



To die, or not to die

My experiences from competing at a
data science competition

Petri Roine | Senior Advisor

"Our commitment to innovation is why last year analysts named SAS a leader in more than 30 reports and what drives our constant curiosity into what the future holds."

Dr. Jim Goodnight

Den lille Havfrue souvenir.
Original statue by Edvard Eriksen.

København

Welcome to Kaggle Competitions

Challenge yourself with real-world machine learning problems



New to Data Science?

Get started with a tutorial on our most popular competition for beginners, [Titanic: Machine Learning from Disaster](#).



Build a Model

Get the data & use whatever tools or methods you prefer to make predictions.



Make a Submission

Upload your prediction file for real-time scoring & a spot on the leaderboard.

✖ Dismiss

11 active competitions		Sort By	Prize
Active	All	Entered	All Categories
	Data Science Bowl 2017 Can you improve lung cancer detection? Featured - 13 days to go		\$1,000,000 1,821 teams
	The Nature Conservancy Fisheries Monitoring Can you detect and classify species of fish? Featured - 13 days to go		\$150,000 2,213 teams
	Intel & MobileODT Cervical Cancer Screening Which cancer treatment will be most effective? Featured - 3 months to go		\$100,000 212 teams
	Google Cloud & YouTube-8M Video Understanding Challenge Can you produce the best video tag predictions? Featured - 2 months to go		\$100,000 370 teams

