

Applied Data Analysis (CS401)



Lecture 9
Applied
machine learning
2018/11/15



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Robert West



Announcements

- Course [evaluations](#) are being collected -- please evaluate us!
- HW 4 due next week, Wed, Nov 21, 23:59
- HW 3 peer reviews due today 23:59
- Project milestone 2 due on Sun, Nov 25, 23:59
- Tomorrow's lab session:

Office hours HW 4 + project [[FAQ](#)]

Teo Stocco (ADAlumnus who founded and sold a startup, now at Smood)



Give us feedback on this lecture here:

<https://go.epfl.ch/ada2018-lec9-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- What's your favorite messaging system?
- ...

Why an extra class on applied ML?

Machine Learning that Matters

Kiri L. Wagstaff

KIRI.L.WAGSTAFF@JPL.NASA.GOV

Jet Propulsion Laboratory, California Institute of Technology, 4800 Oak Grove Drive, Pasadena, CA 91109 USA

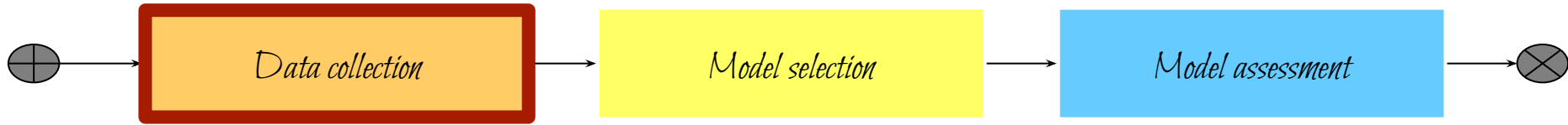
[link](#)

Classic ML
class

ADA

It is easy to sit in your office and run a Weka (Hall et al., 2009) algorithm on a data set you downloaded from the web. It is very hard to identify a problem for which machine learning may offer a solution, determine what data should be collected, select or extract relevant features, choose an appropriate learning method, select an evaluation method, interpret the results, involve domain experts, publicize the results to the relevant scientific community, persuade users to adopt the technique, and (only then) to truly have made a difference (see Figure 1). An ML researcher might well feel fatigued or daunted just contemplating this list of activities. However, each one is a necessary component of any research program that seeks to have a real impact on the world outside of machine learning.

Classification pipeline



Data collection

The first step is collecting data related to the classification task.

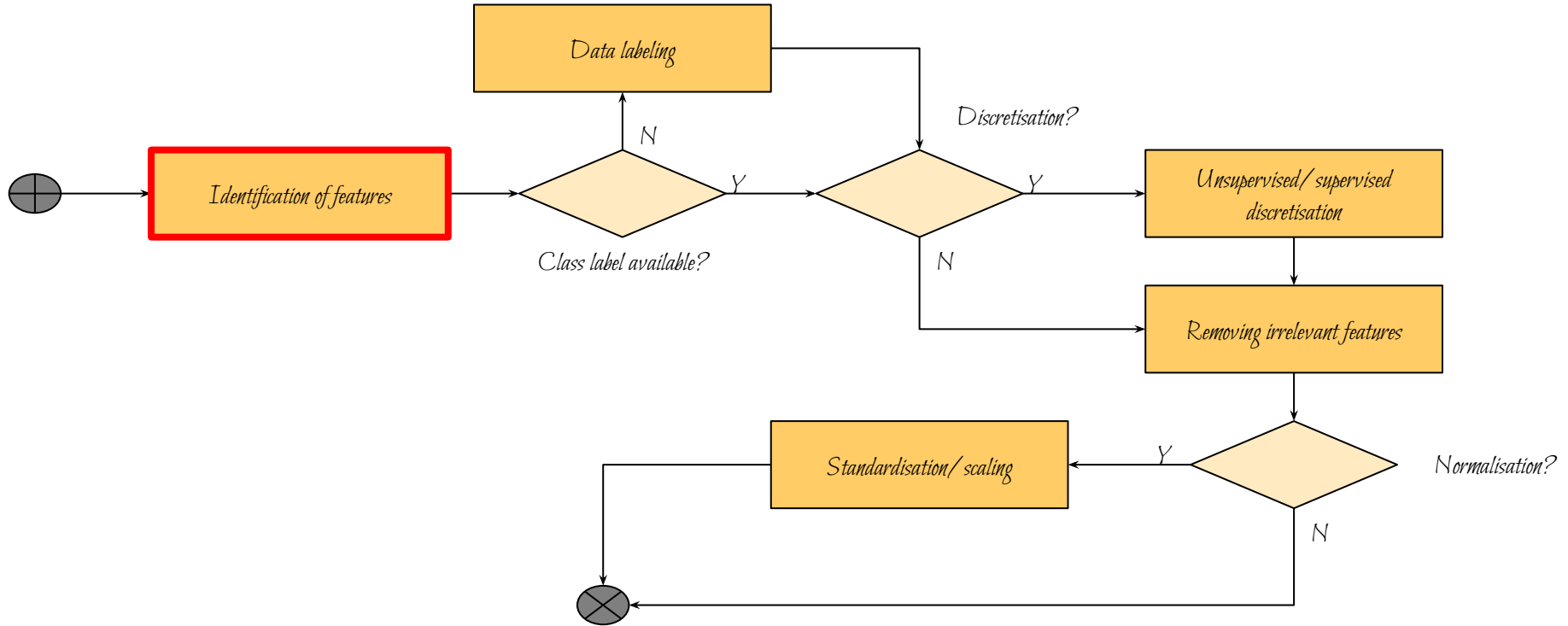
- Definition of the attributes (or features) that describe a data item and the class label.

Domain knowledge is needed.

What if assigning the class label would be too time-consuming or even impossible?

→ Unsupervised methods (e.g., clustering); cf. next lecture!

Data collection



Features

Different types of features [\[more\]](#)

- Continuous (e.g., height, temperature ...)
- Ordinal (e.g., “agree”, “don’t care”, “disagree” ...)
- Categorical (e.g., color, gender ...)

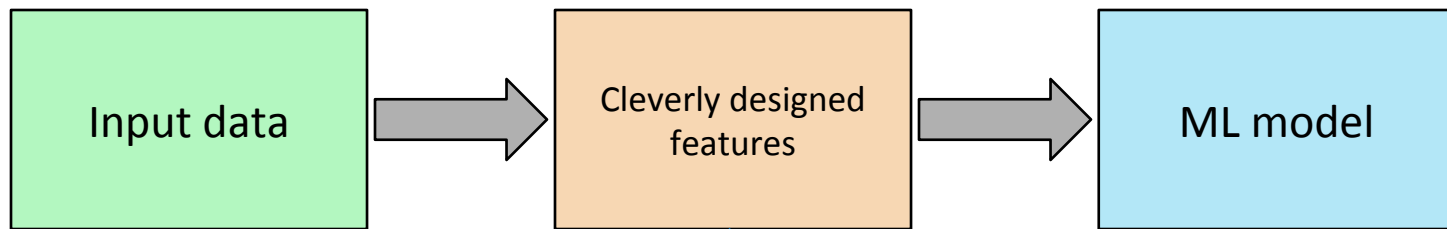
New features can be generated from **simple stats**

- *Feature engineering* is considered a form of art, therefore it is sometimes useful to find what other people did for similar problems

Some classifiers require categorical features => **Discretization**

ML before 2012*

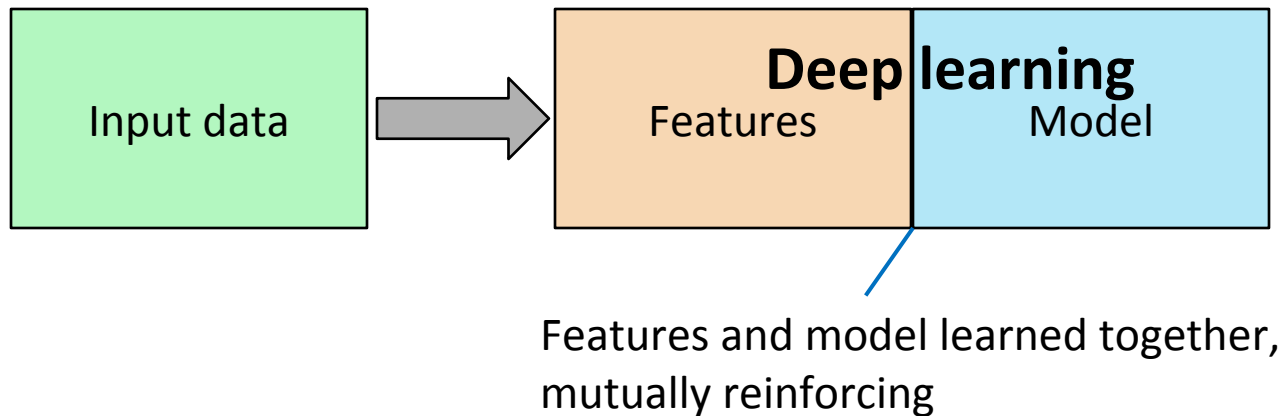
(but still very common today)



