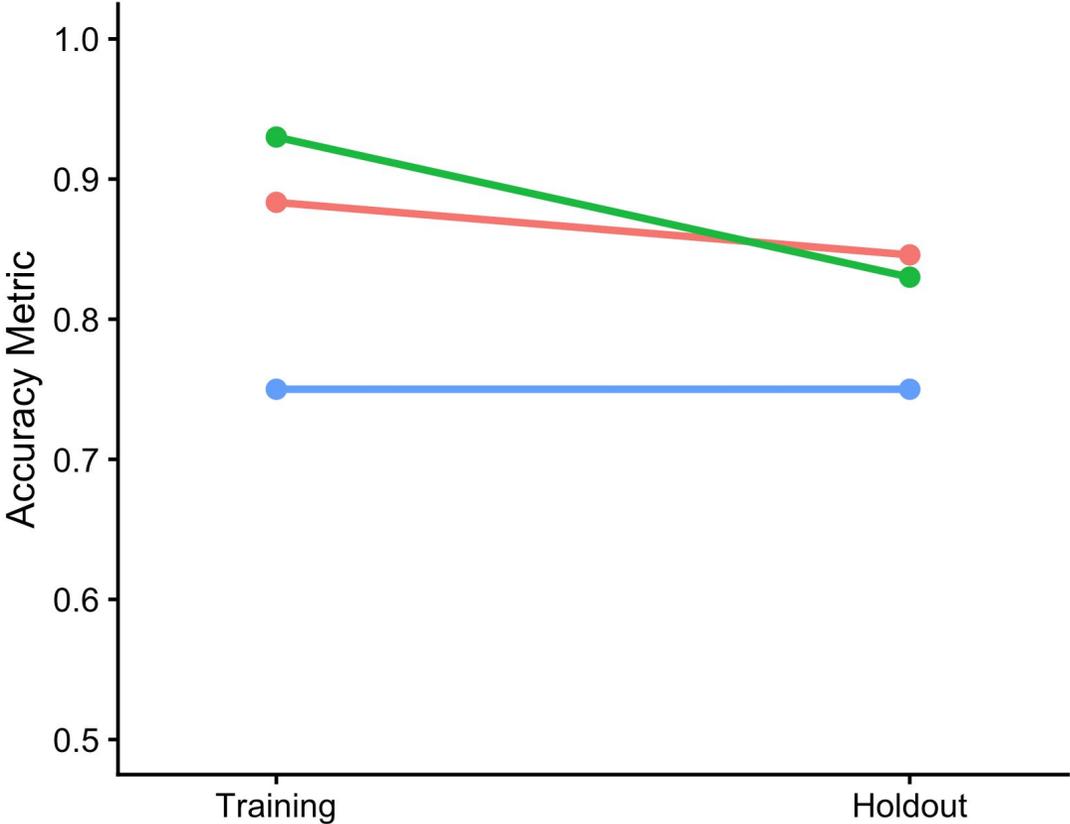


# Effect on model performance

## Traditional evaluation

Model relies on information not available at scoring time

- Model performance decreases for actual prediction
- Traditional evaluation pipeline is not sufficient