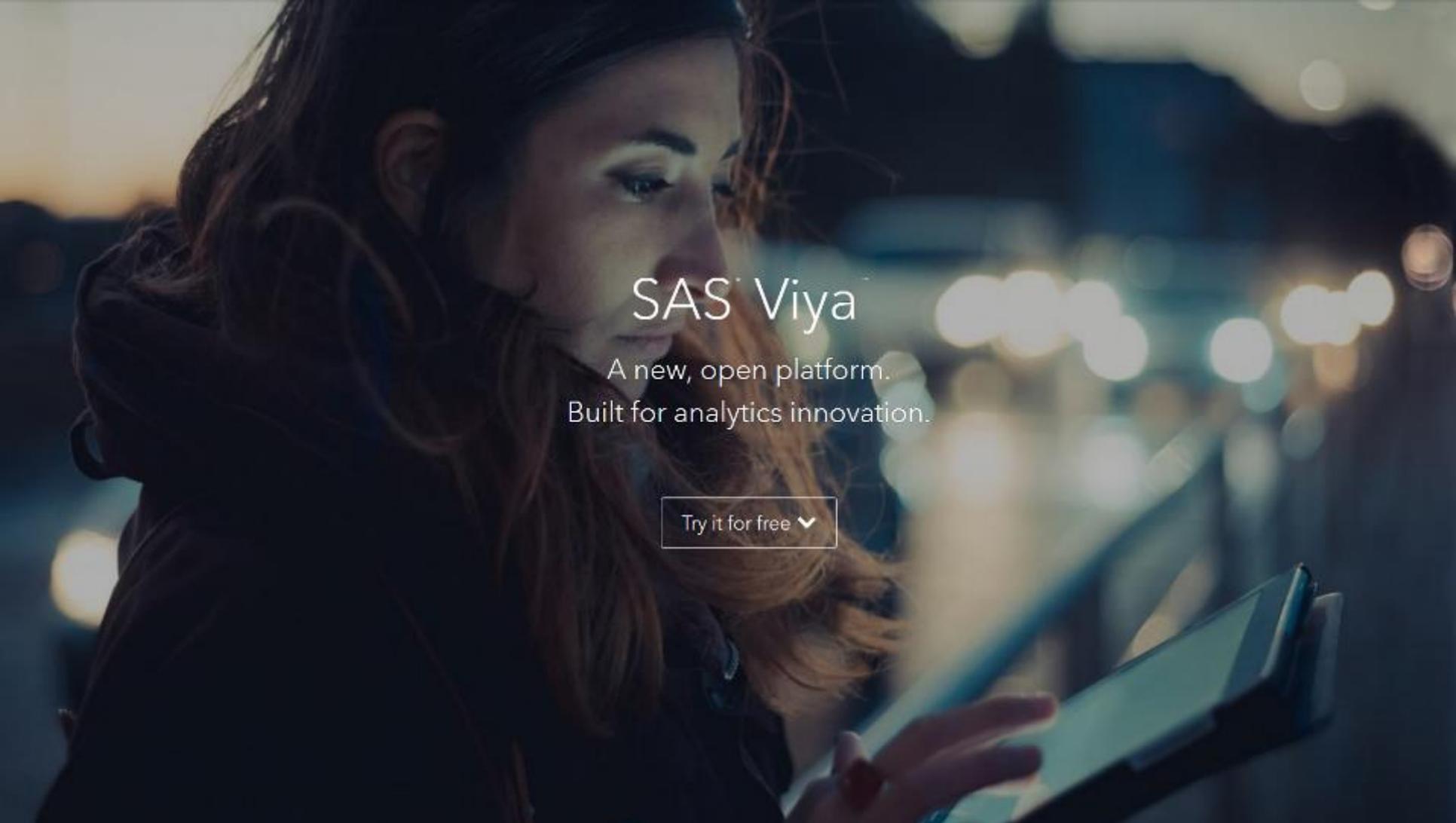


TITANIC

891



418

A woman with long brown hair is looking at a tablet in her hands. The background is a blurred city street at night with warm lights.

SAS Viya

A new, open platform.
Built for analytics innovation.

Try it for free 



SAS[®] Visual Data Mining and Machine Learning

Boost analytical productivity and solve complex problems faster with a single, integrated in-memory environment that's both open and scalable.



SAS[®] Visual Investigator

Address all your intelligence analysis and investigation management needs with a single, cloud-based solution that uses advanced analytics and machine learning technology.



SAS[®] Visual Analytics

Visually explore relationships and patterns in data smartly, quickly and easily – using interactive data discovery. And illuminate critical insights with self-service analytics.



SAS[®] Visual Statistics

Delve deeper into different types of data with a powerful, in-memory solution that lets you create, compare and refine descriptive and predictive models on the fly.



SAS[®] Visual Forecasting

Manage organizational planning challenges and automate large-scale hierarchical forecasting with a solution that supports open source and SAS coding in a single environment.



SAS[®] Optimization

Evaluate alternative actions and scenarios with a powerful array of optimization modeling capabilities and solution techniques.