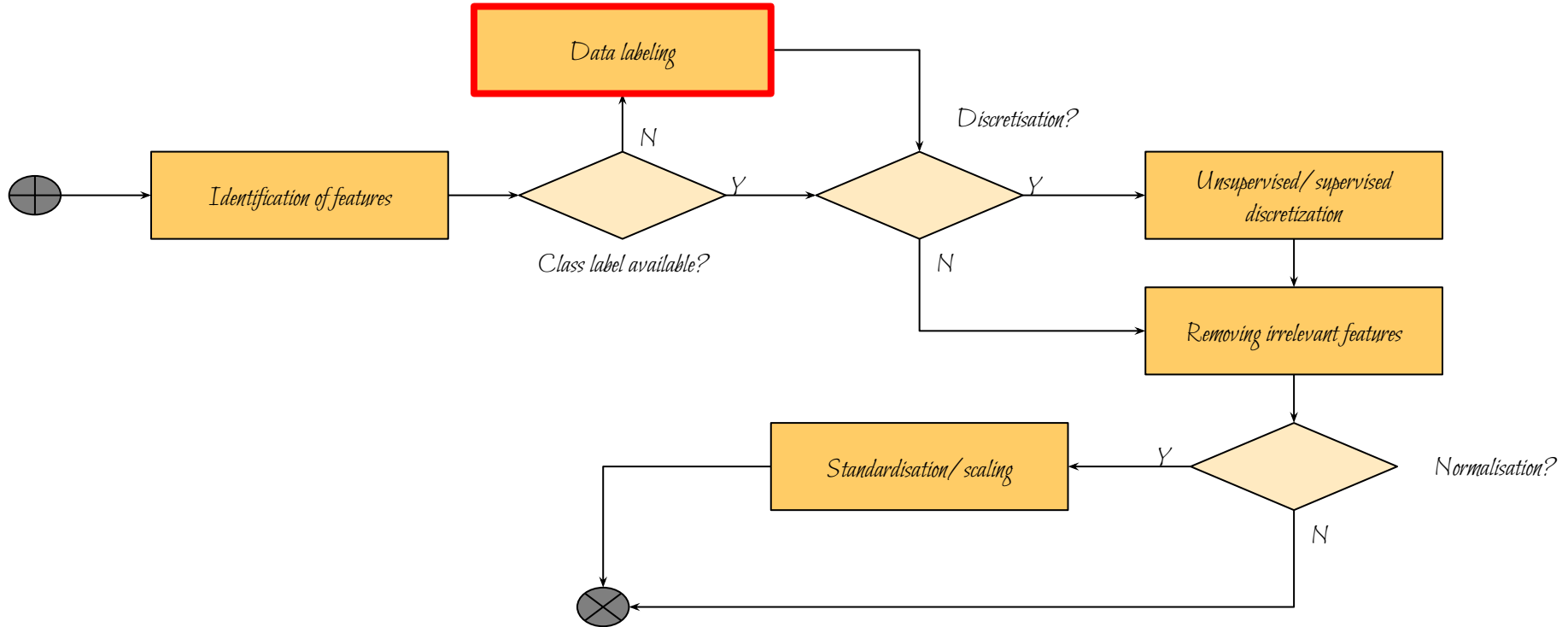
Most of the “heavy lifting” in here.
Final performance only as good as the
feature set.

* Before publication of Krizhevsky et al.’s ImageNet CNN [paper](#)

A typical ML approach after 2012



Data collection

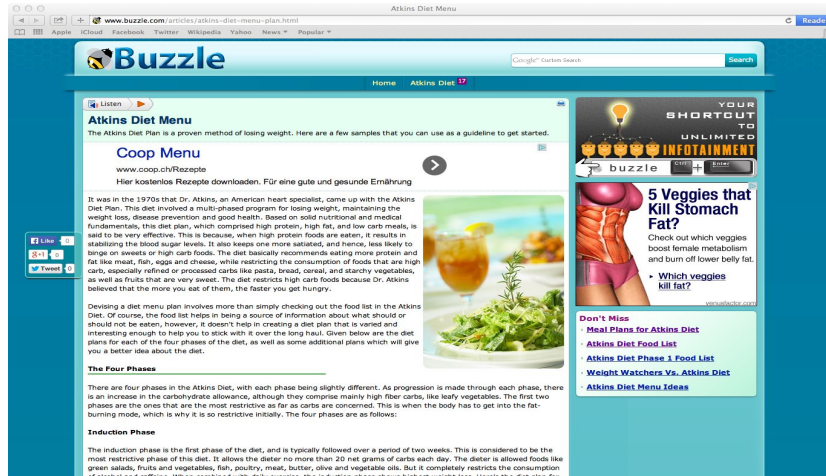


Potential labelers

- You
- Older days:
 - Undergraduate students
 - Domain experts (\$\$\$)
- Now: crowdsourcing
 - Can get both amateurs (~ undergrad students)
 - and experts

Labels

Collecting a lot of data (features) is easy
Labeling data is time consuming, difficult and
sometimes even impossible



Label: “Is page credible?”
Dietary expert is needed



*Is this webpage
credible?*

1. Submit task



Requester

4. Collect answers



CC-CC-C



2. Accept task



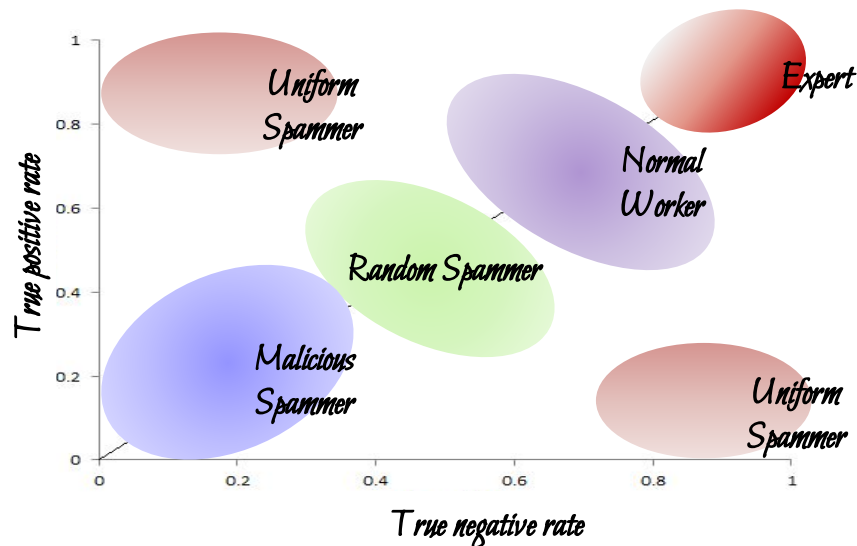
Crowd (workers)

3. Return answers

Crowdsourcing

Different types of workers

- Truthful
 - Expert
 - Normal
- Untruthful
 - Uniform spammer
 - Random spammer
 - Malicious spammer (a.k.a. a\$\$#*!e)



Catching malicious spammers

- Insert obvious examples for which you already know the labels (“honeypot”)
 - Tell workers they won’t be paid if they don’t get those right
 - Filter out workers who don’t get them right
- Aggregate multiple answers
 - p.t.o.

