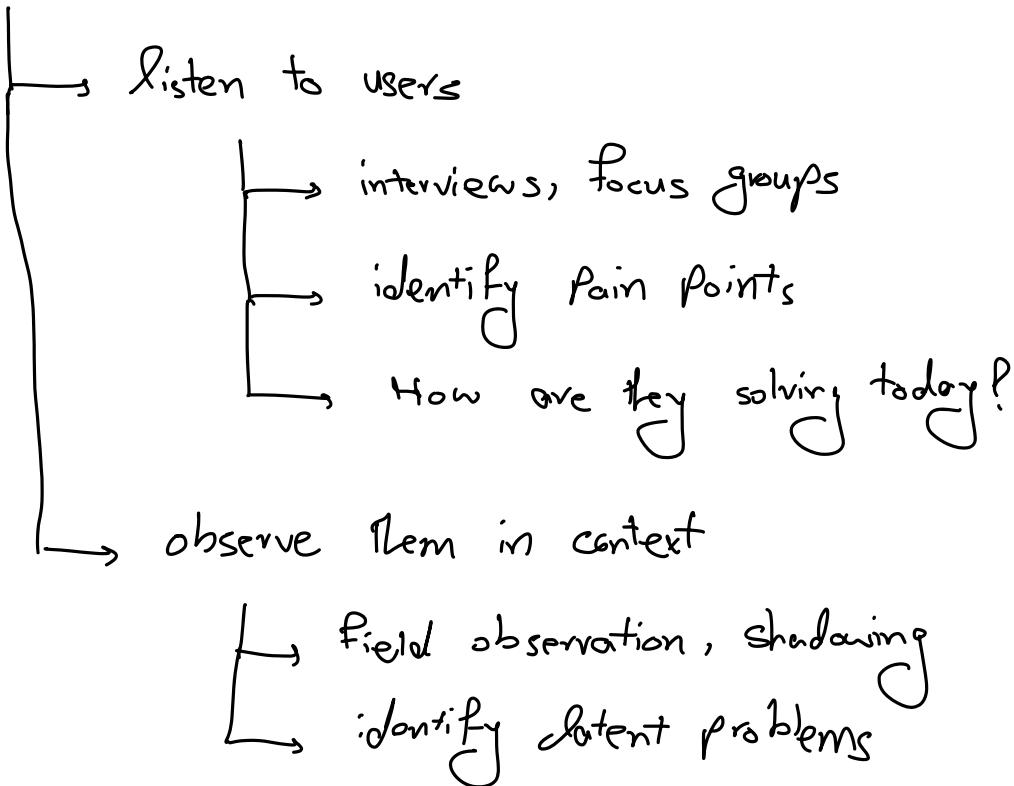


Managing Machine Learning Projects

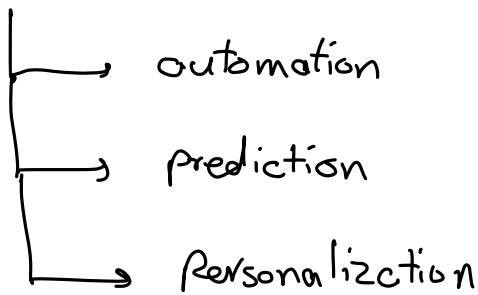
△ How to identify users?



△ what is validation process?

- ① formulate a hypothesis about the solution
- ② test the hypothesis with the users
- ③ Analyze your learnings
- ④ decide to continue or pivot
- ⑤ refine hypothesis & repeat

⚠ Machine Learning adds business value for your users



⚠ Heuristics are methods of solving problems using a simplified set of rules based on past experience

⚠ Heuristics vs ML

Heuristics	Machine Learning
easier to create and maintain	often better performing
minimal computation cost	can evolve with re-training
high interpretability	suitable for a wider range of problems