SAS[®] Econometrics

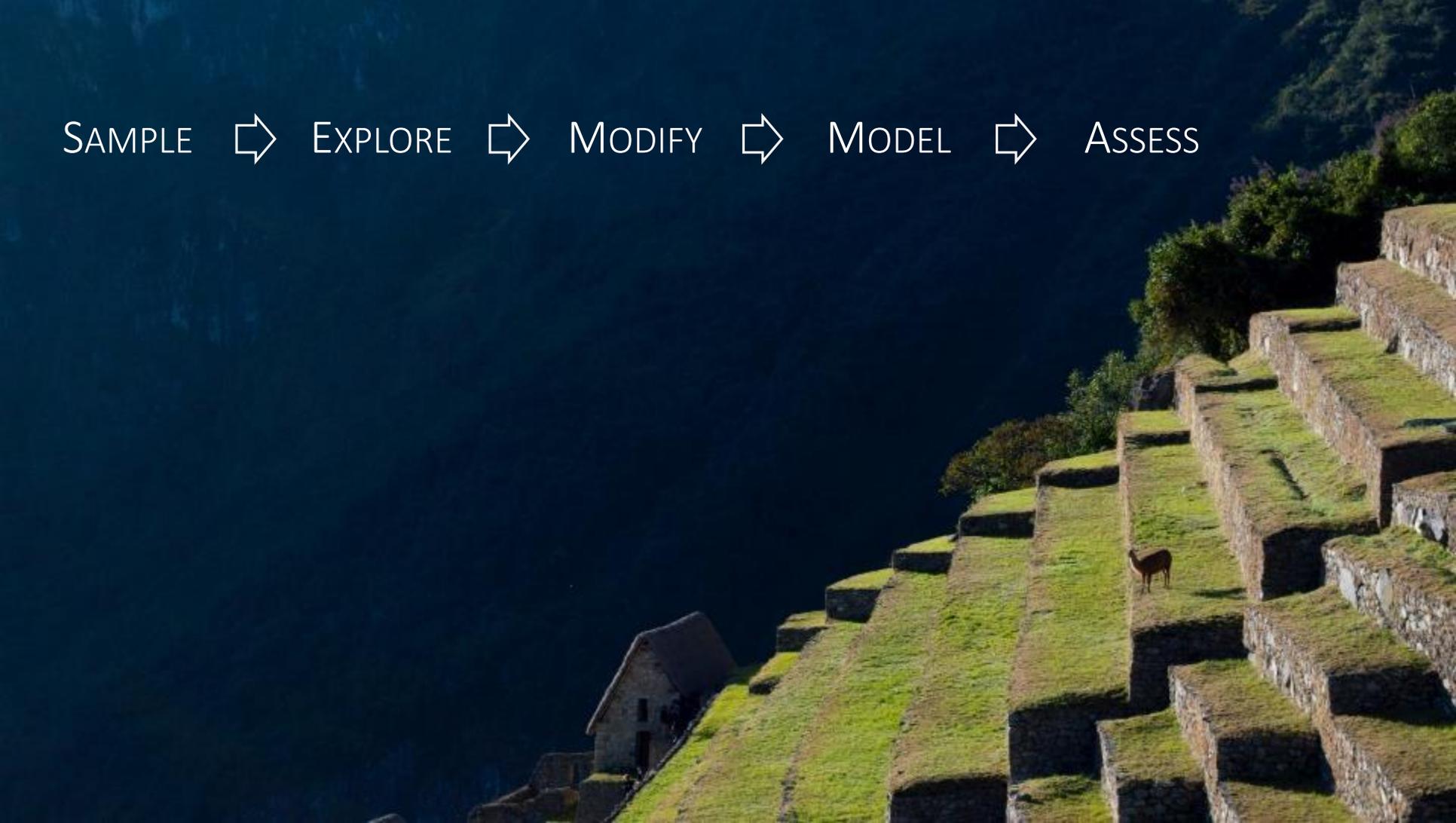
Analyze, describe and predict issues related to economic and financial systems by applying advanced mathematical and statistical methods.



SAS[®] Event Stream Processing

Analyze high-velocity big data in event streams. Train models, score data and take immediate action on what's relevant.

SAMPLE ⇨ EXPLORE ⇨ MODIFY ⇨ MODEL ⇨ ASSESS





SAMPLE

Sampling

DATA DICTIONARY

Variable	Definition	Key
survival	Survival	0 = No, 1 = Yes
pclass	Ticket class	1 = 1st, 2 = 2nd, 3 = 3rd
sex	Sex	Male, female
age	Age in years	0 - 80
sibsp	# of siblings / spouses aboard	0 - 8
parch	# of parents / children aboard	0 - 6
ticket	Ticket number	Unique string
fare	Passenger fare	0 - 512
cabin	Cabin number	C85, E46,...
embarked	Port of Embarkation	C = Cherbourg Q = Queenstown S = Southampton

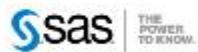
Sign In to SAS®

User ID:

Password:

Sign In

[Sign In](#)



[About](#)