Crowdsourcing

Answer aggregation problem

- Have each example labeled by several workers, aggregate:
 - e.g., majority vote (works if only a minority is malicious)
 - e.g., “peer prediction”



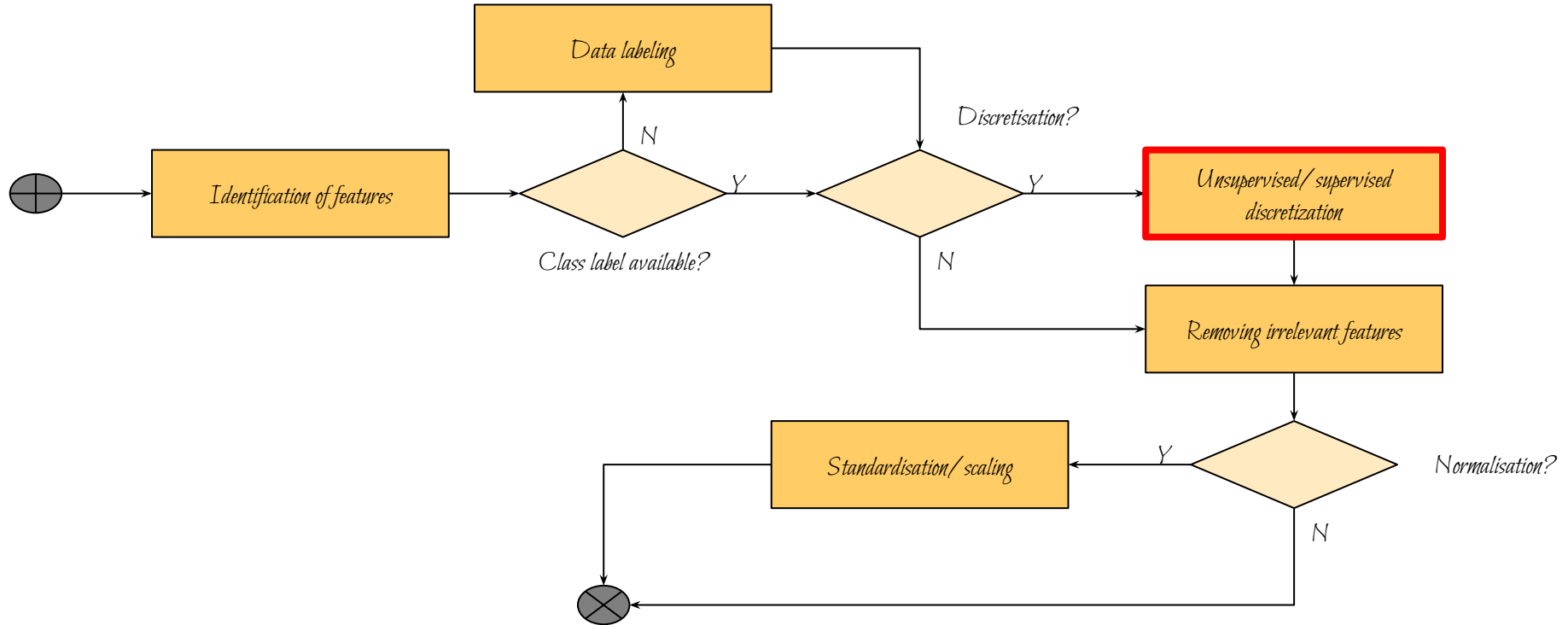
Crowd (workers)

Worker	Webpage	Credible
W_1	www.diet.com	C
W_2	www.diet.com	$\neg C$
W_3	www.diet.com	C
...



www.diet.com	C
--------------	---

Data collection



Discretization

Why?

- Some classifiers want discrete features (e.g., simplest kinds of decision trees)
- Discrete features let a linear classifier learn non-linear decision boundaries
- Certain feature selection methods require discrete (or even binary) features

Discretization

Unsupervised

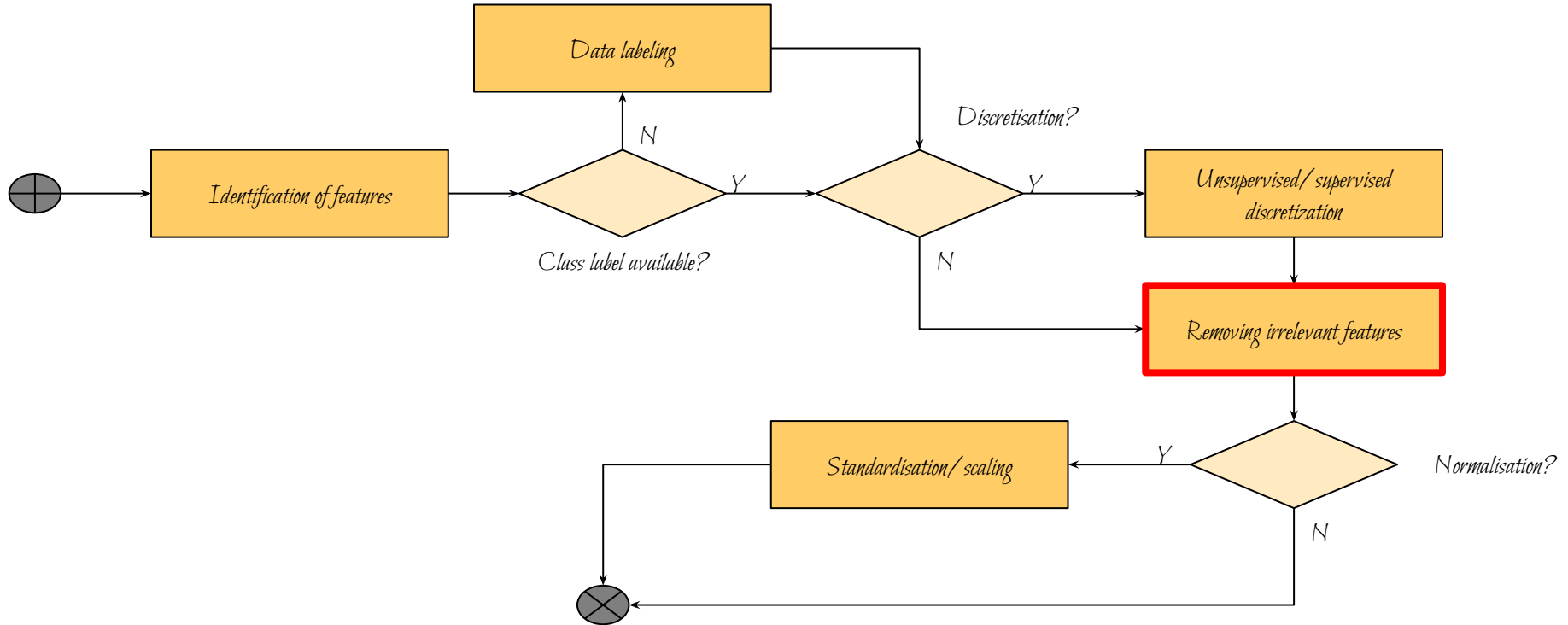
- Equal width
 - Divide the range into a predefined number of bins (bad for skewed data, e.g., from a power law)
- Equal frequency
 - Divide the range into a predefined number of bins so that every interval contains the same number of values
- Clustering

Discretization

Supervised

- Test the hypothesis that membership in two adjacent intervals of a feature is independent of the class
- If they are independent, they should be merged
- Otherwise they should remain separate
- Independence test: χ^2 statistics

Data collection



Feature selection

Reducing the number of N features to a subset of the best size $M < N$

There are 2^N possible subsets

Solutions

- Filtering as preprocessing (“offline”)
- Iterative feature selection (“online”)

Offline feature selection

Rank features according to their individual predictive power; select the best ones

- Pros
 - Independent of the classifier (performed only once)
- Cons
 - Independent of the classifier (ignore interaction with the classifier)
 - Assume features are independent

Ranking of features

Continuous features (and ideally labels)

- Correlation coefficient (Pearson, linear rel.!)