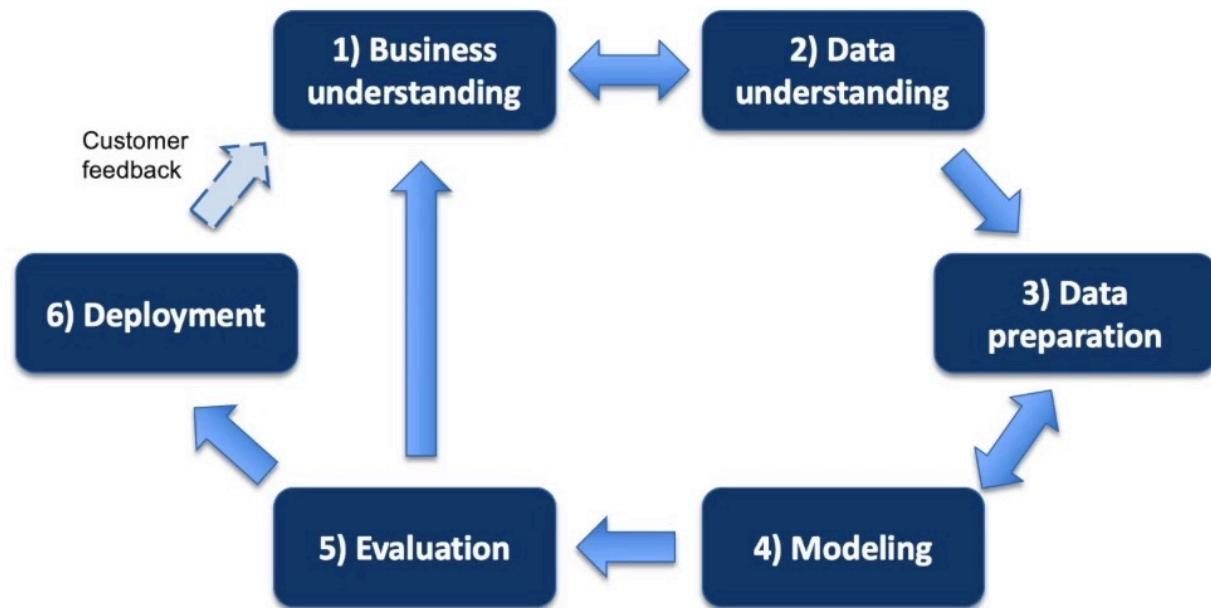
⚠ Relative to normal software projects, ML projects

- require a broader set of skills / team
- have higher technical risk
- Are harder to show progress
- require more ongoing support

⚠ Challenges of ML projects

- probabilistic rather than deterministic
- higher technical risk
- much more upfront work required
- often require change management

⚠ CRISP-DM process



⚠ Business Understanding

1.1 Define the problem	1.2 Define success	1.3 Identify factors
<ul style="list-style-type: none">• Target user• Write the problem statement• Why it matters• How is it solved today?• Gaps in current state	<ul style="list-style-type: none">• Quantify the expected business impact• Identify constraints• Translate impact into metrics – outcome & output metrics• Define success targets for metrics	<ul style="list-style-type: none">• Gather domain expertise• Identify potentially relevant factors

Data Understanding

2.1 Gather data	2.2 Validate data	2.3 Explore the data
<ul style="list-style-type: none">Identify data sources for each factorLabel dataCreate features	<ul style="list-style-type: none">Quality control dataResolve data issues – missing, erroneous, outliers	<ul style="list-style-type: none">Statistical analysis and visualizationDimensionality reductionIdentify relationships & patterns

Data Preparation

3.1 Split data	3.2 Determine feature set	3.3 Prepare for modeling
<ul style="list-style-type: none">Split data for training and test	<ul style="list-style-type: none">Feature engineeringFeature selection	<ul style="list-style-type: none">Encoding categorical featuresScale/standardize dataResolve class imbalance

Modeling

4.1 Model selection	4.2 Model tuning
<ul style="list-style-type: none">Evaluate algorithms via cross-validationDocumentation and versioning	<ul style="list-style-type: none">Hyperparameter optimizationDocumentation and versioningModel re-training