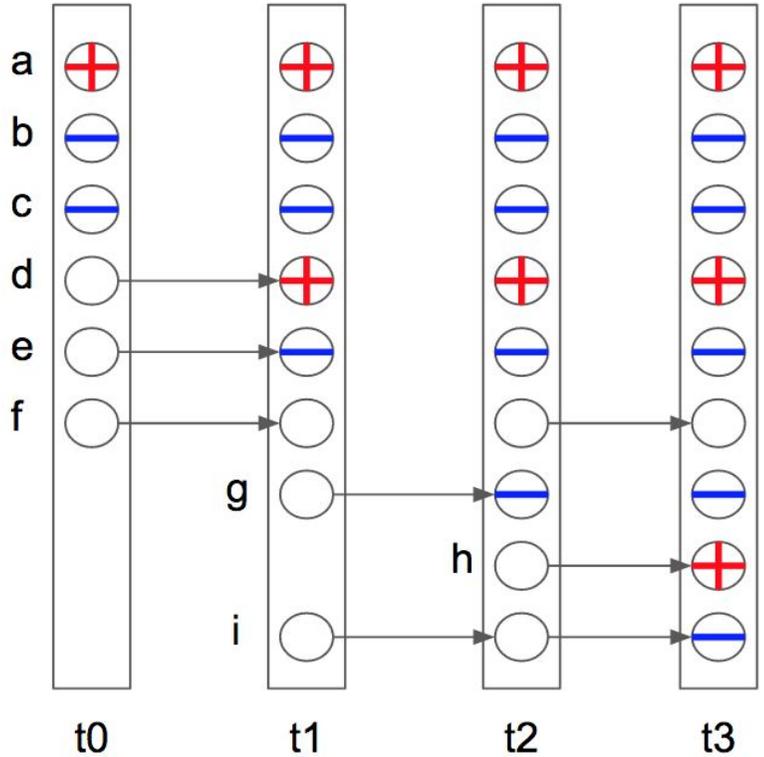
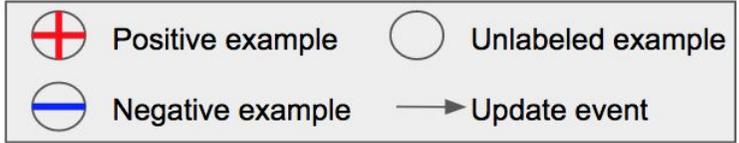
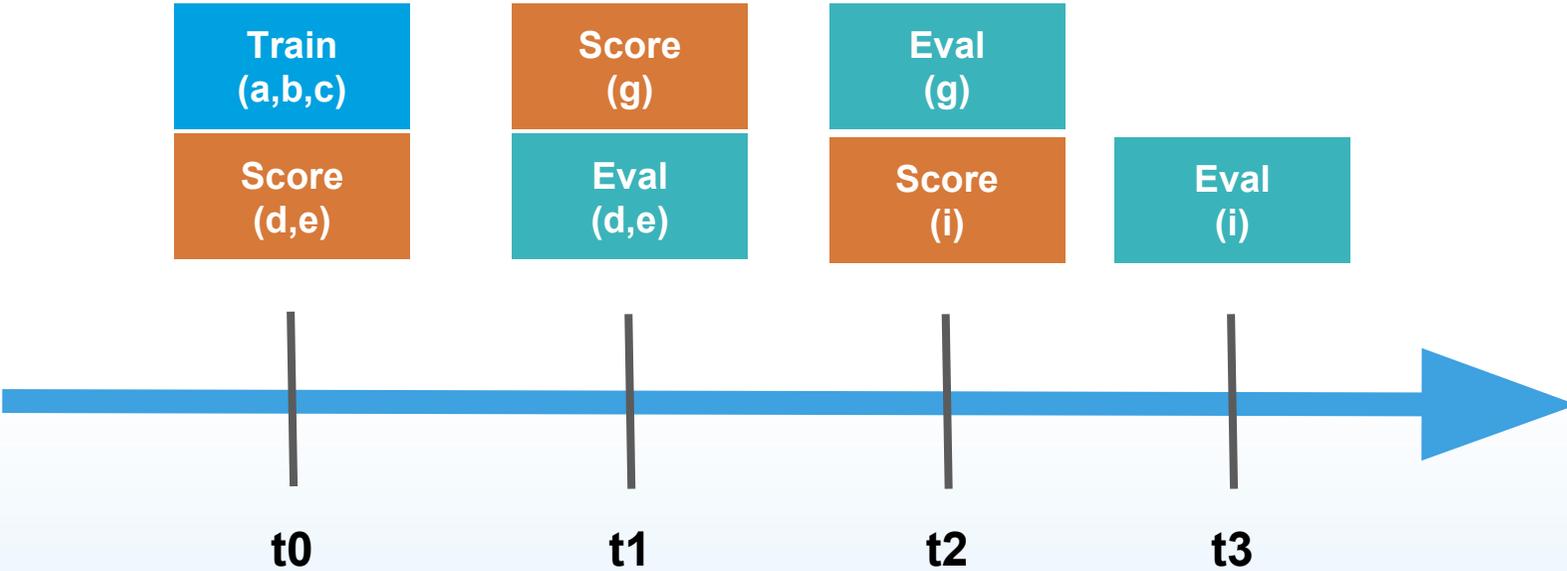