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Categorical features and labels

- Mutual information

$$I(F; C) = H(C) - H(C | F) = H(F) + H(C) - H(F, C)$$

$$H(F) = - \sum_i P(f_i) \log_2 P(f_i)$$

$$H(F, C) = - \sum_i \sum_j P(f_i, c_j) \log_2 P(f_i, c_j)$$

Ranking of features

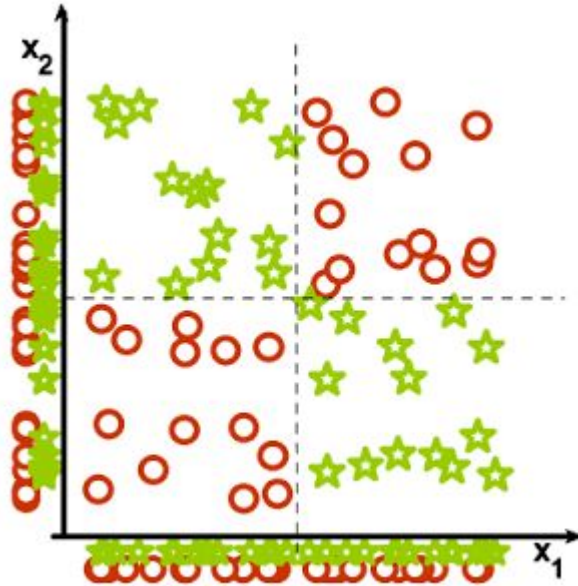
Categorical features and labels (cont'd)

- **χ^2 method**

Difference w.r.t. mutual information: the chi-square test checks the independence of the class and the feature, without indicating the strength or direction of any existing relationship (you just get a significance, a.k.a. p-value)

Ranking of features

Beware: collectively relevant features may look individually irrelevant!



Online feature selection

Forward feature selection: iteratively add features, using cross-validation to guide feature inclusion; stop when no improvement

- **Pros**

- Interact with the classifier
- No feature-independence assumption

- **Cons**

- Computationally intensive

Online feature selection

Backward selection (a.k.a. ablation): iteratively **remove** features, using cross-validation to guide feature removal; stop when no improvement