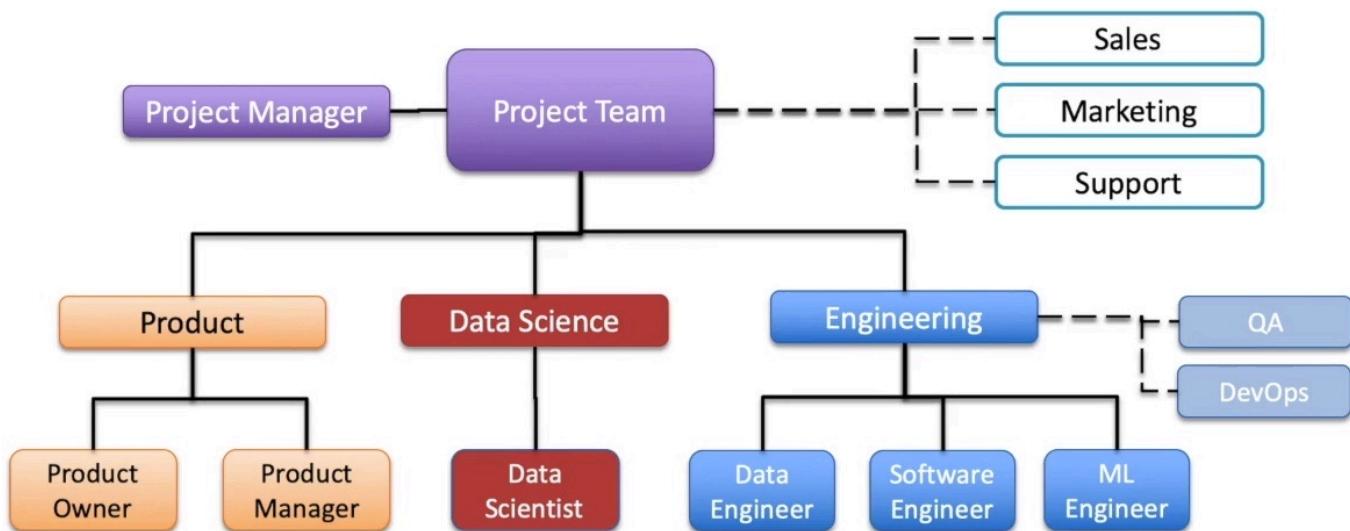
⚠ Evaluation

5.1 Evaluate results	5.2 Test solution
<ul style="list-style-type: none">• Model scoring on test set• Interpretation of model outputs and performance	<ul style="list-style-type: none">• Software unit & integration tests• Model testing – unit tests, directional expectation• User tests

⚠ Deployment

6.1 Deploy	6.2 Monitor
<ul style="list-style-type: none">• API framework• Product integration• Scaling infrastructure• Security• Software deployment process	<ul style="list-style-type: none">• Model performance monitoring• Model retraining

⚠ Typical team roles

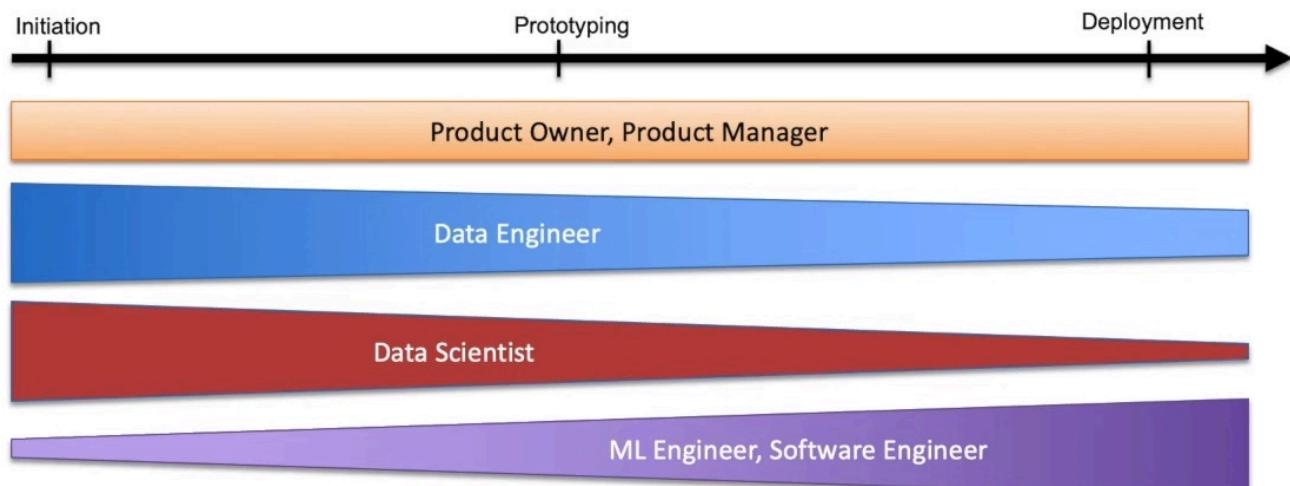


⚠ Data Scientist vs. ML Engineer

Data Scientist	ML Engineer / MLOps
statistical / data science background plus programming skills & domain expertise	computer science or eng. background plus ML training
gather, process & derive insights from data	develop production data pipelines and ML system
determining of ML approach and prototyping	work with software engineering & DevOps on model integration and deployment