# Effect on model performance

## Time-based evaluation

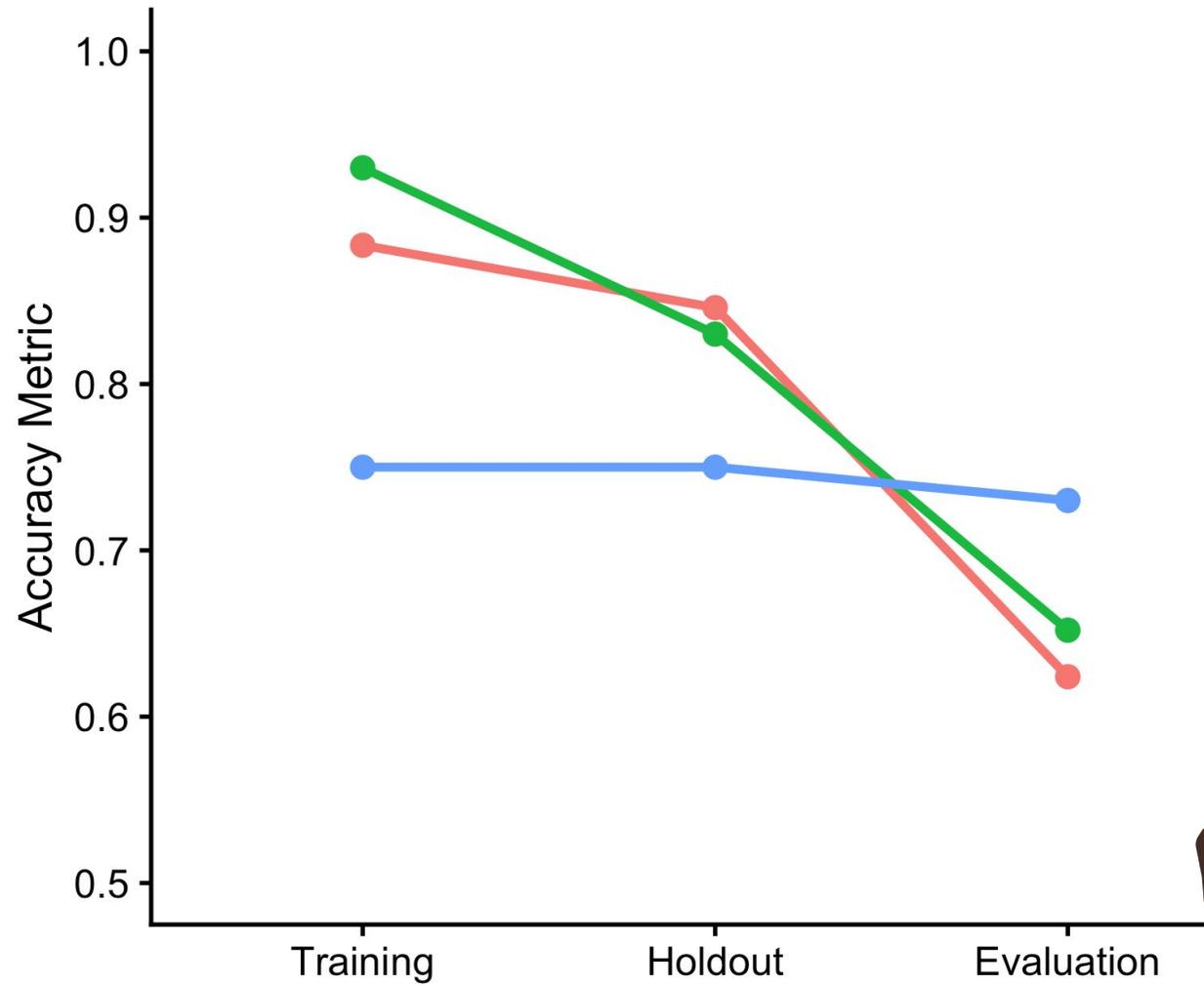
Need to treat each record separately

- Score and evaluate at different times



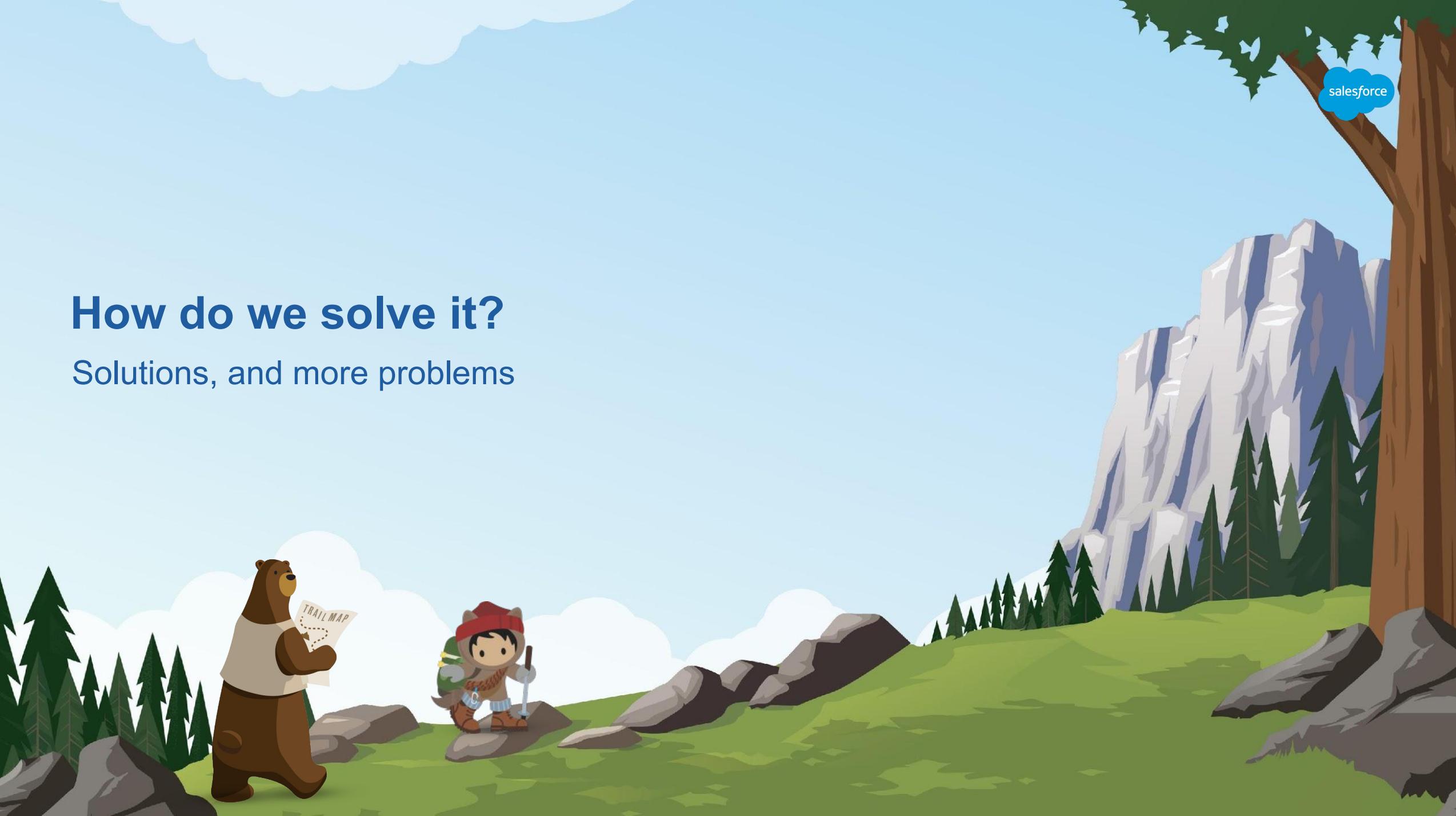
# Effect on model performance

Time based evaluation



# How do we solve it?

Solutions, and more problems



# A simple start



What are some problems with this data?