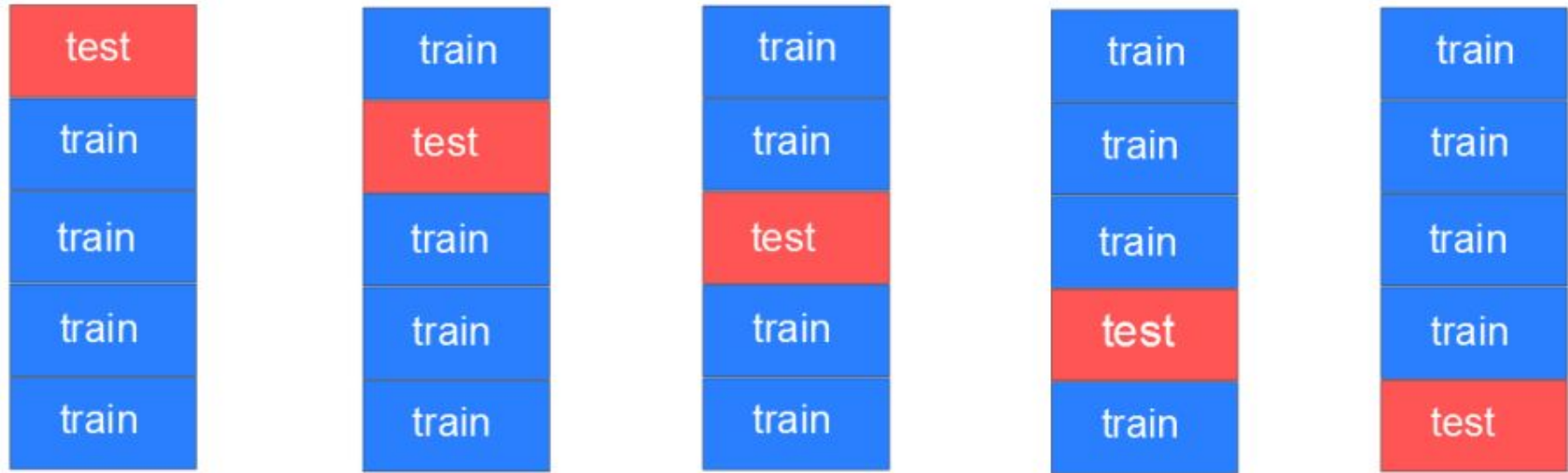
- **Pros**

- Interact with the classifier
- No feature-independence assumption

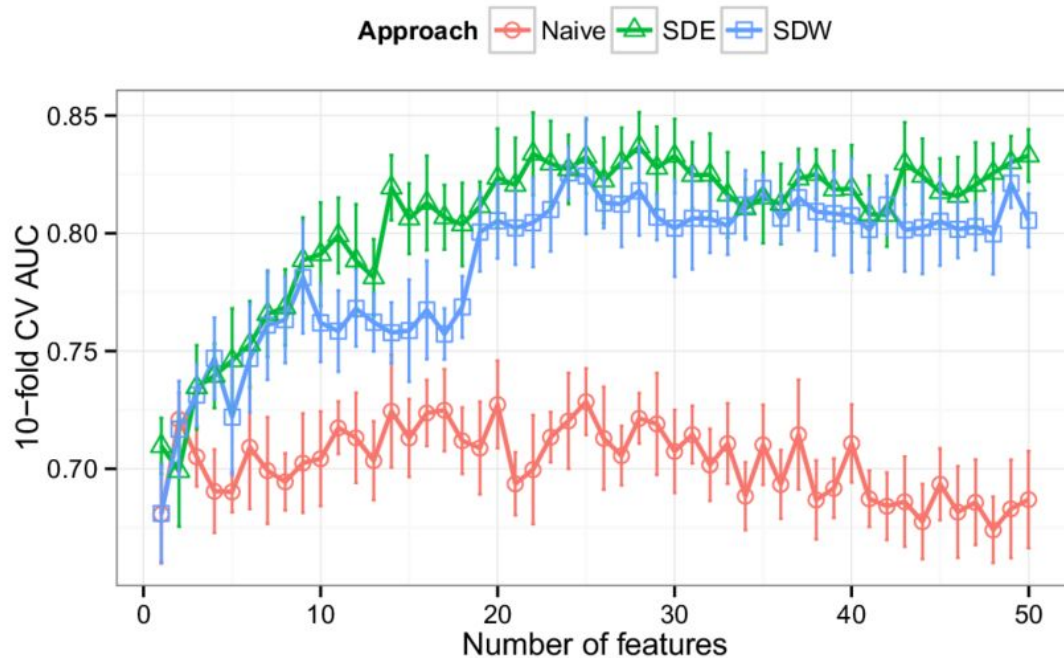
- **Cons**

- Computationally intensive

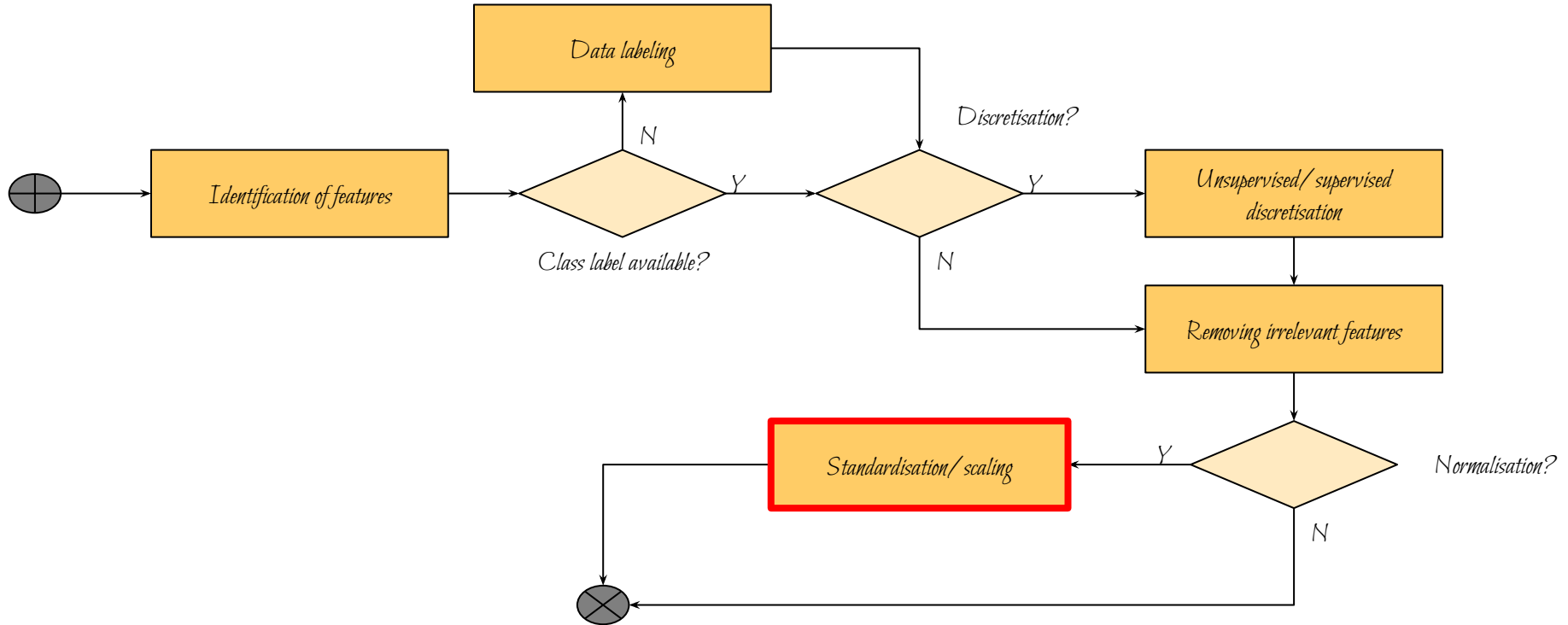
Cross-validation



Feature count vs. performance



Data collection



Feature normalization

Some classifiers do not manage well features with very different scales

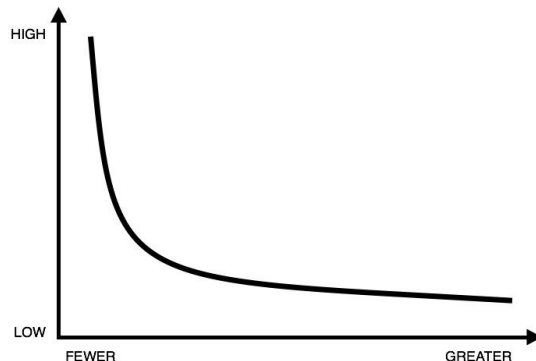
- Revenue: 10,000,000 (CHF)
- # of employees: 300

Features with large values dominate the others, and the classifier tends to over-optimize them

Logarithmic scaling

$$x_i' = \log(x_i)$$

- Consider order of magnitude, rather than direct value
- Good for heavy-tailed features (e.g., from power laws)



Min-max scaling

$$x_i' = (x_i - m_i) / (M_i - m_i)$$

where M_i and m_i are the max and min values of feature x_i respectively

The new feature x_i' lies in the interval $[0,1]$

Standardization

$$x_i' = (x_i - \mu_i) / \sigma_i$$

where μ_i is the mean value of feature x_i , and σ_i is the standard deviation

The new feature x_i' has mean 0 and standard deviation 1

Dangers of standardization and scaling

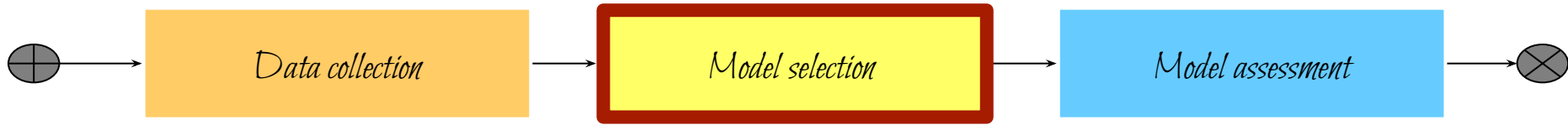
Standardization:

- Assume that the data has been generated by a Gaussian
- Uses mean and std → not meaningful for heavy-tailed data (may be mitigated by log-scaling)

Min-max scaling:

- If the data has outliers, they scale the typical values to a very small interval

Classification pipeline



Model selection: high level

Need to choose type of model

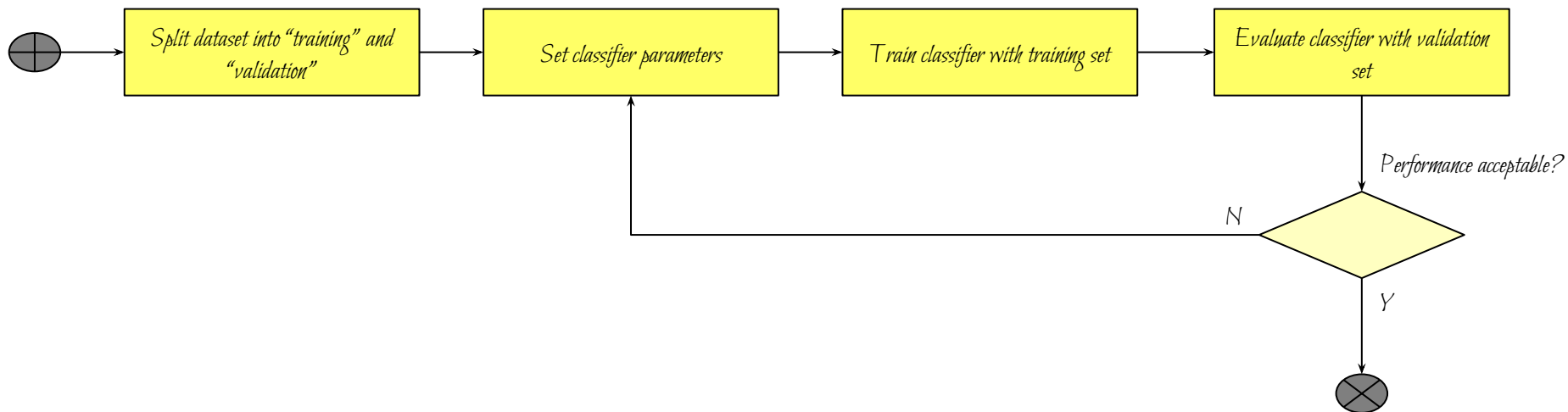
- Logistic regression?
- Decision trees?
- Random forest?
- Gradient-boosted trees?
- Support vector machine?
- Deep learning?
- ...

Model selection: low level

Usually a classifier has some “hyperparameters” to be tuned

- Threshold (e.g., logistic regression)
- Regularization factor
- Distance function (e.g., k-NN)
- Number of neighbors (e.g., k-NN)
- Number of trees (e.g., random forest)

Model selection



Need evaluation metric!

Loss function (more of them later!)

Categorical output

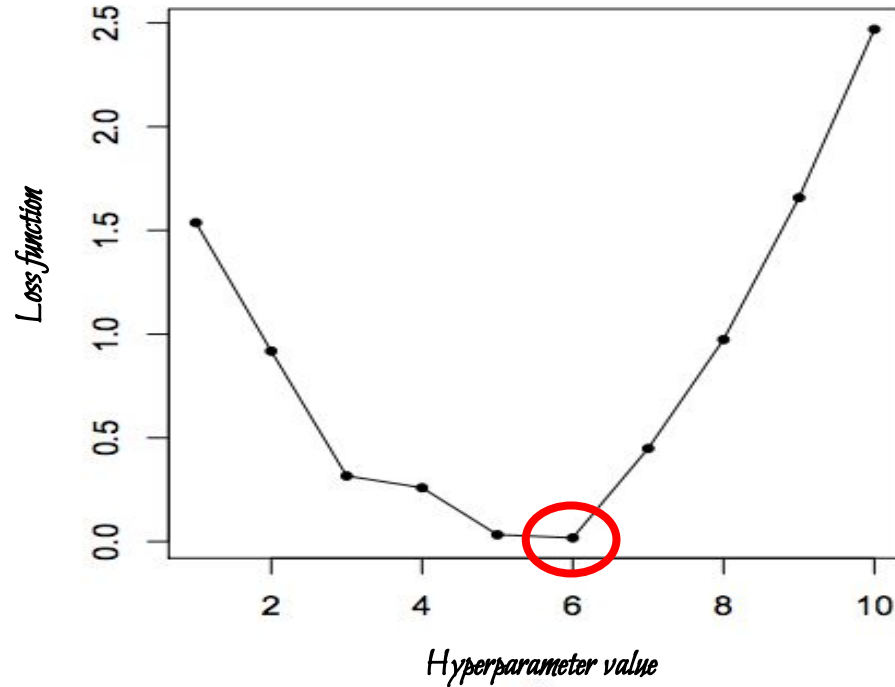
- e.g., 0-1 loss function (= 1 minus accuracy):

$$J = \frac{1}{n} \sum_{i=1}^n \#(y \neq f(x_i))$$

Real-valued output

- e.g., squared error: $J = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$
- e.g., absolute error: $J = \frac{1}{n} \sum_{i=1}^n |y_i - f(x_i)|$

Model selection



Performance metrics for binary classification

For categorical binary classification, the usual metrics are based on the **confusion matrix**, which has 4 values:

- True Positives (positive examples classified as positive)
- True Negatives (negative examples classified as negative)
- False Positives (negative examples classified as positive)
- False Negatives (positive examples classified as negative)

		Class	
		A	B
Classified	A	TP	FP
	B	FN	TN

Accuracy

$$A = \frac{TP+TN}{TP+TN+FP+FN} = \frac{TP+TN}{N}$$

Appropriate metric when

- classes are not skewed
- Errors (FP, FN) have the same importance

Accuracy (skewed example)

Classifier 1		Class	
		Fraud	¬Fraud
Classified	Fraud	5	10
	¬Fraud	5	80

$$A = 85/100 = 0.85$$

Always ¬Fraud		Class	
		Fraud	¬Fraud
Classified	Fraud	0	0
	¬Fraud	10	90

$$A = 90/100 = 0.90$$

Question time

		Class	
		A	B
Classified	A	45	20
	B	5	30

		Class	
		A	B
Classified	A	40	10
	B	10	40

Which is the “best” classifier?