⚠ Involvement over project cycle

Project lifecycle



⚠ Agile approach to ML

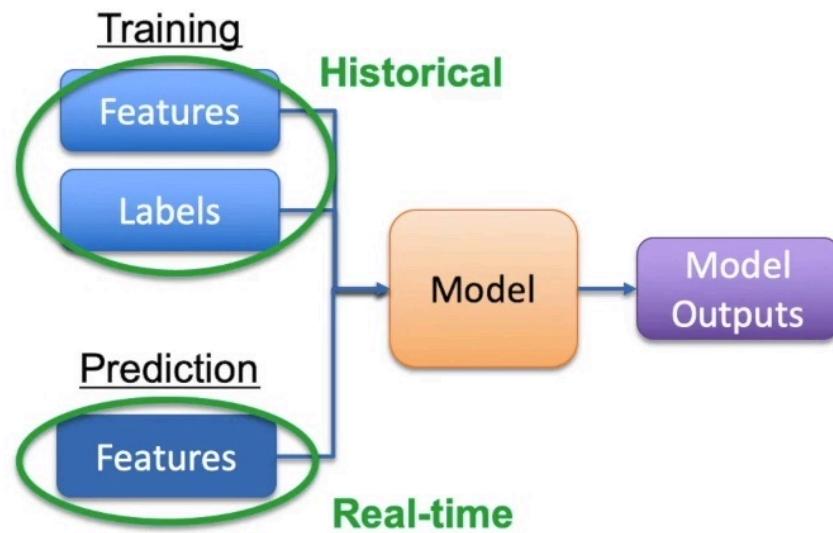
↳ sequence of iterative experiments

- explore or hypothesis
- build it, using more of CRISP-DM each time
- observe it in action, get feedback
- analyze results and repeat

⚠ Out come metrics vs output metrics

outcome metrics	output metrics
refers to the desired business impact on your organization or for your customer	refers to the desired output from the model & will set after setting the desired outcome
stated in terms of the expected impact (\$)	measured in terms of a model performance metric
does not contain model performance metrics or other technical metrics	typically not communicated to the customer

A Model training and evaluation



A factors that influence data requirements

- ↳ number of features
- ↳ complexity of feature-target relationships
- ↳ data quality (missing and noisy data)
- ↳ desired model performance

A Data can have several issues

- ↳ missing data
- ↳ anomalous data
- ↳ incorrectly mapped data

⚠ types of missing data

	Missing Completely at Random	Missing at Random	Missing Not at Random
Description	No pattern in missing data or association to values of other attributes	Probability of missing-ness relates to another feature of the data	Probability of missing-ness relates to values of the feature itself
Example	Power outages of sensors	Males are less likely to answer survey questions about depression	Purchased item ratings skew towards people who hated the product
Potential for Bias	Low	High	High

⚠ Dealing with missing data

- drop (remove)
- flag (+treat as special value of feature)
- replace (with mean value / median)
- backfill or forward-fill
- infer it

⚠ Outliers

- points which fall far from the rest, either in a feature value or the target value