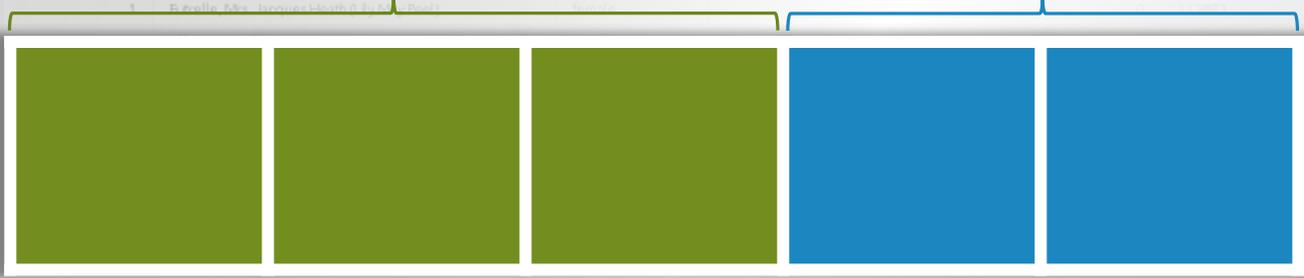
TRAIN  Server: cas-shared-default      

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
1	0	3	Braund, Mr. Owen Harris	male	22	1	0	A/5 21171	7.25
2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	1	0	PC 17599	71.2833
3	1	3	Hakkinen, Miss. Laina	female	26	0	0	STON/O2. 3101282	7.925
4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1	0	113803	53.1
5	0	3	Allen, Mr. William Henry	male	35	0	0	373450	8.05
6	0	3	Moran, Mr. James	male	.	0	0	330877	8.4583
7	0	1	McCarthy, Mr. Timothy J	male	54	0	0	17463	51.8625
8	0	3	Palsson, Master. Gosta Leonard	male	2	3	1	349909	21.075
9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27	0	2	347742	11.1333
10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14	1	0	237736	30.0708
11	1	3	Sandstrom, Miss. Marguerite Rut	female	4	1	1	PP 9549	16.7
12	1	1	Bonnell, Miss. Elizabeth	female	58	0	0	113783	26.55
13	0	3	Saunderscock, Mr. William Henry	male	20	0	0	A/S. 2151	8.05
14	0	3	Andersson, Mr. Anders Johan	male	39	1	5	347082	31.275
15	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	14	0	0	350406	7.8542
16	1	2	Hewlett, Mrs. (Mary D Kingcome)	female	55	0	0	248706	16
17	0	3	Rice, Master. Eugene	male	2	4	1	382652	29.125
18	1	2	Williams, Mr. Charles Eugene	male	.	0	0	244373	13
19	0	3	Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele)	female	31	1	0	345763	18
20	1	3	Maselman, Mrs. Fatima	female	.	0	0	2649	7.225
21	0	2	Fynney, Mr. Joseph J	male	35	0	0	239865	26
22	1	2	Beesley, Mr. Lawrence	male	34	0	0	248698	13
23	1	3	McGowan, Miss. Anna "Annie"	female	15	0	0	330923	8.0292
24	1	1	Sloper, Mr. William Thompson	male	28	0	0	113788	35.5
25	0	3	Palsson, Miss. Torborg Danira	female	8	3	1	349909	21.075
26	1	3	Asplund, Mrs. Carl Oscar (Selma Augusta Emilia Johansson)	female	38	1	5	347077	31.3875
27	0	3	Emir, Mr. Farred Chehab	male	.	0	0	2631	7.225
28	0	1	Fortune, Mr. Charles Alexander	male	19	3	2	19950	263
29	1	3	O'Dwyer, Miss. Ellen "Nellie"	female	.	0	0	330959	7.8792
30	0	3	Todoroff, Mr. Lalo	male	.	0	0	349216	7.8958