Id	Name	Address	Phone	ClosedBy	ReasonLost	Amount	<b>Converted</b>	...
342	...	...	...	32212	-	\$41k	True	
221	...	...	...	-	-	-	False	
098	...	...	...	86721	Unknown	-	False	
462	...	...	...	32212	-	\$23k	True	
140	...	...	...	-	Competitor	-	False	



# A simple start



What are some problems with this data?

- *ReasonLost* filled out means no conversion

Id	Name	Address	Phone	ClosedBy	ReasonLost	Amount	Converted	...
342	...	...	...	32212	-	\$41k	True	
221	...	...	...	-	-	-	False	
098	...	...	...	86721	Unknown	-	False	
462	...	...	...	32212	-	\$23k	True	
140	...	...	...	-	Competitor	-	False	



# A simple start



What are some problems with this data?

- *ReasonLost* filled out means no conversion
- *Amount* filled out means conversion

Id	Name	Address	Phone	ClosedBy	ReasonLost	Amount	Converted	...
342	...	...	...	32212	-	\$41k	True	
221	...	...	...	-	-	-	False	
098	...	...	...	86721	Unknown	-	False	
462	...	...	...	32212	-	\$23k	True	
140	...	...	...	-	Competitor	-	False	



# A simple start



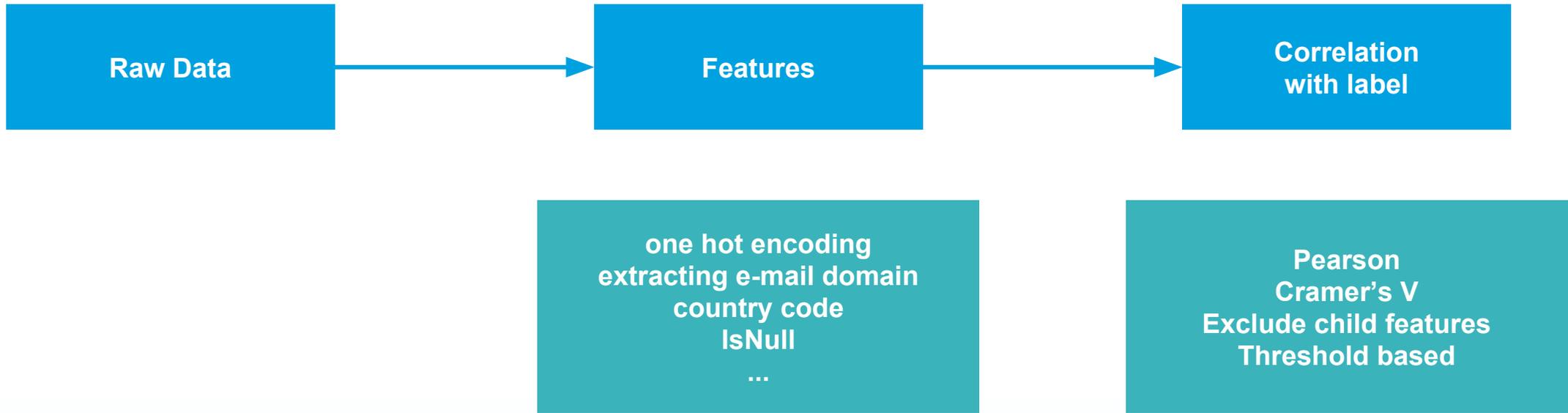
What are some problems with this data?

- *ReasonLost* filled out means no conversion
- *Amount* filled out means conversion
- *ClosedBy* filled out, more likely to have conversion

Id	Name	Address	Phone	ClosedBy	ReasonLost	Amount	Converted	...
342	...	...	...	32212	-	\$41k	True	
221	...	...	...	-	-	-	False	
098	...	...	...	86721	Unknown	-	False	
462	...	...	...	32212	-	\$23k	True	
140	...	...	...	-	Competitor	-	False	



# Catching features that are *too good*



# Does not solve everything



Data behaves in mysterious ways

Id	Name	Address	Phone	Expected Revenue	<b>Converted</b>	...
342	...	...	...	0		
221	...	...	...	0	False	
098	...	...	...	0		
462	...	...	...	15,000	True	
140	...	...	...	12,000	True	



# Does not solve everything



Data behaves in mysterious ways

- Default value is not always *null*

Id	Name	Address	Phone	Expected Revenue	Converted	...
342	...	...	...	0		
221	...	...	...	0	False	
098	...	...	...	0		
462	...	...	...	15,000	True	
140	...	...	...	12,000	True	



# Does not solve everything



## Data behaves in mysterious ways

- Default value is not always *null*
- A value  $> 0$  indicates conversion

Id	Name	Address	Phone	Expected Revenue	Converted	...
342	...	...	...	0		
221	...	...	...	0	False	
098	...	...	...	0		
462	...	...	...	15,000	True	
140	...	...	...	12,000	True	



# Does not solve everything



## Data behaves in mysterious ways

- Default value is not always *null*
- A value  $> 0$  indicates conversion
- Auto-bucketizing can catch these cases

