Supervised Learning





Mirco Schönfeld University of Bayreuth

mirco.schoenfeld@uni-bayreuth.de @TWIyY29









3

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0











Supervised Machine Learning

- Objective:
- learn the *class label* y for each value of x, the *feature vector* consisting of multiple *features* (categorical or numerical)
- Result:
- A function f(x) = y that best predicts y for each value of x

If y is

- a real number, a regression model is learned
- a Boolean value (true/false, +1/-1), we speak of binary classification
- a nominal value of some finite set, it is a *multiclass classification*





Supervised Machine Learning: Bird View



https://exeter-data-analytics.github.io/MachineLearning/supervised-learning-1.html



- Deal with missing data (ignore/impute)
- Use unsupervised methods to accentuate the data's structure
- Remove uninformative variables (crude filtering)
- Hand-crafted features
- Dimensionality reduction techniques
- Features are typically normalized/standardized









Discrete vs. Continuous



7



Dichotomous data:

Data points can only take up 2 values

Discrete data:

Data points can only take up values from a set of possible values.

Either finite or infinite but countable

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0





Continuous data:

Data points can be measured to an arbitrary level of exactness



Scales

8

Nominal

Values divide space of possible values into different segments

No order between segments

Example: gender, nationality

Ordinal

Allows for rank order – data can be sorted

No relative degree or relative difference between groups

Example: school grades

S. S. Stevens: On the Theory of Scales of Measurement. In: Science. 1946, 103, S. 677–680.



Interval

Order and difference between groups is defined

No natural "zero" & ratio between points is undefined

Example: temperature in Celsius

Ratio

Meaningful zero, i.e. unique and non-arbitrary

Example: temperature in Kelvin, age, weight, height

















Regression

11



Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0





Classification







Supervised Machine Learning...in other words

Supervised Machine Learning aims at forecasting class labels for measured data

Correct class labels are known for training data

Training means to search for a good function mapping measured artifacts to known class labels

What you need:

- Classificator
- Measured data \bullet
- Class labels





Weight

Chihuahuas

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0









15 An Example (<2 lbs, 5 in>, Chihuahua) Height Weight: 2 punds Height: 5 inches

Weight

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0







This is a training-set pair consisting of a *feature vector* and a *class label*.







Weight

An Example

Weight

The classifier *f*

Height

17

if (height > 7) print Beagle else if (weight < 3) print Chihuahua else print Dachshund;

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0



Relation to Clustering

Classification:

- Class label is discrete
- Enough training data for each class Regression:
 - Target variable is numeric

Clustering:

- no class label / target variable
- no training input data without pre-defined classes
- produces partitioning of data

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0



Clusters could be used to identify new classes





Major Classes of Classifiers

- Decision trees Suitable for binary and multiclass classification with a limited number of features
- Perceptrons Applies weights to components of vectors. Output +1 if sum exceeds a threshold, otherwise -1.
- Neural networks Acyclic networks of perceptrons
- Instance-based learning Compares instances of data to the entire training set. Example: k-nearest neighbor.
- Support-vector machines Maps training examples to points in space so as to maximise the gap between categories











Perceptrons

21



Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

Components of instances

Neural Networks

)
$$\xrightarrow{\text{Output}} 1$$

Neural Networks

23

https://towardsdatascience.com/designing-your-neural-networks-a5e4617027ed

k-Nearest Neighbors

24

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

Support Vector Machines

25

https://en.wikipedia.org/wiki/Support-vector_machine

Some Classification Algorithms Compared

Supervised Learning: Training and Testing

27

Leskovec, J., Rajaraman, A. and Ullman, J.D., 2020. *Mining of massive data sets*. Cambridge university press.

Model is represented as classification rules, decision trees, or mathemical formulas

The model is used to classify the data whose class is *unknown*! Its purpose is *generalization*!

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

Model

How to Design a Model: Feature Selection

Formulate characteristics that help distinguishing between classes.

For spam-detection: find words or combinations of words that indicate a mail being spam.

Classification using a model

Spam: Wholesale Fashion Watches -57% today. Designer watches for cheap ... Spam: You can buy Viagra Fr\$1.85 All Medications at unbeatable prices! ... Spam: WE CAN TREAT ANYTHING YOU SUFFER FROM JUST TRUST US ...

Ham: The practical significance of hypertree width in identifying more ... Ham: Abstract: We will motivate the problem of social identity clustering: ... Ham: Good to see you my friend. Hey Peter, It was good to hear from you. ...