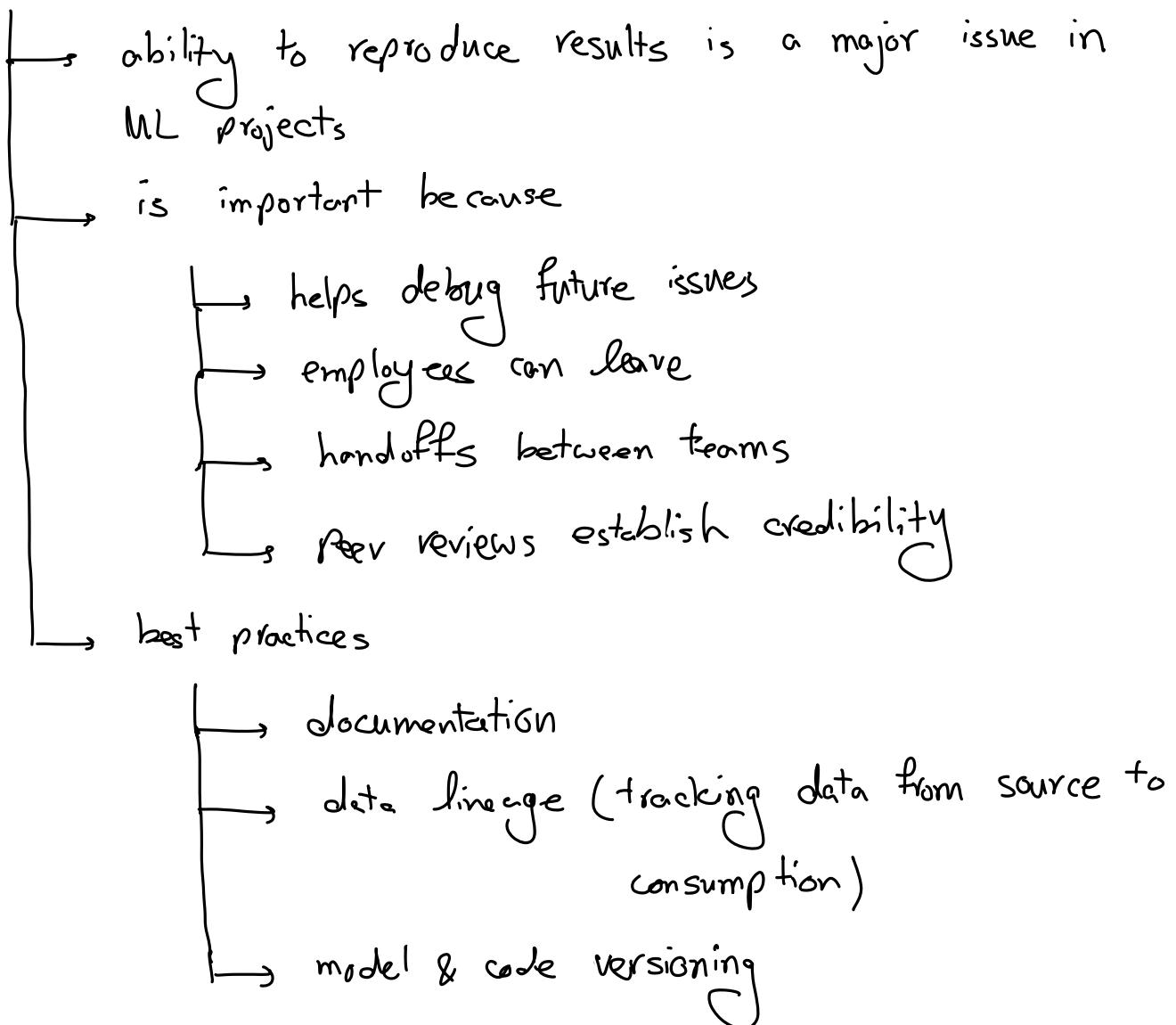
⚠ Exploratory data analysis (EDA)