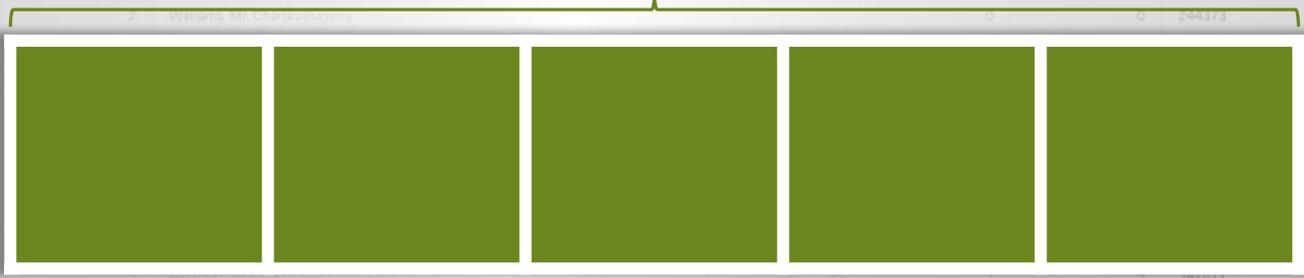
PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	ParCh	Fare
1	0	3	Braund, Mr. Owen Harris	male	22	1	0	7.25
2	1	1	Cumings, Mrs. John Bradley	female	35	1	0	53.1
3	1	3	Häkkinen, Miss. Suoma	female	26	0	0	7.925
4	1	1	Swabe, Mrs. James (Hatch) Road	female	32	0	0	53.1
5	0	3	Allen, Mr. William Victor	male	29	0	0	8.05
6	0	3	Baker, Mr. John H. Jr.	male	34	1	0	8.4583
7	0	3	Bell, Mrs. Josephine	female	38	1	0	51.8625
8	0	3	Berglund, Mr. Carl Olof	male	26	0	0	21.075
9	1	1	Bonnell, Mr. Edward Davis	male	35	0	0	11.1333
10	1	1	Booth, Mr. Walter	male	19	0	0	30.0708
11	1	1	Bray, Mr. James	male	29	0	0	16.7
12	1	1	Bray, Mrs. James (Lynch) Esq.	female	24	0	0	26.55
13	0	3	Saunderscock, Mr. William Henry	male	20	0	0	8.05
14	0	3	Andersson, Mr. Anders Johan	male	29	1	5	31.275
15	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	14	0	0	7.8542
16	1	2	Hewlett, Mrs. (May D) Kingcome	female	65	0	0	16
17	0	3	Rice, Master Eugene	male	2	4	1	29.125
18	1	2	Wilsons, Mr. Charles Eugene	male	0	0	0	13
19	0	3	Woodward, Mr. Robert Foy	male	30	0	0	18
20	1	1	Yamamoto, Mr. Tadamasa	male	35	0	0	7.225
21	0	3	Young, Mr. William	male	32	0	0	26
22	1	1	Young, Mrs. Lavinia	female	32	0	0	13
23	1	1	Young, Miss. Lavinia	female	22	0	0	8.0292
24	1	1	Young, Miss. Lavinia	female	22	0	0	35.5
25	0	3	Young, Mr. William	male	32	0	0	21.075
26	1	1	Young, Mrs. Lavinia	female	32	0	0	31.3875
27	0	3	Emy, Mr. Parred Oskarab	male	25	0	0	7.225
28	0	1	Fortune, Mr. Charles Alexander	male	19	3	2	263
29	1	3	O'Dwyer, Miss. Ellen "Nellie"	female	22	0	0	7.8792
30	0	3	Todoroff, Mr. Lalio	male	26	0	0	7.8958

Training set

Validation set



Training set



Server: cas-shared-default

Fare	Cabin	Embarked	Validate
7.25		S	1
71.2833	C85	C	0
7.925		S	1
53.1	C123	S	0
			1
			0
			1
			0
			1
			0
16.7	G6	S	1
26.55	C103	S	0
8.05		S	1
31.275		S	0

$$f(x) = \text{Train/Validate}$$

$$f(x) = \text{Floor} \left(\text{PassengerId} - x * \text{Floor} \left(\frac{\text{PassengerId}}{x} \right) \right)$$

$$f(2) \approx 50/50$$

$$f(1.4) \approx 70/30$$

$$f(1.25) \approx 80/20$$

Plan

Add calculated column

Input table:

TRAIN

Output table:

TRAIN

expression:

\'PassengerId\' - (2 * Floor ((\...

Output column:

Validate

inColumns:

[\'PassengerId\', (2 * Floor ((Passe...

Column format:



TRAIN_NEW



Server: cas-shared-default

Table Profile Column Profile

Rename column:

Validate

Change data type:

Double

Change column format: ?

COMMAS.



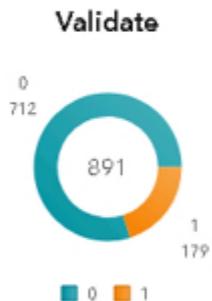
Unique Groups

Column Metrics

Standard Metrics

Advanced Metrics

Distribution



Unique (n):	2	Missing:	0
Minimum:	0	Maximum:	1
Mean:	0.20	Sum:	179.00

SibSp	Parch	Ticket	Fare	Cabin	Embarked	Validate
1	0	A/5 21171	7.25		S	1
1	0	PC 17599	71.2833	C85	C	0
0	0	STON/O2. 3101282	7.925		S	0
1	0	113803	53.1	C123	S	0
0	0	373450	8.05		S	0



DATA EXPLORATION

Visualization

Variable Identification

Missing Values Imputation

Outliers Treatment

Page 1

Page 2



Data

TRAIN_NEW

Add



Category

- Age - 89
- Cabin - 148
- Embarked - 4
- