Data Preparation & Building Numerical Predictors

KNIME AG





Data Preparation

Motivation

Real world data is "dirty"

→ Contains missing values, noises, outliers, inconsistencies

Comes from different information sources

 \rightarrow Different attribute names, values expressed differently, related tuples

Different value ranges and hierarchies

→ One attribute range may overpower another

Huge amount of data

3

→ Makes analyis difficult and time consuming



Data Preparation

- Data Cleaning & Standardization (domain dependent)
- Aggregations (often domain dependent)
- Normalization
- Dimensionality Reduction
- Outlier Detection
- Missing Value Imputation
- Feature Selection
- Feature Engineering
- Sampling
- Integration of multiple Data Sources

4



Data Preparation: Normalization

Normalization: Motivation

Example:

- Lengths in cm (100 200) and weights in kilogram (30 150) fall both in approximately the same scale
- What about lengths in m (1-2) and weights also in gram (30000 150000)?
 The weight values in mg dominate over the length values for the similarity of records!

Goal of normalization:

Transformation of attributes to make record ranges comparable



6

Normalization: Techniques

min-max normalization

$$y = \frac{x - x_{min}}{x_{max} - x_{min}} (y_{max} - y_{min}) + y_{min}$$



z-score normalization

$$y = \frac{x - mean(x)}{stddev(x)}$$

7



normalization by decimal scaling

$$y = \frac{x}{10^{j}}$$
 where j is the smallest integer for $\max(y) < 1$
Here [ymin, ymax] is [-1,1]



Data Preparation: Outlier Detection

Outlier Detection

 An outlier could be, for example, rare behavior, system defect, measurement error, or reaction to an unexpected event





Outlier Detection: Motivation

- Why finding outliers is important?
 - Summarize data by statistics that represent the majority of the data
 - Train a model that generalizes to new data
 - Finding the outliers can also be the focus of the analysis and not only data cleaning

Outlier Detection Techniques

- Knowledge-based
- Statistics-based
 - Distance from the median
 - Position in the distribution tails
 - Distance to the closest cluster center
 - Error produced by an autoencoder
 - Number of random splits to isolate a data point from other data







Material

Open for Innovation	Hub Blog Forum Events Careers Contact Download Q								
	SOFTWARE / SOLUTIONS / LEARNING / PARTNERS / COMMUNITY / ABOUT								
Home > About > Blog									
/ News	Four Techniques for Outlier Detection								
/ Blog	•								
/ Team	Mon, 10/01/2018 - 10:00 — admin								
/ Careers	Authors: Maarit Widmann and Moritz Heine								
/ Contact Us	Ever been skewed by the presence of outliers in your set of data? Anomalies, or outliers, can be a								
/ Travel Information	serious issue when training machine learning algorithms or applying statistical techniques. They								
/ KNIME Open Source Story	are often the result of errors in measurements or exceptional system conditions and therefore do not describe the common functioning of the underlying system. Indeed, the best practice is to implement an outlier removal phase before proceeding with further analysis.								
/ Open for Innovation									
	But hold on there! In some cases, outliers can give us information about localized anomalies in the whole system; so the detection of outliers is a valuable process because of the additional information they can provide about your dataset.								
	There are many techniques to detect and optionally remove outliers from a dataset. In this blog post, we show an implementation in KNIME Analytics Platform of four of the most frequently used - traditional and novel - techniques for outlier detection.								

https://www.knime.com/blog/four-techniques-for-outlier-detection

16



Data Preparation: Dimensionality Reduction

Is there such a thing as "too much data"?