- A. Classifier 1
- B. Classifier 2
- C. Both are equally good

Question time

		Class	
		Cancer	¬Cancer
Classified	Cancer	45	20
	¬Cancer	5	30

		Class	
		Cancer	¬Cancer
Classified	Cancer	40	10
	¬Cancer	10	40

Which is the “best” classifier?

- A. Classifier 1
- B. Classifier 2
- C. Both are equally good

Precision and recall

Precision: what fraction of positive predictions are actually positive?

$$P = \frac{TP}{TP + FP}$$

Recall: what fraction of positive examples did I recognize as such?

$$R = \frac{TP}{TP + FN}$$

Precision and recall

Classifier 1		Class	
		Cancer	¬Cancer
Classified	Cancer	45	20
	¬Cancer	5	30

$$P_1 = 45/65 = 0.69$$

$$R_1 = 45/50 = 0.9$$

Classifier 2		Class	
		Cancer	¬Cancer
Classified	Cancer	40	10
	¬Cancer	10	40

$$P_2 = 40/50 = 0.8$$

$$R_2 = 40/50 = 0.8$$

Everybody has cancer		Class	
		Cancer	¬Cancer
Classified	Cancer	50	50
	¬Cancer	0	0

$$P = 50/100 = 0.5$$

$$R = 50/50 = 1$$

F-score

$$P = \frac{TP}{TP + FP}$$

Sometimes it's necessary to have a unique metric to compare classifiers

$$R = \frac{TP}{TP + FN}$$

F-score (or F1-score): harmonic mean of precision and recall

$$F1 = 1 / (0.5 * (1/P + 1/R)) = 2 \cdot \frac{P \cdot R}{P + R}$$

Precision and recall can be **differently weighted**, if one is more important than the other

Precision and recall

Classifier 1		Class	
		Cancer	¬Cancer
Classified	Cancer	45	20
	¬Cancer	5	30

$$F_1 = 2 * (0.69 * 0.9) / (0.69 + 0.9) \\ = 0.78$$

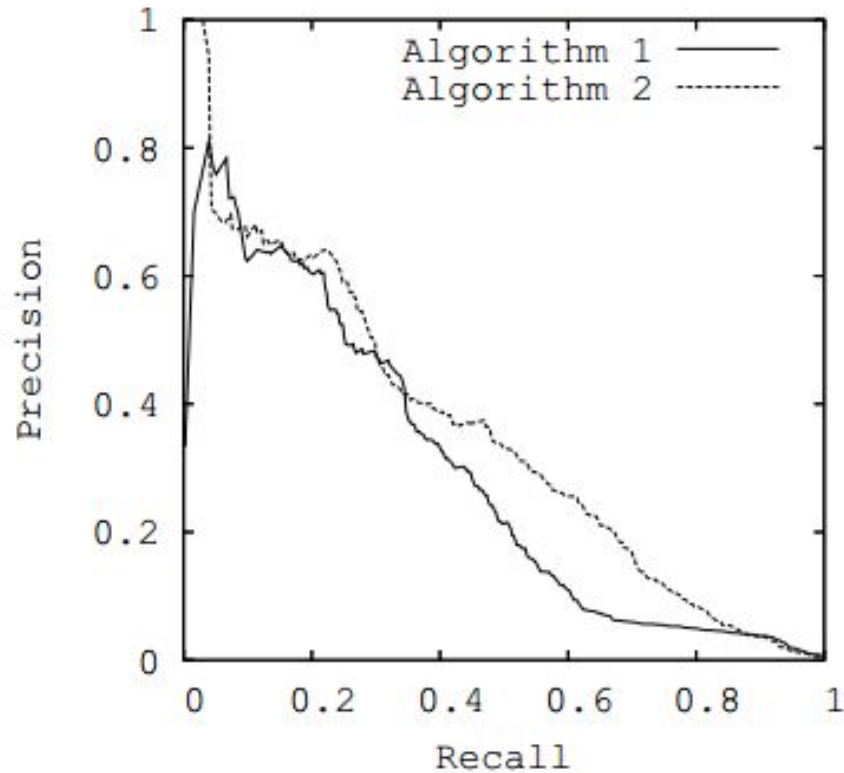
Classifier 2		Class	
		Cancer	¬Cancer
Classified	Cancer	40	10
	¬Cancer	10	40

$$F_2 = 2 * (0.8 * 0.8) / (0.8 + 0.8) \\ = 0.8$$

Everybody has cancer		Class	
		Cancer	¬Cancer
Classified	Cancer	50	50
	¬Cancer	0	0

$$F = 2 * (0.5 * 1) / (0.5 + 1) = 0.66$$

Precision/recall curve



Recap



- Model selection:
 - Use training data in cross-validation
 - Need evaluation metric
 - Typically based on confusion matrix
 - e.g., accuracy, precision, recall



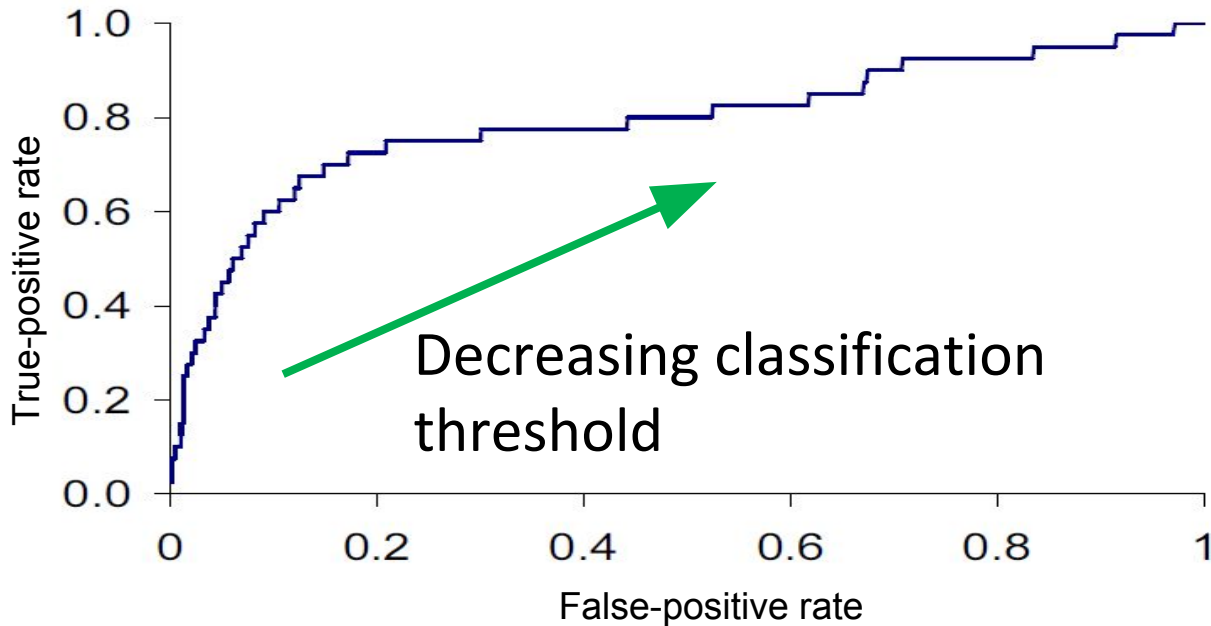
		Class	
		A	B
Classified	A	TP	FP
	B	FN	TN

ROC plots

ROC = Receiver-Operating Characteristic (WTF?!)