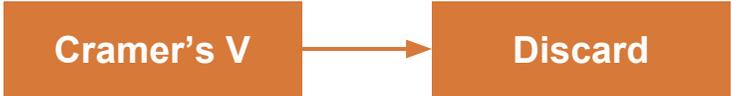
Id	Name	Address	Phone	Expected Revenue	Converted	...
342	...	...	...	0		
221	...	...	...	0	False	
098	...	...	...	0		
462	...	...	...	15,000	True	
140	...	...	...	12,000	True	



# Does not solve everything

## Data behaves in mysterious ways

- Default value is not always *null*
- A value > 0 indicates conversion
- Auto-bucketizing through decision tree can catch these cases



Id	Name	Address	Phone	Expected Revenue	Bucketized	Converted	...
342	...	...	...	0	[1, 0, 0]		
221	...	...	...	0	[1, 0, 0]	False	
098	...	...	...	0	[1, 0, 0]		
462	...	...	...	15,000	[0, 1, 0]	True	
140	...	...	...	12,000	[0, 1, 0]	True	

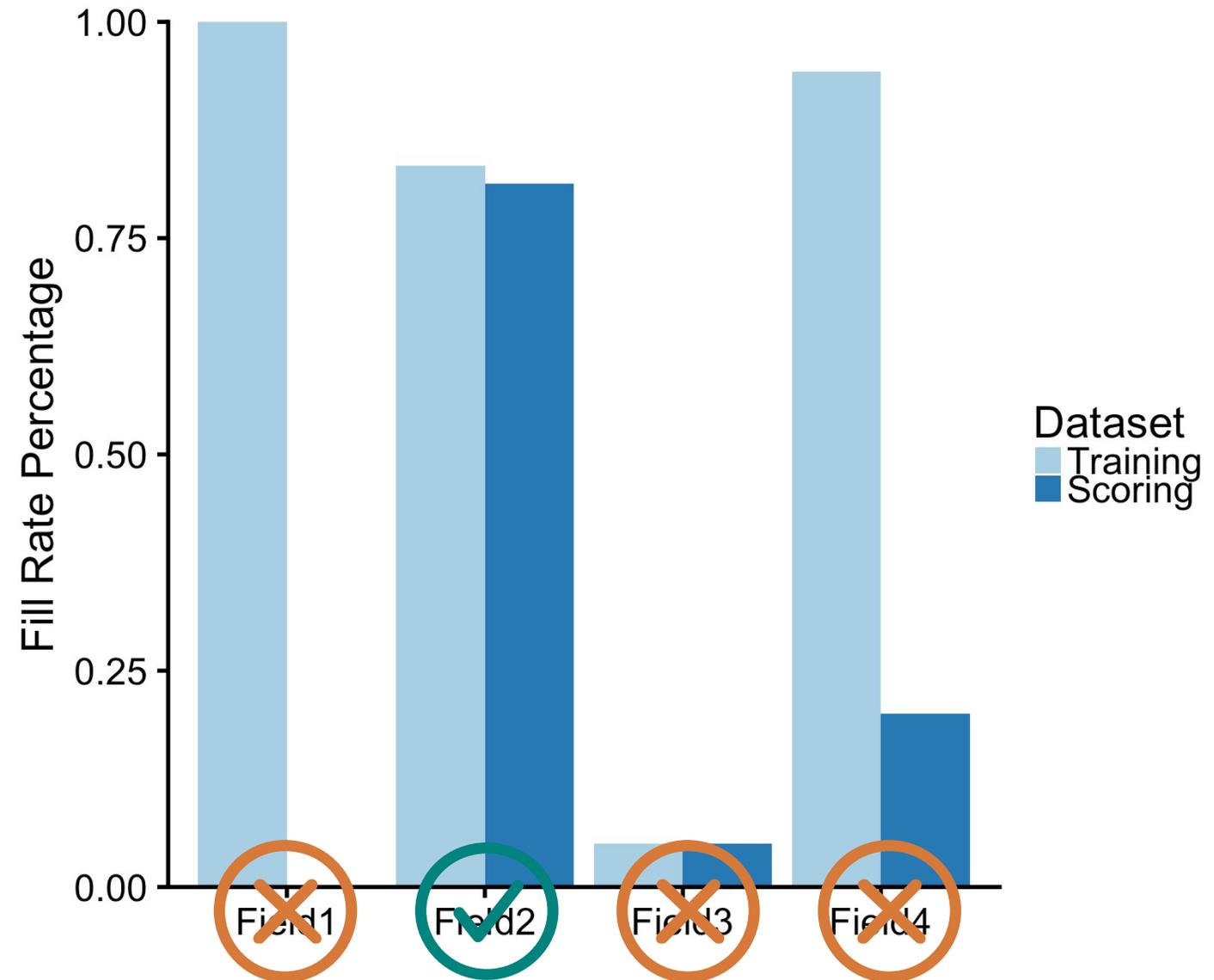


# Change over time

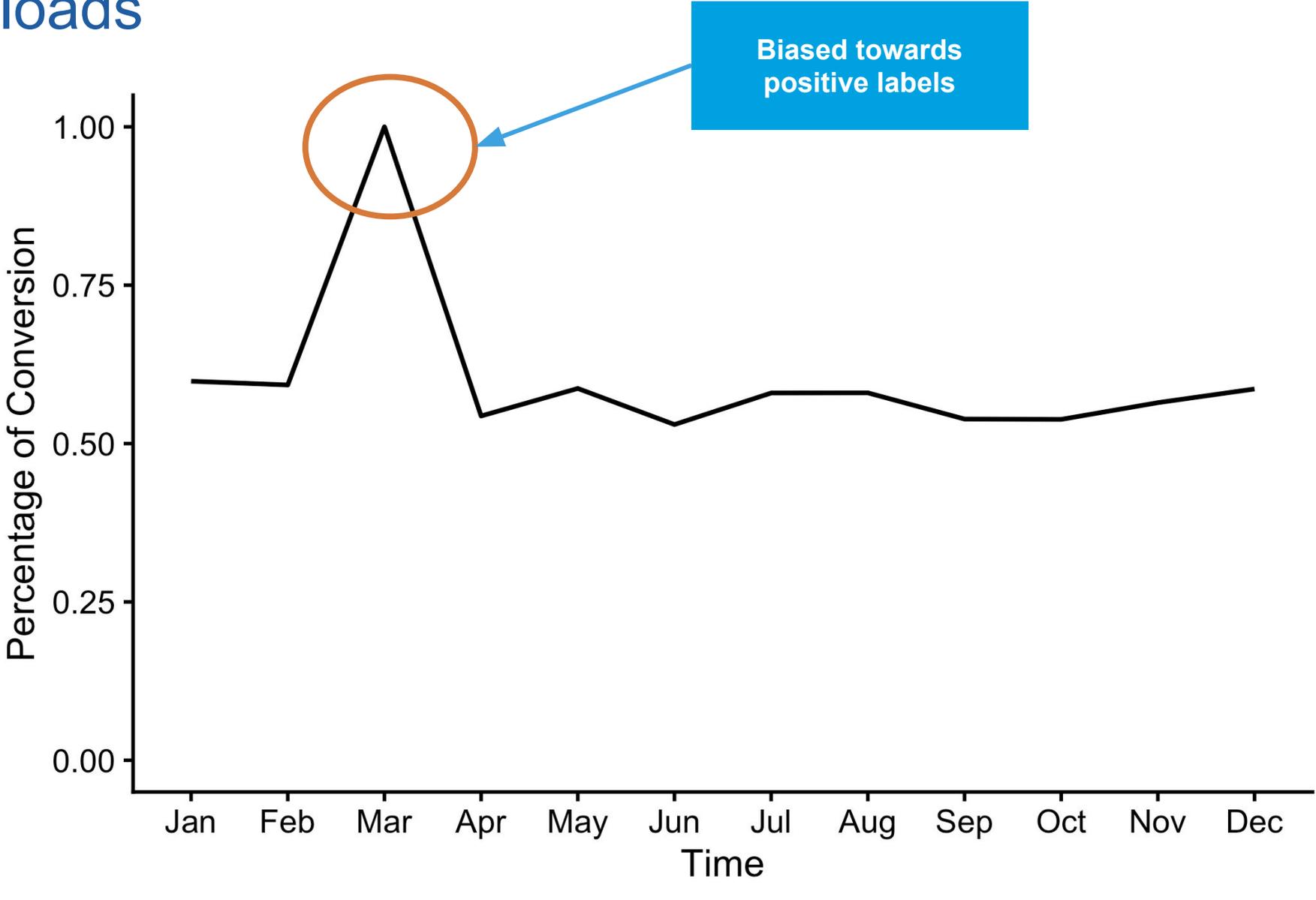


So far, we have only talked about data at the same point in time

- But training and scoring data are rarely produced at the same time
- Training data is historical, scoring data is more current