Curse of Dimensionality:

Including more features will improve classification conceptually but will render computation increasingly difficult.

Spam vs. Ham

- Spam: Sta.rt earn*ing the salary yo,u d-eserve by o'btaining the prope,r crede'ntials!
- Ham: PDS implies convexity of the resulting optimization problem (Kernel Ridge ...

/kæp.tʃə/ or Creating a Training Set

Label Data Manually?

Often time consuming and costly process for large training sets

Select all images with a store front. Click verify once there are none left.

Active learning:

kick start a classifier only with some training examples, but leave it primarily with unclassified data, which it must classify. If the classifier is unsure of the classification (e.g., the newly arrived example is very close to the boundary), then the classifier can ask for ground truth at some significant cost

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

How to Choose a Model

31

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

Classification using a model

Spam vs. Ham

Model A:		Correct cla	ass		Modal F	
		Hn	Spam			D .
am-Filter	Harr	200	38			
orts	S _F	0	762		Correct cla	ISS
					Ham	Sp
	Sp	pam-Filter		Ham	189	
	re	ports	5	Spam	11	7

Beware of Overfitting

Two types of classification errors:

- 1. Training error misclassification on training data
- 2. Generalization error expected error on *unseen* data

Overfitting:

32

Good results on training data (low training error) and bad results with test/validation data (high generalization error)

Error significantly *underestimated* – severe problem in application scenarios

Detecting overfitting: Evaluation of training with *new* data – NOT using training data!

Training Test Split

Split training data at least in training and testing (Popular splits: 2:1 / 90:10)

Recommended: Split data in training, testing and validation Splits: 80:10:10

Choose best classificator *only* on training and test data Estimate accuracy & tune parameters of model

Keep validation-data secret! Use that only once to estimate the generalization power of the model!

Terminology of test and validation data is often mixed up.

What if...

- Prepare data
- **Chose classificator**
- Train it

34

- Test it
- Validate it
- Results do not look good
- Repeat

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

What's the problem?

Repeated testing leads to overfitting

Once the validation data is used, do not go back to improve classification!

Cross-Validation

Cross-validation is a technique to help choosing classificator and optimal parameters

Partition data in k non-overlapping parts of equal size During *i*th iteration, use data in partition D_i for validation, all other data as training data

Quality of classificator: mean over all k iterations

Iteration 1
Iteration 2
Iteration 3
Iteration k

Exhaustive Cross-Validation

Cross-validation methods which learn and test on all possible combinations to divide the original sample into training and test set.

Leave-p-out cross-validation:

Use p observations as the test set and the remaining observations as training set.

Leave-one-out cross-validation:

Leave-p-out cross validation with p=1 Means finding one classificator for each instance – N classificators!

Imbalanced data

Imbalance:

Number of samples of different classes are diverging significantly.

Consequences of building models using imbalanced data:

- Bias Classifiers are more sensitive to detecting the majority class
- Optimization metrics Metrics like accuracy may not report true performance

Has implications for sampling for cross validation!

- Often, collecting samples of a certain class is difficult because these are rare events.

Resampling

39

adding more examples from the minority class (over-sampling)

Various strategies, e.g. under-sampling by generating cluster-cendroids, over-sampling by synthesizing elements (SMOTE), ...

Balancing classes by removing samples from the majority class (under-sampling) and/or

Oversampling

Original dataset

Overfitting can occur on subtle ways

- Evaluation and task are adapted to solution
- Preprocessing of data tells something about solution
- Unbalanced data

40

- Little variety of data
- New observations
- Insufficient data

ution

No Free Lunch Theorem

Choosing an appropriate algorithm requires making assumptions

Wolpert, David H., and William G. Macready. "No free lunch theorems for optimization." IEEE transactions on evolutionary computation 1.1 (1997): 67-82.

With no assumptions, there will be no universal algorithm "better" than random choice

[...] what an algorithm gains in performance on one class of problems is necessarily offset by its performance on the remaining problems;

Good Classifications: The Confusion Matrix

Model A:		Correct class	
		Ham	Spam
Spam-Filter	Ham	200	38
reports	Spam	0	762

R/	Correct class		
		Ham	Spam
Spam-Filter	Ham	189	1
reports	Spam	11	799

Classification

using a model

Spam vs. Ham

		Correct	
		Positive (P)	Negative (
cted	Positive	True Positive	False Positive
Pred	Negative	False Negative	True Negative

Measuring Goodness

43

Precision: Proportion of predicted positives that are truly positive \bullet good choice when we need to be very sure of prediction

Recall: Proportion of actual positives that are correctly classified \bullet good choice when as many positives as possible should be captured

TP

- TP
- TP + FP

ΓΡ		Correct	
+ FN		Positive (P)	Negative (
cted	Positive	True Positive	False Positive
Pred	Negative	False Negative	True Negative

Measuring Goodness

- Accuracy: Proportion of true results among total number of cases good choice when classes are balanced
 - TP + TN
 - TP + FP + FN + FN
- F_1 Score: harmonic mean between precision & recall a number between 0 and 1 good choice when we want a model with both good precision and recall
 - 2 * precisi precisi
 - Important variant F_{β} allows to apply a custom weight to precision & recall

ion * recall		Correct	
ion + re	call	Positive (P)	Negative (
icted	Positive	True Positive	False Positive
Pred	Negative	False Negative	True Negative

45

Measuring Goodness & more

Prevalence P	е	
$\overline{\mathrm{P}+\mathrm{N}}$		
accuracy	(ACC)	
ACC =	$\frac{\mathrm{TP} + \mathrm{TN}}{\mathrm{P} + \mathrm{N}} =$	$\frac{\mathrm{TP} + \mathrm{TN}}{\mathrm{TD} + \mathrm{TN} + \mathrm{FD} + \mathrm{FN}}$
balanced	P + N accuracy (B/	$\mathbf{1P} + \mathbf{1N} + \mathbf{FP} + \mathbf{FN}$
$BA = -\frac{7}{2}$	$\frac{TPR + TNR}{2}$	

F1 score

is the harmonic mean of precision and sensitivity: $\mathrm{F}_{1} = 2 \times \frac{\mathrm{PPV} \times \mathrm{TPR}}{\mathrm{PPV} + \mathrm{TPR}} = \frac{2\mathrm{TP}}{2\mathrm{TP} + \mathrm{FP} + \mathrm{FN}}$

phi coefficient (ϕ or r_{ϕ}) or Matthews correlation coefficient (MCC)

$$\mathrm{MCC} = rac{\mathrm{TP} imes \mathrm{TN} - \mathrm{FP} imes \mathrm{FN}}{\sqrt{(\mathrm{TP} + \mathrm{FP})(\mathrm{TP} + \mathrm{FN})(\mathrm{TN} + \mathrm{FP})(\mathrm{TN} + \mathrm{FN})}}$$

Fowlkes-Mallows index (FM)

$$\mathrm{FM} = \sqrt{\frac{TP}{TP + FP}} \times \frac{TP}{TP + FN} = \sqrt{PPV \times TPR}$$

informedness or bookmaker informedness (BM)

BM = TPR + TNR - 1

markedness (MK) or deltaP (Δp)

MK = PPV + NPV - 1

Diagnostic odds ratio (DOR)

$$\mathrm{DOR} = rac{\mathrm{LR}+}{\mathrm{LR}-}$$

sensitivity, recall, i

$$TPR = \frac{TP}{P} = \frac{TP}{TP}$$
specificity, selective

$$TNR = \frac{TN}{N} = \frac{TP}{TP}$$
precision or positive

$$PPV = \frac{TP}{TP + FP}$$
negative predictive

$$NPV = \frac{TN}{TN + FN}$$
miss rate or false r

$$FNR = \frac{FN}{P} = \frac{FN}{FN}$$
fall-out or false positive

II, hit rate, or true positive rate (TPR) $rac{\mathrm{TP}}{\mathrm{P}+\mathrm{FN}} = 1-\mathrm{FNR}$ vity or true negative rate (TNR) $rac{\mathrm{TN}}{\mathrm{TN}+\mathrm{FP}} = 1-\mathrm{FPR}$ ve predictive value (PPV) = 1 - FDRe value (NPV) = 1 - FORnegative rate (FNR) $\frac{\mathrm{FN}}{\mathrm{N} + \mathrm{TP}} = 1 - \mathrm{TPR}$ sitive rate (FPR) $\mathrm{FPR} = rac{\mathrm{FP}}{\mathrm{N}} = rac{\mathrm{FP}}{\mathrm{FP} + \mathrm{TN}} = 1 - \mathrm{TNR}$