"Too much data":

- Consumes storage space
- Eats up processing time
- Is difficult to visualize
- Inhibits ML algorithm performance
- Beware of the model: Garbage in \rightarrow Garbage out



Dimensionality Reduction Techniques

- Measure based
 - Ratio of missing values
 - Low variance
 - High Correlation
- Transformation based
 - Principal Component Analysis (PCA)
 - Linear Discriminant Analysis (LDA)
 - t-SNE
- Machine Learning based
 - Random Forest of shallow trees
 - Neural auto-encoder



Missing Values Ratio

•	Row	pclass Number (integer)	survived Number (integer)	$_{\rm Bring}^{\rm name} \sim$	sex v	$_{\rm Number(double)} \sim$	sibsp Number (integer)	parch Number (Integer) ~	ticket \checkmark	fare _{Number (double)} ~	cabin	embarked	boat \smile	$_{_{Number(integer)}}^{\rm body} \sim$	home.dest String
29	Row28	8 1	1	Bissette, M.,	female	35	0	0	PC 17760	135.633	C99	s	8	0	0
30	Row29	9 1	1	Bjornstrom	male	28	0	0	110564	26.55	C52	s	D	0	Stockholm, Sweden / Washington, DC
31	Row30	0 1	0	Blackwell,	male	45	0	0	113784	35.5	т	s	•	0	Trenton, NJ
32	Row31	1 1	1	Blank, Mr	male	40	0	0	112277	31	A31	с	7	0	Glen Ridge, NJ
33	Row32	2 1	1	Bonnell, Mi	female	30	0	0	36928	164.867	C7	S	8	0	Youngstown, OH
34	Row33	1.1	1	Bonnell, Mi	female	58	0	0	113783	26.55	C103	S	8	0	Birkdale, England Cleveland, Ohio
35	Row34	4 1	0	Borebank,	male	42	0	0	110489	26.55	D22	S	C	0	London / Winnipeg, MB
36	Row35	5 1	1	Bowen, Mi	female	45	0	0	PC 17608	262.375	O	с	4	0	Cooperstown, NY
37	Row36	5 1	1	Bowerman,	female	22	0	1	113505	55	E33	S	6	0	St Leonards-on-Sea, England Ohio
38	Row37	7 1	1	Bradley, Mr	male	۲	0	0	111427	26.55	O	s	9	0	Los Angeles, CA
39	Row38	8 1	0	Brady, Mr	male	41	0	0	113054	30.5	A21	S	C	0	Pomeroy, WA
40	Row39	9 1	0	Brandeis,	male	48	0	0	PC 17591	50.496	810	с	0	208	Omaha, NE
41	Row40	0 1	0	Brewe, Dr	male	۲	0	0	112379	39.6	O	C	C	0	Philadelphia, PA
42	Row41	1 1	1	Brown, Mrs	female	44	0	0	PC 17610	27.721	84	c	6	0	Derver, CO
43	Row42	2 1	1	Brown, Mrs	female	59	2	0	11769	51.479	C101	S	D	0	Belmont, MA
44	Row43	1.1	1	Bucknell,	female	60	0	0	11813	76.292	D15	c	8	0	Philadelphia, PA
45	Row44	4 1	1	Burns, Mis	female	41	0	0	16966	134.5	E40	C	3	0	0
46	Row45	5 1	0	Butt, Major	male	45	0	0	113050	26.55	838	\$	0	0	Washington, DC
47	Row46	5 1	0	Cairns, Mr	male	0	0	0	113798	31	O	s	0	0	0
48	Row47	7 1	1	Calderhea	male	42	0	0	PC 17476	26.288	E24	S	5	0	New York, NY
49	Row48	8 1	1	Candee, Mr.	. female	53	0	0	PC 17606	27.446	O	С	6	0	Washington, DC
50	Row49	9 1	1	Cardeza, M.	male	36	0	1	PC 17755	512.329	851 853 855	c	3	0	Austria-Hungary / Germantown, Philadelphia, P
51	Row50	0 1	1	Cardeza, M.	female	58	0	1	PC 17755	512.329	B51 B53 B55	С	3	0	Germantown, Philadelphia, PA
52	Row51	1 1	0	Carlsson,	male	33	0	0	695	5	851 853 855	S	0	0	New York, NY
53	Row52	2 1	0	Carrau, Mr	male	28	0	0	113059	47.1	O	8	0	0	Montevideo, Uruguay
54	Row53	1.1	0	Carrau, Mr	male	17	0	0	113059	47.1	0	\$	0	0	Montevideo, Uruguay



THEN remove column

IF (% missing value > threshold)



20

Low Variance

29 Rov 30 Rov	w28 1 w29 1	1				sectore (condex)	Number (integer)	String	Number (double)	String	String	String	Number (integer)	String	Low variance Filler
30 Rov	w29 1		Bissette, M	female	35	0	0	PC 17760	135.633	C99	s	8	•	0	
		1	Bjornstrom	male	28	0	0	110564	26.55	C52	\$	D	0	Stockholm, Sweden / Washington, DC	
31 Rov	w30 1	0	Blackwell,	male	45	0	0	113784	35.5	т	s	C	0	Trenton, NJ	10 M
32 Rov	w31 1	1	Blank, Mr	male	40	0	0	112277	31	A31	с	7	0	Glen Ridge, NJ	
33 Rov	w32 1	1	Bonnell, Mi	female	30	0	0	36928	164.867	C7	s	8	0	Youngstown, OH	
34 Rov	w33 1	1	Bonnell, Mi	female	58	0	0	113783	26.55	C103	s	8	0	Birkdale, England Cleveland, Ohio	
35 Rov	w34 1	0	Borebank,	male	42	0	0	110489	26.55	D22	s	0	0	London / Winnipeg, MB	
36 Rov	w35 1	1	Bowen, Mi	female	45	0	0	PC 17608	262.375	0	c	4	0	Cooperstown, NY	
37 Rov	w36 1	1	Bowerman,	female	22	0	1	113505	55	E33	s	6	0	St Leonards-on-Sea, England Ohio	
38 Rov	w37 1	1	Bradley, Mr	male	0	0	0	111427	26.55	0	\$	9	0	Los Angeles, CA	
39 Rov	w38 1	0	Brady, Mr	male	41	0	0	113054	30.5	A21	s	O	0	Pomeroy, WA	
40 Rov	w39 1	0	Brandeis,	male	48	0	0	PC 17591	50.496	810	c	0	208	Omaha, NE	
41 Rov	w40 1	0	Brewe, Dr	male	0	0	0	112379	39.6	0	С	O	0	Philadelphia, PA	
42 Rov	w41 1	1	Brown, Mrs	female	44	0	0	PC 17610	27.721	84	c	6	0	Deriver, CO	
43 Rov	w42 1	1	Brown, Mrs	female	59	2	0	11769	51.479	C101	s	D	0	Belmont, MA	
44 Rov	w43 1	1	Bucknell,	female	60	0	0	11813	76.292	D15	c	8	0	Philadelphia, PA	
45 Rov	w44 1	1	Burns, Mis	female	41	0	0	16966	134.5	E40	C	3	0	0	
46 Rov	w45 1	0	Butt, Major	male	45	0	0	113050	26.55	838	\$	0	0	Washington, DC	Note: requires min
47 Rov	w46 1	0	Cairns, Mr	male	0	0	0	113798	31	0	s	0	0	0	l Note. requires mir
48 Rov	w47 1	1	Calderhea	male	42	0	0	PC 17476	26.288	E24	\$	5	0	New York, NY	
49 Rov	w48 1	1	Candee, Mr	female	53	0	0	PC 17606	27.446	0	С	6	0	Washington, DC	max-normalizatior
50 Rov	w49 1	1	Cardeza, M	male	36	0	1	PC 17755	512.329	851 853 855	c	3	0	Austria-Hungary / Germantown, Philadelphia, PA	
51 Rov	w50 1	1	Cardeza, M	female	58	0	1	PC 17755	512.329	B51 B53 B55	С	3	0	Germantown, Philadelphia, PA	and only works to
52 Rov	w51 1	0	Carlsson,	male	33	0	0	695	5	851 853 855	\$	0	0	New York, NY	
53 Rov	w52 1	0	Carrau, Mr	male	28	0	0	113059	47.1	0	S	O	0	Montevideo, Uruguay	numeric columns
54 Rov	w53 1	0	Carrau, Mr	male	17	0	0	113059	47.1	0	\$	0	0	Montevideo, Uruguay	

- If column has constant value (variance = 0), it contains no useful information
- In general: IF (variance < threshold) THEN remove column</p>



High Correlation

- Two highly correlated input variables probably carry similar information
- IF (corr(var1, var2) > threshold) => remove var1







Principal Component Analysis (PCA)

 PCA is a statistical procedure that **orthogonally** transforms the original *n* coordinates of a data set into a new set of *n* coordinates, called principal components.

 $(PC_1, PC_2, \cdots PC_n) = PCA(X_1, X_2, \cdots X_n)$

- The first principal component PC₁ follows the direction (eigenvector) of the largest possible variance (largest eigenvalue of the covariance matrix) in the data.
- Each succeeding component PC_k follows the direction of the **next largest possible variance** under the constraint that it is orthogonal to (i.e., uncorrelated with) the preceding components $(PC_1, PC_2, \cdots PC_{k-1})$.

If you're still curious, there's LOTS of different ways to think about PCA: <u>https://stats.stackexchange.com/questions/2691/making-sense-of-principal-component-analysis-eigenvectors-eigenvalues</u>



Image from Wikipedia



Principal Component Analysis (PCA)

- PC₁ describes most of the variability in the data, PC₂ adds the next big contribution, and so on. In the end, the last PCs do not bring much more information to describe the data.
- Thus, to describe the data we could use only the top m < n (i.e., $PC_1, PC_2, \cdots PC_m$) components with little if any loss of information
- Caveats:
 - Results of PCA are quite difficult to interpret
 - Normalization required
 - Only effective on numeric columns

Dimensionality Reduction





Linear Discriminant Analysis (LDA)

 LDA is a statistical procedure that **orthogonally** transforms the original *n* coordinates of a data set into a new set of *k-1* coordinates, called linear discriminants, where *k* is the number of classes in the class variable

$$(LD_1, LD_2, \cdots LD_{k-1}) = LDA(X_1, X_2, \cdots X_n)$$

Here, however, discriminants (components)
 maximize the separation between classes



- PCA : unsupervised
- LDA : supervised





Linear Discriminant Analysis (LDA)

- LD₁ describes best the class separation in the data, LD₂ adds the next big contribution, and so on. In the end, the last LDs do not bring much more information to separate the classes.
- Thus, for our classification problem we could use only the top m < k 1 (i.e., $LD_1, LD_2, \cdots LD_m$) discriminants with little if any loss of information
- Caveats:
 - Results of LDA are quite difficult to interpret
 - Normalization required
 - Only effective on numeric columns

Dimensionality Reduction



Ensembles of Shallow Decision Trees

- Often used for classification, but can be used for feature selection too
- Generate a large number (we used 2000) of trees that are very shallow (2 levels, 3 sampled features)
- Calculate the statistics of candidates and selected features.
 The more often a feature is selected in such trees, the more likely it contains predictive information
- Compare the same statistics with a forest of trees trained on a random dataset.







Autoencoder

 Feed-Forward Neural Network architecture with encoder / decoder structure. The network is trained to reproduce the input vector onto the output layer.



- That is, it compresses the input vector (dimension n) into a smaller vector space on layer "code" (dimension m<n) and then it reconstructs the original vector onto the output layer.
- If the network was trained well, the reconstruction operation happens with minimal loss of information.



Data Preparation: Feature Selection

Feature Selection vs. Dimensionality Reduction

- Both methods are used for reducing the number of features in a dataset. However:
- Feature selection is simply selecting and excluding given features without changing them.
- Dimensionality reduction **might transform** the features into a lower dimension.
- Feature selection is often a somewhat more aggressive and more computationally expensive process.
 - Backward Feature Elimination
 - Forward Feature Construction



Backward Feature Elimination (greedy top-down)

- 1. First train one model on *n* input features
- 2. Then train *n* separate models each on n 1 input features and remove the feature whose removal produced the least disturbance
- 3. Then train n 1 separate models each on n 2 input features and remove the feature whose removal produced the least disturbance
- 4. And so on. Continue until desired maximum error rate on *training* data is reached.

Backward Feature Elimination

△ Dialog - 6:3 - Feature Selection Filter			-		×			
hile								
Column Selection Flow Variables Job Ma	mager Selection Memory Policy					-		
Include static columns Select features manually								
						redictor Scorer	Feature Selection Loop End	
Select features automatically by score	threshold							
Prediction score threshold 0.85								
Ontimization October The space is being	li internet de la constante de					000		
Opumization Criterion: The score is being i	naximizeo.				- II.		Choose the variable	n Fil
Accuracy	Nr. of features	S Sentiment Analysis					to optimize	•
0	.885	6 I SentimentRating						
	878	S MaritalStatus					000	
0	.864	4 S Gender						
0	.861	8 L Acc						
c	.854	3 D ChuroScore						
	.752	2 S Products						
	0.47	13						
L								
		OK Apply Car	cel (2				



Forward Feature Construction (greedy bottom-up)

- 1. First, train *n* separate models on one single input feature and keep the feature that produces the best accuracy.
- 2. Then, train n 1 separate models on 2 input features, the selected one and one more. At the end keep the additional feature that produces the best accuracy.
- 3. And so on ... Continue until an acceptable error rate is reached.



Material



https://thenewstack.io/3-new-techniques-for-data-dimensionality-reduction-in-machine-learning/



Data Preparation: Feature Engineering

Feature Engineering: Motivation

Sometimes transforming the original data allows for better discrimination by ML algorithms.





Feature Engineering: Techniques

- Coordinate Transformations
 Remember PCA and LDA?
 Polar coordinates , ...
- Distances to cluster centres, after data clustering
- Simple math transformations on single columns (e^x, x², x³, tanh(x), log(x), ...)
- Combining together multiple columns in math functions
 (f(x₁, x₂, ... xn), x₂ x₁, ...)
- The whole process is domain dependent





Feature Engineering in Time Series Analysis

- Second order differences: y = x(t) x(t-1) & y'(t) = y(t) y(t-1)
- Logarithm: log(y'(t))





Regression Problems

Supervised Learning: Classification vs. Regression

- $X = (x_1, x_2)$ and $y = \{label 1, ..., label n\}$ or $y \in \mathbb{R}$
- A training set with many examples of (X, y)
- The model learns on the examples of the training set to produce the right value of y for an input vector X

Classification

y = {yellow, gray}
y = {churn, no churn}
y = {increase, unchanged, decrease}
y = {blonde, gray, brown, red, black}
y = {job 1, job 2, ..., job n}

Numerical Predictions (Regression)

y = temperature y = number of visitors y = number of kW y = price y = number of hours



Regression Overview

- Goal: Explain how target attribute depends on descripitive attributes
 - Target attribute
 - Descriptive attribute(s)
- → Response variable
- → Regressor variable(s)
- As a parameterized function class *f*
 - Estimate parameters to describe the relationship
 - Must be simple enough for interpolation and extrapolation purposes
 - Example:

Line (black) v.s. Polynomial (blue) with degree 7





Regression

Predict numeric outcomes on existing data (supervised)