- helps us catch issues in our data & understand it better

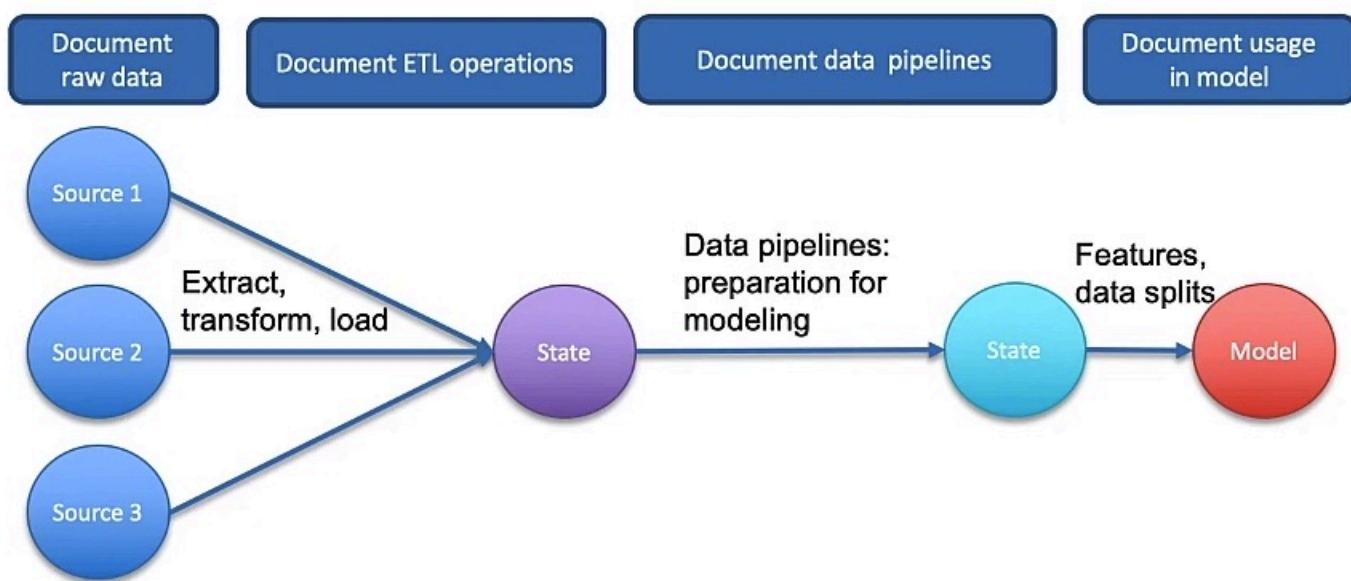
⚠ Feature Selection Methods

	Filter Methods	Wrapper Methods	Embedded Methods
Description	<ul style="list-style-type: none">Statistical tests which rely on characteristics of the data only	<ul style="list-style-type: none">Train model on subsets of features	<ul style="list-style-type: none">Extracts features which contribute most to training of a model
Pros & Cons	<ul style="list-style-type: none">Computationally inexpensiveOften used before modeling to remove irrelevant features	<ul style="list-style-type: none">Computationally very expensiveOften unfeasible for real-world modeling	<ul style="list-style-type: none">Leverage model training with minimal additional computation

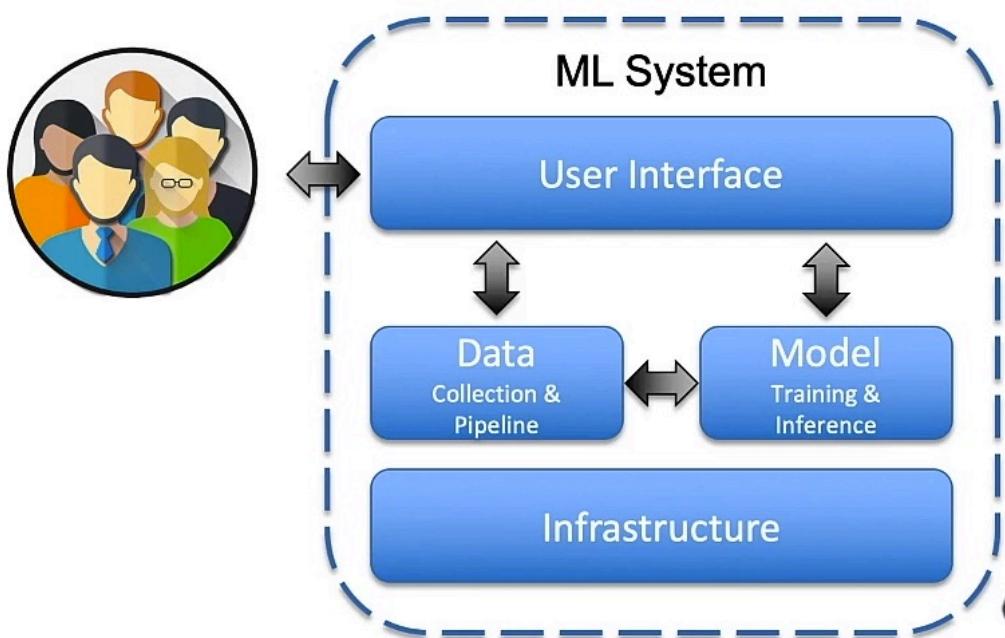
⚠ Reproducibility



⚠ Data lineage



⚠ what is a ML system?