false discovery rate (FDR) $\mathrm{FDR} = rac{\mathrm{FP}}{\mathrm{FP} + \mathrm{TP}} = 1 - \mathrm{PPV}$ false omission rate (FOR) $\mathrm{FOR} = rac{\mathrm{FN}}{\mathrm{FN} + \mathrm{TN}} = 1 - \mathrm{NPV}$ Positive likelihood ratio (LR+) $\mathrm{LR}+=rac{\mathrm{TPR}}{\mathrm{FPR}}$ Negative likelihood ratio (LR-) $LR-=rac{FNR}{TNR}$ prevalence threshold (PT) $\mathrm{PT} = \frac{\sqrt{\mathrm{TPR}(-\mathrm{TNR}+1)} + \mathrm{TNR} - 1}{(\mathrm{TPR} + \mathrm{TNR} - 1)} = \frac{\sqrt{\mathrm{FPR}}}{\sqrt{\mathrm{TPR}} + \sqrt{\mathrm{FPR}}}$ threat score (TS) or critical success index (CSI) $\mathrm{TS} = \frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN} + \mathrm{FP}}$

Correct

		Positive (P)	Negative (
cted	Positive	True Positive	False Positive
Pred	Negative	False Negative	True Negative

Interpretability

Large project set out to evaluate ML application to problems in healthcare Scenario: predicting pneumonia risk Goal: predict probability of death for patients with pneumonia

The most accurate model of the study was a multitask neural net. Outperformed other models by wide margin but was still dropped. Why?

One rule-based system learned the rule "patient has asthma \rightarrow lower risk" Reflected a true pattern in training data

The best model was the least intelligible one – was deemed to risky No way of checking the features that were picked up

Cooper, Gregory F., et al. "An evaluation of machine-learning methods for predicting pneumonia mortality." Artificial intelligence in medicine 9.2 (1997): 107-138. Caruana, Rich, et al. "Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission." Proceedings of the 21th ACM SIGKDD. 2015.

Goals of Interpretable Models

- Trust: Identify and mitigate *bias* Recognizing bias in a black-box algorithm is *very* hard
- Causality: Account for context Helps you understand how the factors included in the model led to the prediction
- Informativeness: Extract knowledge Helps you determine if patterns that appear to be present in the model are really there. Rather learning from the model than evaluating it (compared to identifying bias)
- Transferability: Generalize Models are trained on carefully collected datasets to solve narrowly defined problems. Interpretable models should help you determine if and how they can be generalized
- Fair and Ethical Decision-Making algorithmic decision-making mediates more and more of our interactions. Need a way to make sure that decisions conform to ethical standards

Properties of Interpretable Models

Transparency 1.

- Simulatability Transparency at the level of the entire model
- Decomposability / Intelligibility Transparency at the level of the individual components, e.g. parameters
- Algorithmic Transparency Transparency at the level of the training algorithm

Properties of Interpretable Models

- 2. Post-hoc Interpretability
 - **Text Explanations**
 - Visualization

- Local Explanations
- Explanation by Example

Ensemble Models

would interact in unintended ways

Helps overcoming

- low accuracy: single models/algorithms might not be good enough
- high variance: single estimators are very sensitive to inputs to the learned features lacksquare• features noise and bias: single estimators might rely heavily on one or few features

Various ways to combine models

- Inspectable models ease debugging problems in data collection, feature engineering, etc
- Ensemble models provide ways to restrict features to separate models that otherwise

Ensemble Models

Using multiple algorithms to diversify model predictions

https://towardsdatascience.com/ensemble-models-5a62d4f4cb0c

Ensemble Models

Using multiple weak instances of the same algorithm

https://towardsdatascience.com/ensemble-models-5a62d4f4cb0c

Predictions

Ensemble Models: Bagging & Bootstrapping

Combine predictions from multiple models

https://towardsdatascience.com/ensemble-models-5a62d4f4cb0c

Prof. Dr. Mirco Schönfeld | Data Modeling & Knowledge Generation | v1.0

Ensemble Models: Boosting

Ensemble of algorithms that builds models on top of several weak learners Here: Sequental Adaptive Boosting (AdaBoost)

Ensemble Models: Stacking

Stacking intermediate predictions to make a final prediction

Algorithm 1

https://towardsdatascience.com/ensemble-models-5a62d4f4cb0c

How to select a model?

- Quality of predictions i.e. performance in terms of a quality metric
- Speed

56

i.e. training time, prediction time

Robustness

i.e. handling noise or missing values and still classify correctly

Scalability •

i.e. computational efficiency

- Interpretability lacksquaresubjective means
- Other

Thanks. mirco.schoenfeld@uni-bayreuth.de https://xkcd.com/1838/