Applications

- Forecasting
- Quantitative Analysis

Methods

- Linear
- Polynomial
- Regression Trees
- Partial Least Squares



Statistics on Linear Regression

Variable	Coeff.	Std. Err.	t-value	P>ltl
Petal.Length	0.4158	0.0096	43.3872	0.0
Intercept	-0.3631	0.0398	-9.1312	4.44E-16







Linear Regression Algorithm

Regression Line

- Given a data set with two continuous attributes, x and y
- There is an approximate linear dependency between x and y





Regression Line

- Given a data set with two continuous attributes, x and y
- There is an approximate linear dependency between x and y



- We find a **regression line** (i.e., determine the parameters *a* and *b*) such that the line fits the data as well as possible
- Examples:
 - Trend estimation (e.g., oil price over time)

45

- Epidemiology (e.g., cigarette smoking vs. lifespan)
- Finance (e.g., return on investment vs. return on all risky assets)
- Economics (e.g., spending vs. available income)



Linear Regression

Predicts the values of the target variable y based on a linear combination of the values of the input feature(s) x_j Two input features: $\hat{y} = a_0 + a_1x_1 + a_2x_2$ p input features: $\hat{y} = a_0 + a_1x_1 + a_2x_2 + \dots + a_px_p$

- Multiple regression: several input features \rightarrow regression hyper-plane
- Residuals: differences between observed and predicted values (errors)
 Use the residuals to measure the model fit



Simple Linear Regression

Optimization goal: minimize sum of squared residuals



KNIME

Simple Linear Regression

• Think of a straight line $\hat{y} = f(x) = a + bx$

• Find a and b to model all observations (x_i, y_i) as close as possible

→ SSE $F(a, b) = \sum_{i=1}^{n} (f(x) - y_i)^2 = \sum_{i=1}^{n} (a + bx_i - y_i)^2$ should be minimal

That is:

$$\frac{\partial F}{\partial a} = \sum_{i=1}^{n} 2(a + bx_i - y_i) = 0$$
$$\frac{\partial F}{\partial b} = \sum_{i=1}^{n} 2(a + bx_i - y_i) x_i = 0$$

 \rightarrow A unique solution exists for *a* and *b*



Linear Regression

Optimization goal: minimize the squared residuals

$$\sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \sum_{j=0}^{n} a_j x_{j,i})^2 = (y - aX)^T (y - aX)$$

Solution:

$$\hat{a} = (X^T X)^{-1} X^T y$$



- Computational issues:
 - X^TX must have full rank, and thus be invertible (Problems arise if linear dependencies between input features exist)
 - Solution may be unstable, if input features are almost linearly dependent



Linear Regression: Summary

- Positive:
 - Strong mathematical foundation
 - Simple to calculate and to understand (For moderate number of dimensions)
 - High predictive accuracy (In many applications)
- Negative:
 - Many dependencies are non-linear (Can be generalized)
 - Model is global and cannot adapt well to locally different data distributions But: Locally weighted regression, CART



Regression vs. Time Series Analysis

- Regression
 - Targets & set of input features
 - Describing the relationship between the target and input features
 - Model → interpolation

- Time series analysis
 - Sequence of observations
 - Predicting future obs from
 - Existing time series
 - Accompanying time series
 - Model → extrapolation





KNIMF

Polynomial Regression

Predicts the values of the target variable ybased on a polynomial combination of degree d of the values of the input feature(s) x_j

$$\tilde{y} = a_0 + \sum_{j=1}^p a_{j,1} x_j + \sum_{j=1}^p a_{j,2} x_j^2 + \dots + \sum_{j=1}^p a_{j,d} x_j^d$$

- Simple regression: one input feature \rightarrow regression curve
- Multiple regression: several input features \rightarrow regression hypersurface
- Residuals: differences between observed and predicted values (errors)
 Use the residuals to measure the model fit



Evaluation of Regression Models

Numeric Errors: Formulas

Error Metric	Formula	Notes
R-squared	$1 - \frac{\sum_{i=1}^{n} (y_i - f(x_i))^2}{\sum_{i=1}^{n} (y_i - \overline{y})^2}$	Universal range: the closer to 1 the better
Mean absolute error (MAE)	$\frac{1}{n}\sum_{i=1}^{n} y_i - f(x_i) $	Equal weights to all distances Same unit as the target column
Mean squared error (MSE)	$\frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$	Common loss function
Root mean squared error (RMSE)	$\sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2}$	Weights big differences more Same unit as the target column
Mean signed difference	$\frac{1}{n}\sum_{i=1}^{n} (y_i - f(x_i))$	Only informative about the direction of the error
Mean absolute percentage error (MAPE)	$\frac{1}{n} \sum_{i=1}^{n} \frac{ y_i - f(x_i) }{ y_i }$	Requires non-zero target column values



MAE (Mean Absolute Error) vs. RMSE (Root Mean Squared Error)

MAE	RMSE
Easy to interpret - mean absolute error	Cannot be directly interpreted as the average error
All errors are equally weighted	Larger errors are weighted more
Generally smaller than RMSE	Ideal when large deviations need to be avoided

Example: Actual values = [2,4,5,8], Case 1: Predicted Values = [4, 6, 8, 10] Case 2: Predicted Values = [4, 6, 8, 14]





R-squared vs. RMSE

R-squared	RMSE
Relative measure : Proportion of variability explained by the model	Absolute measure : How much deviation at each point
Range: Usually between 0 and 1. 0 = no variability explained 1 = all variability explained	Same scale as the target

Examp	le:

Actual values = [2,4,5,8],

Case 1: Predicted Values = [3, 4, 5, 6]

Case 2: Predicted Values = [3, 3, 7, 7]

	R-sq	RMSE
Case 1	0.96	1.12
Case 2	0.65	1.32



Numeric Scorer

- Similar to scorer node, but for nodes with *numeric* predictions
- Compare dependent variable values to predicted values to evaluate model quality.
- Report R², RMSE, MAPE, etc.



A Statistics - 3:30	_		\times			
File						
R2:		0.786				
Mean absolute error:		23,535.03				
Mean squared error:		1,311,070,583.795				
Root mean squared error:		36,208.709				
Mean signed difference:	signed difference: 1,000.836					
Mean absolute percentage erro	or:	0.145				
Adjusted R ² :		0.786				



Regression Tree

Regression Tree: Goal







Regression Tree: Initial Split





Regression Tree: Initial Split





Regression Tree: Growing the Tree





Regression Tree: Final Model





Regression Tree: Algorithm

Start with a single node containing all points.

- 1. Calculate c_i and E_i .
- 2. If all points have the same value for feature x_i , stop.



- 3. Otherwise, find the best binary splits that reduces $E_{i,s}$ as much as possible.
 - $E_{j,s}$ doesn't reduce as much \rightarrow stop
 - A node contains less than the minimum node size \rightarrow stop
 - Otherwise, take that split, creating two new nodes.
 - In each new node, go back to step 1.



Regression Trees: Summary

- Differences to decision trees:
 - Splitting criterion: minimizing intra-subset variation (error)
 - Pruning criterion: based on numeric error measure
 - Leaf node predicts average target values of training instances reaching that node
- Can approximate piecewise constant functions
- Easy to interpret



Regression Trees: Pros & Cons

- Finding of (local) regression values (average)
- Problems:
 - No interpolation across borders
 - Heuristic algorithm: unstable and not optimal.
- Extensions:
 - Fuzzy trees (better interpolation)
 - Local models for each leaf (linear, quadratic)

Exercises

Exercises

- Dataset: Sales data of individual residential properties in Ames, Iowa from 2006 to 2010
- One of the columns is the price for which the house was sold
- Goal: Predicting the house price
- Data Preparation:
 - 01_Missing_Value_Handling
 - 02_Outlier_Detection
- Regression:
 - 03_Linear_Regression
 - 04_Regression_Trees