⚠ There are several system design decisions which impact the choice of technology

- ① cloud vs Edge
- ② offline learning vs online learning
- ③ Batch predictions vs online predictions

⚠ Cloud vs Edge

	Cloud ML	Edge ML
Description	Computations done on cloud and result delivered to end device	Computations done directly on device (phone, sensor, etc)
Requirements	Network connectivity	Sufficient compute power, memory
Benefits	High throughput	Low latency, privacy, no need for connectivity
Examples	<ul style="list-style-type: none"> Chatbots Demand prediction 	<ul style="list-style-type: none"> Quality control Autonomous driving

⚠ offline vs online Models

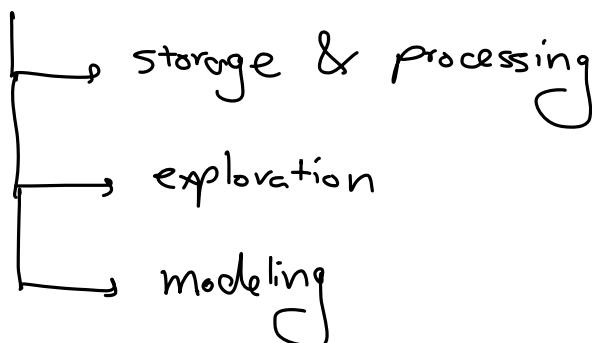
	Scheduled	Real-time
Model re-training	Offline learning	Online learning
Prediction	Batch prediction	Online prediction

	Offline learning	Online learning
Description	Model re-training done on a schedule (weeks/months) using datapoints in many iterations	Continual re-training as new data arrives (mins/hours) using each new datapoint once
Benefits	<ul style="list-style-type: none"> Easier to implement in production Easier to evaluate 	<ul style="list-style-type: none"> Handles big data Real-time adaptation to changing environment
Challenges	<ul style="list-style-type: none"> Slower to adapt to changes in environment or data distribution 	<ul style="list-style-type: none"> Harder to implement & evaluate performance
Examples	<ul style="list-style-type: none"> Most current applications 	<ul style="list-style-type: none"> Flagging spam in social media

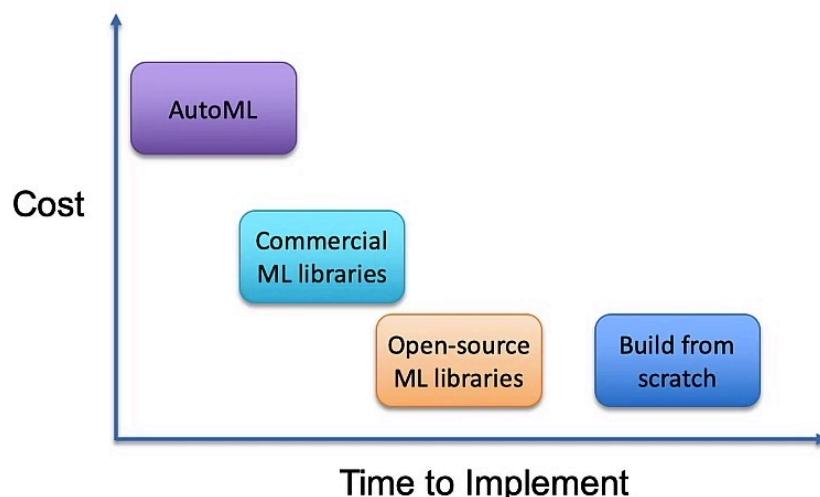
⚠ Batch vs online prediction

	Batch prediction	Online prediction
Description	Generate predictions on batch of observations on a recurring schedule	Real-time predictions generated upon request
Benefits	<ul style="list-style-type: none">Leverage more efficient operations and technologiesEasier monitoring of drift	<ul style="list-style-type: none">Predictions available immediately
Challenges	<ul style="list-style-type: none">Predictions not immediately available for new data	<ul style="list-style-type: none">Minimizing latencyMonitoring of model drift
Examples	<ul style="list-style-type: none">Recommendation systemsDemand prediction	<ul style="list-style-type: none">Translation appAutonomous vehicles