# Bulk uploads



# Criteria to exclude

## Low overall fill ratio

- No point in keeping a feature that is mostly null

## Big discrepancy between training and scoring

- Convert to probability distribution and compare with Jensen-Shannon Divergence

## Skewed dates and ratios

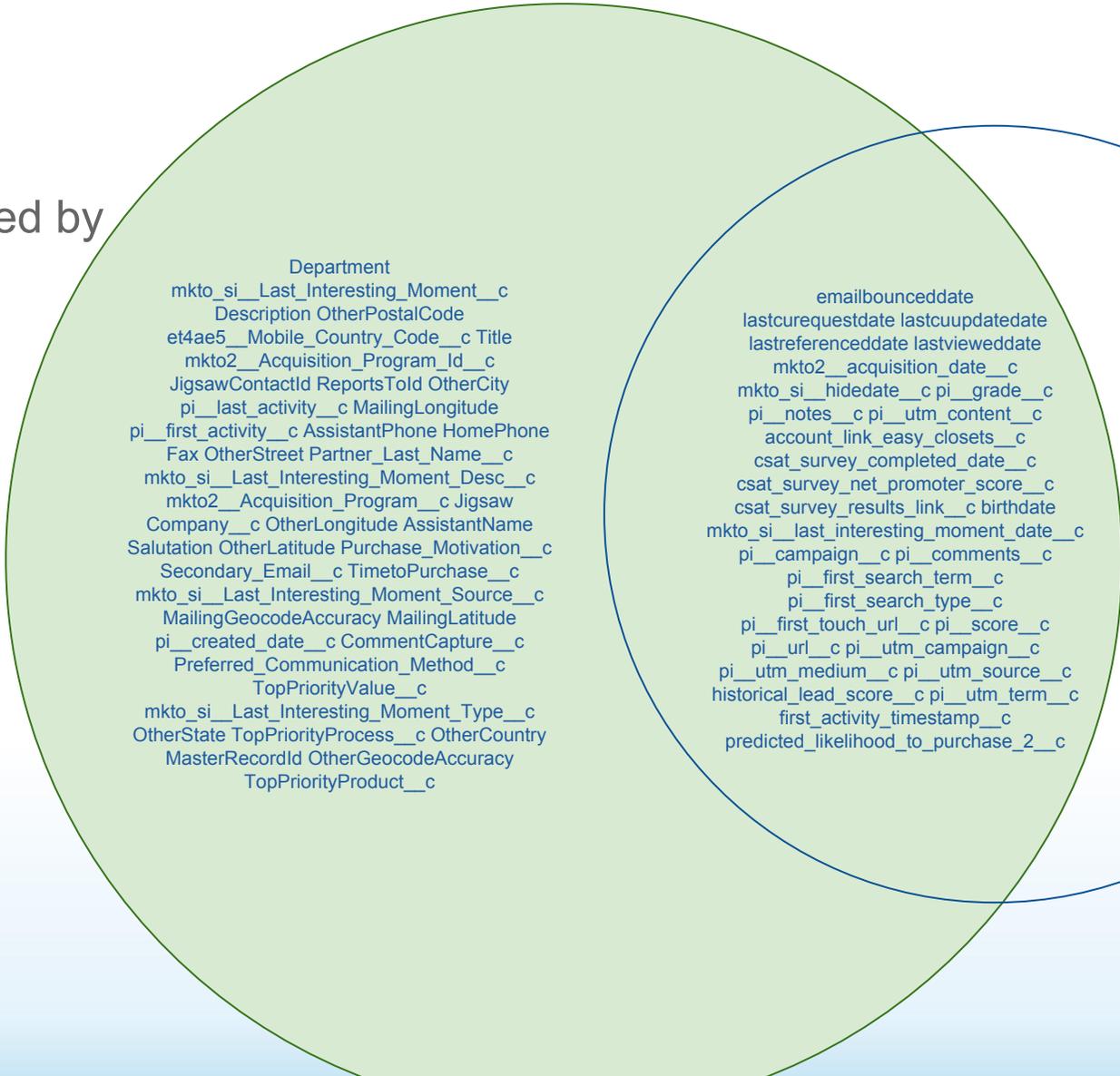
- Be careful about including date features that might be inherently biased



# AutoML vs Hand Tuning



Leakers removed by AutoML: 73



Department  
mkto\_si\_\_Last\_Interesting\_Moment\_\_c  
Description OtherPostalCode  
et4ae5\_\_Mobile\_Country\_Code\_\_c Title  
mkto2\_\_Acquisition\_Program\_Id\_\_c  
JigsawContactId ReportsTold OtherCity  
pi\_\_last\_activity\_\_c MailingLongitude  
pi\_\_first\_activity\_\_c AssistantPhone HomePhone  
Fax OtherStreet Partner\_Last\_Name\_\_c  
mkto\_si\_\_Last\_Interesting\_Moment\_Desc\_\_c  
mkto2\_\_Acquisition\_Program\_\_c Jigsaw  
Company\_\_c OtherLongitude AssistantName  
Salutation OtherLatitude Purchase\_Motivation\_\_c  
Secondary\_Email\_\_c TimetoPurchase\_\_c  
mkto\_si\_\_Last\_Interesting\_Moment\_Source\_\_c  
MailingGeocodeAccuracy MailingLatitude  
pi\_\_created\_date\_\_c CommentCapture\_\_c  
Preferred\_Communication\_Method\_\_c  
TopPriorityValue\_\_c  
mkto\_si\_\_Last\_Interesting\_Moment\_Type\_\_c  
OtherState TopPriorityProcess\_\_c OtherCountry  
MasterRecordId OtherGeocodeAccuracy  
TopPriorityProduct\_\_c