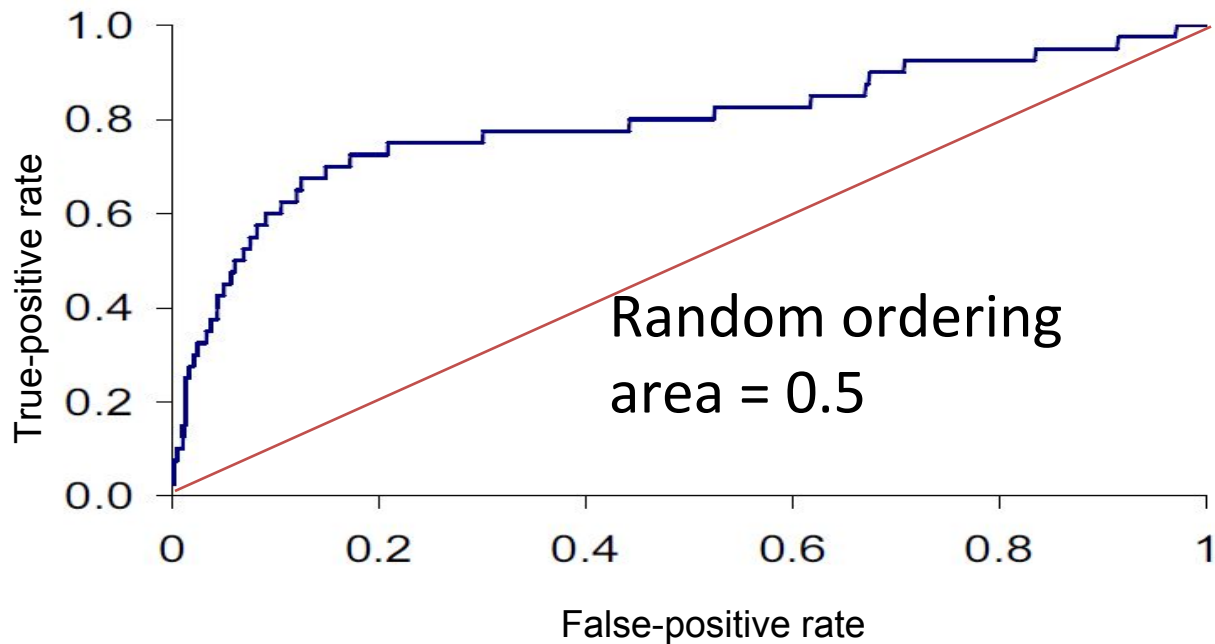
Y-axis: true-positive rate = $TP / (TP + FN)$, a.k.a. recall

X-axis: false-positive rate = $FP / (FP + TN)$

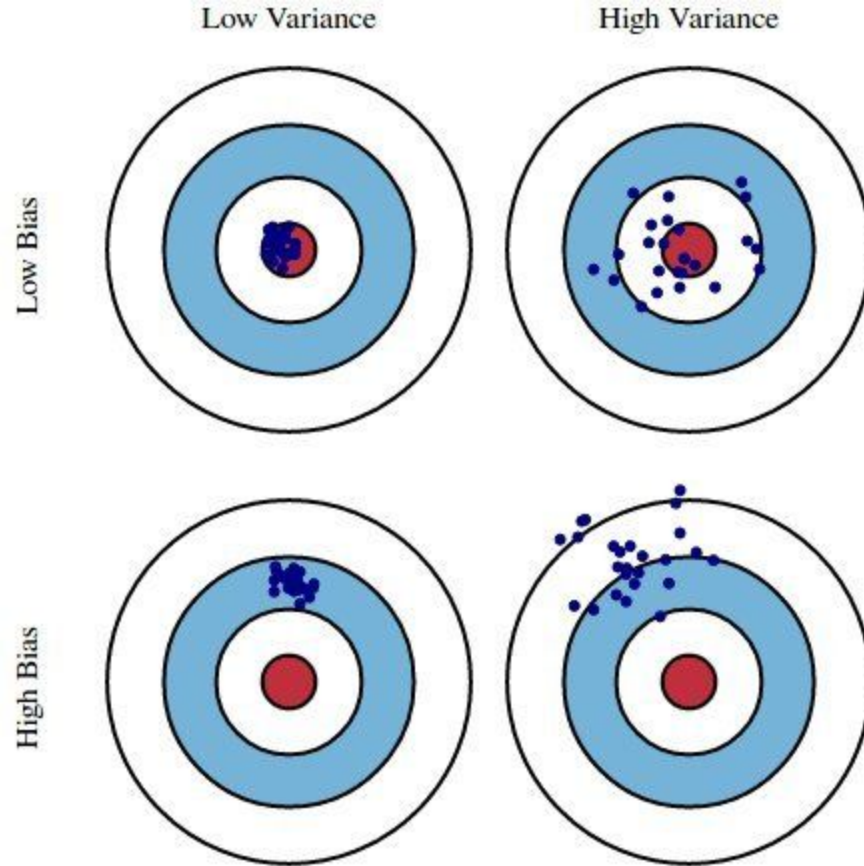


ROC AUC

ROC AUC is the “Area under the curve” – a single number that captures the overall quality of the classifier. It should be between 0.5 (random classifier) and 1.0 (perfect).

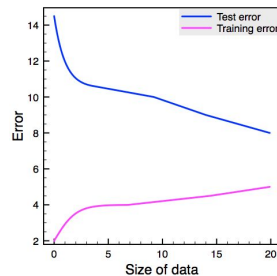


Bias and variance



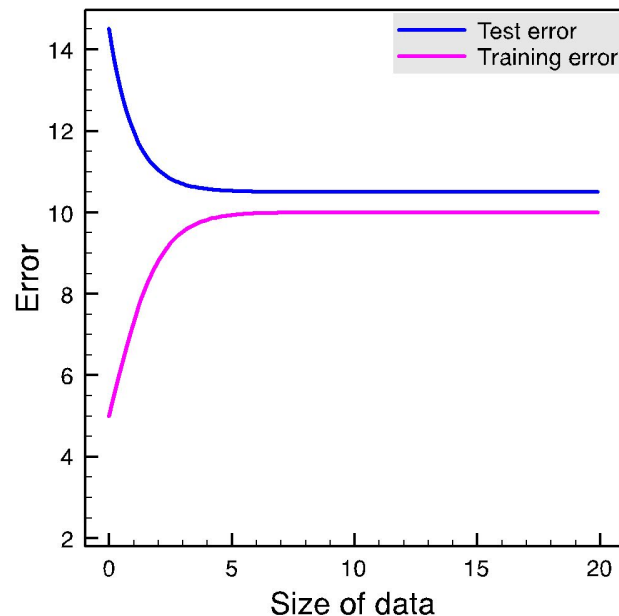
Bias and variance

Bias and variance can be assessed by comparing the error metric on the training set and the test set => always **plot learning curves** (training set size vs. error)

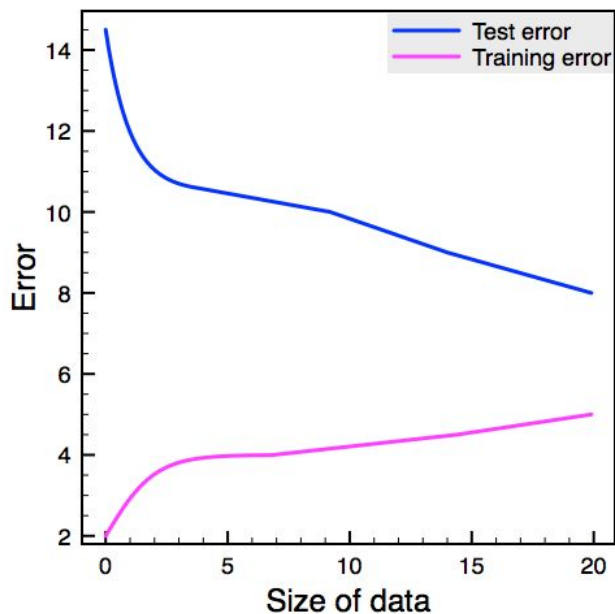


In the ideal case, we want low bias (small training error) and low variance (small test error)

When more data helps



High bias



High variance

More data often beats better algorithms



EXPERT OPINION

Contact Editor: **Brian Brannon**, bbrannon@computer.org

The Unreasonable Effectiveness of Data

Alon Halevy, Peter Norvig, and Fernando Pereira, *Google*

Large Language Models in Machine Translation

Thorsten Brants Ashok C. Popat Peng Xu Franz J. Och Jeffrey Dean

Google, Inc.