⚠ Working with Big Data has unique challenges

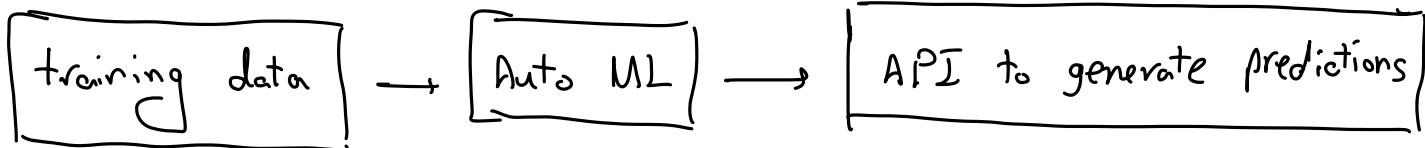


⚠ ML technology options



A What is Auto ML?

↳ enables developers with limited ML expertise to quickly build models with little / no code



A ML system failure

↳ training - serving skew

↳ mismatch between the training data and the input data while in production

↳ excessive latency

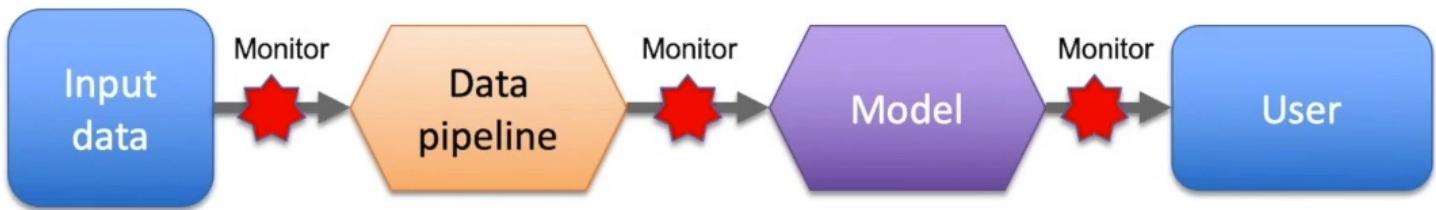
↳ data drift

↳ model is trained on a static training set, but the environment changes over time

↳ concept drift

↳ distributions of data may stay same, but the patterns that the model learned no longer apply

⚠ ML system monitoring



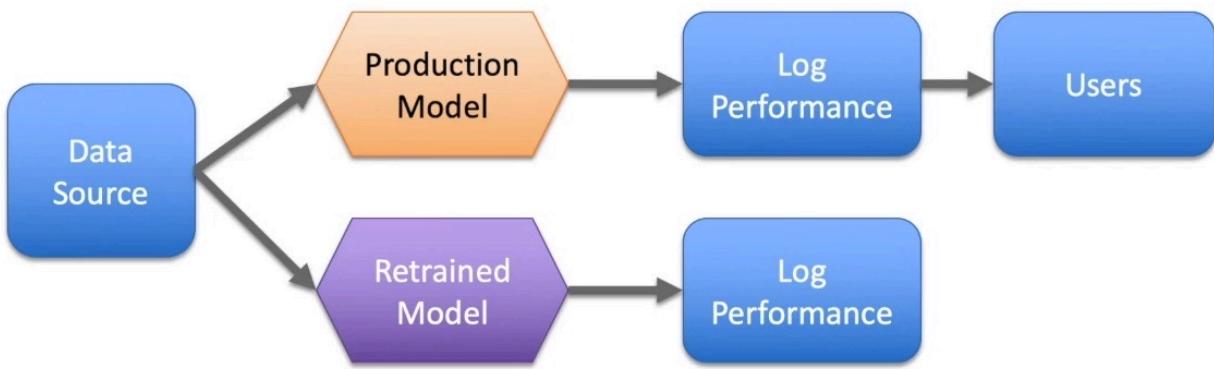
⚠ Retraining vs. updating models

Retraining	Updating
Using data collected since previous training to re-train model weights	Uses new data to re-do the modeling process
enables model to adapt to changes to environment	Allows for adjustments to model form
can be done on a schedule or triggered	can lead to identification of higher quality model

⚠ why version models?

- evaluate performance across iterations
- track dependencies - ensure reproducibility
- facilitates collaboration
- build in rollback capability
- enable testing multiple models

⚠ shadow releasing



If retrained model performance exceeds production model, move retrained model into production.

⚠ champion - challenger testing