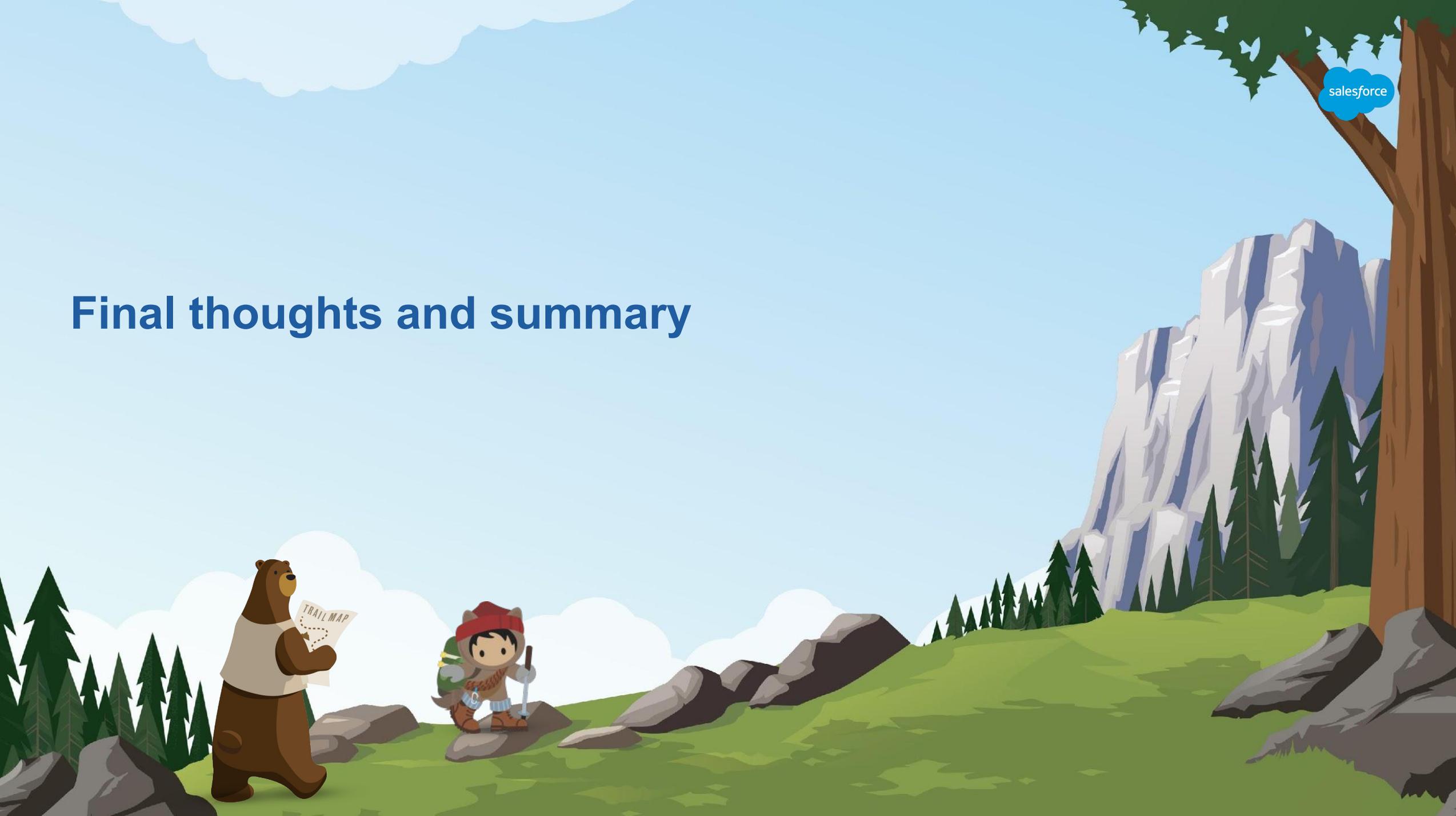
emailbounceddate  
lastcurequestdate lastcuupdatedate  
lastreferenceddate lastvieweddate  
mkto2\_\_acquisition\_date\_\_c  
mkto\_si\_\_hidedate\_\_c pi\_\_grade\_\_c  
pi\_\_notes\_\_c pi\_\_utm\_content\_\_c  
account\_link\_easy\_closets\_\_c  
csat\_survey\_completed\_date\_\_c  
csat\_survey\_net\_promoter\_score\_\_c  
csat\_survey\_results\_link\_\_c birthdate  
mkto\_si\_\_last\_interesting\_moment\_date\_\_c  
pi\_\_campaign\_\_c pi\_\_comments\_\_c  
pi\_\_first\_search\_term\_\_c  
pi\_\_first\_search\_type\_\_c  
pi\_\_first\_touch\_url\_\_c pi\_\_score\_\_c  
pi\_\_url\_\_c pi\_\_utm\_campaign\_\_c  
pi\_\_utm\_medium\_\_c pi\_\_utm\_source\_\_c  
historical\_lead\_score\_\_c pi\_\_utm\_term\_\_c  
first\_activity\_timestamp\_\_c  
predicted\_likelihood\_to\_purchase\_2\_\_c

Leakers removed by data scientist hand tuning: 42

best\_time\_to\_call\_date\_\_c  
total\_lead\_score\_\_c  
csat\_customer\_service\_survey\_disallowed\_\_c  
referral\_credit\_applied\_\_c  
referral\_days\_til\_purchase\_\_c  
predicted\_likelihood\_to\_purchase\_\_c  
createdbyid  
createddate  
lastactivitydate  
lastmodifieddate  
last\_activity\_date\_\_c  
systemmodstamp



# Final thoughts and summary



# Solve for all customers, not just one

Thresholds are tricky to choose

- What is a good feature and what is a bad leaker?

Easy to optimize for one model, but not for thousands

- Choosing a threshold that perfects one model, but makes hundreds worse is not good!

“Smart” decisions based on data shape preferred

- for example, auto-bucketizing - let the algorithm figure out a smart way

Lots of experimentation

- to learn heuristics that can be translated into algorithms



# Key Takeaways

## Enterprise data is very messy

- Often leads to hindsight bias/label leakage
- “Too good to be true” is a real problem

## Standard Machine Learning pipeline is not sufficient

- Time based evaluation is needed to know how your models are doing
- You cannot simply optimize for best model at training time

## Novel approaches needed to detect and remove leakage

- both on raw and transformed data
- choosing the right threshold to satisfy all customers



# TransmogrifAI



All the methods discussed here are part of our open-source library, *TransmogrifAI*

- Built on top of SparkML
- <https://github.com/salesforce/TransmogrifAI>

We are hiring more data scientists!



thank you