1600 Amphitheatre Parkway

Mountain View, CA 94303, USA

`{brants,popat,xp,och,jeff}@google.com`

Abstract

This paper reports on the benefits of large-scale statistical language modeling in machine translation. A distributed infrastructure is proposed which we use to train on up to 2 trillion tokens, resulting in language models having up to 300 billion n -grams. It is capable of providing smoothed probabilities for fast, single-pass decoding. We introduce a new smoothing method, dubbed *Stupid Backoff*, that is inexpensive to train on large data sets and approaches the quality of Kneser-Ney Smoothing as the amount of training data increases.

How might one build a language model that allows scaling to very large amounts of training data? (2) How much does translation performance improve as the size of the language model increases? (3) Is there a point of diminishing returns in performance as a function of language model size?

This paper proposes one possible answer to the first question, explores the second by providing learning curves in the context of a particular statistical machine translation system, and hints that the third may yet be some time in answering. In particular, it proposes a *distributed* language model training and deployment infrastructure, which allows direct and efficient integration into the hypothesis-search algorithm rather than a follow-on re-scoring phase.

Classification pipeline



Model assessment

Model assessment is the goal of estimating the performance of a fixed model (i.e., the best model found during model selection)

Use held-out test set that you've never seen during training

Important reads

Machine Learning that Matters

Kiri L. Wagstaff

KIRI.L.WAGSTAFF@JPL.NASA.GOV

Jet Propulsion Laboratory, California Institute of Technology, 4800 Oak Grove Drive, Pasadena, CA 91109 USA

A Few Useful Things to Know about Machine Learning

Pedro Domingos
Department of Computer Science and Engineering
University of Washington
Seattle, WA 98195-2350, U.S.A.
pedrod@cs.washington.edu



[slides](#)

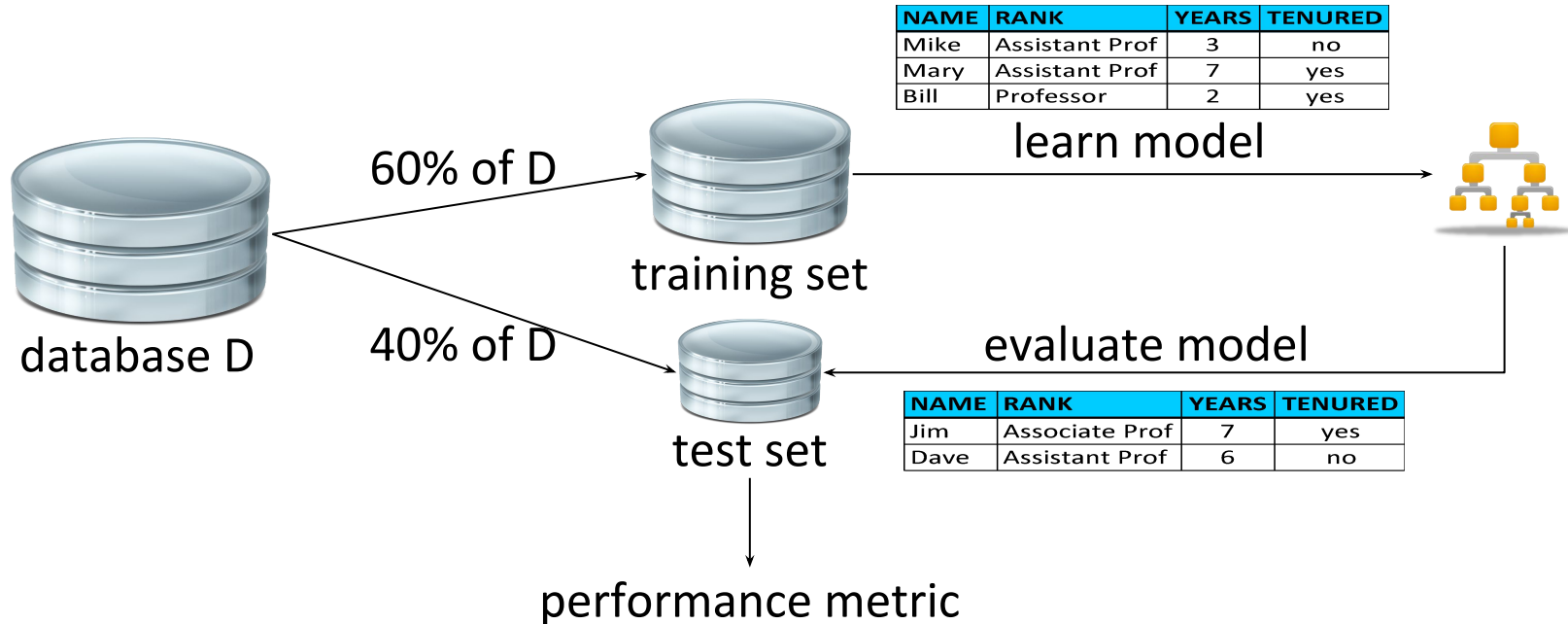
Feedback

Give us feedback on this lecture here:

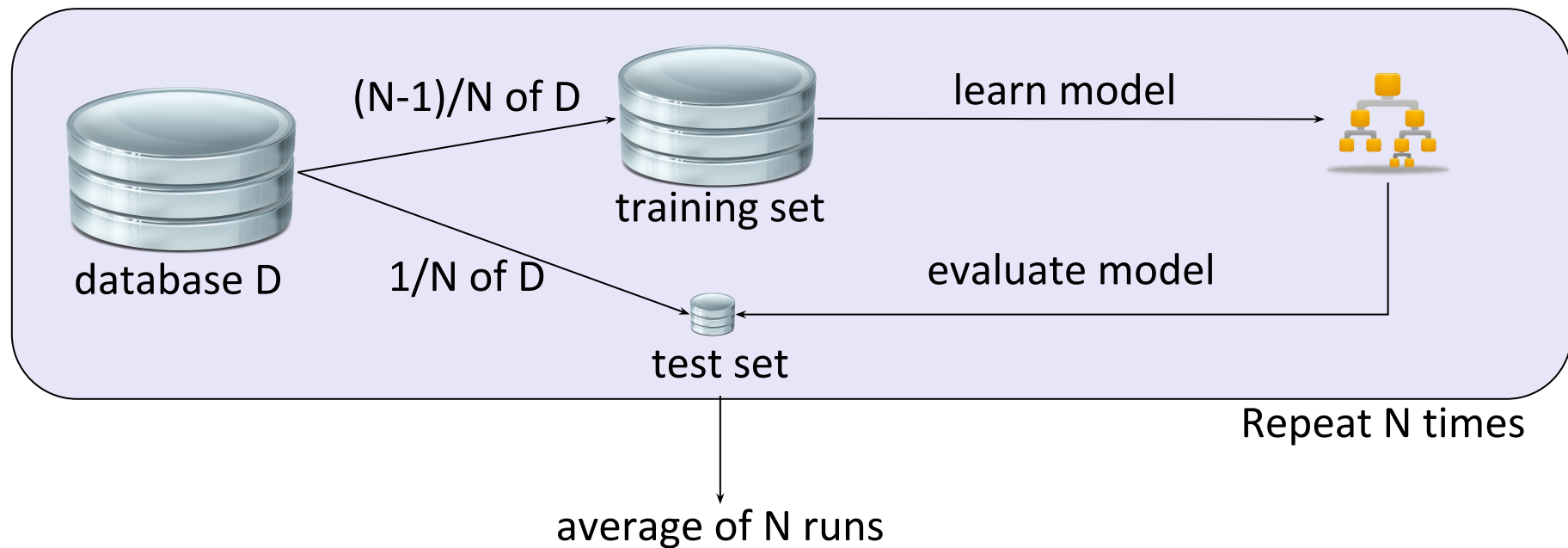
<https://go.epfl.ch/ada2018-lec9-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- What's your favorite messaging system?
- ...

Training and testing with heaps of data



Data-efficient training and testing: Leave-one-out cross-validation



Data-efficient training and testing: k-fold cross validation

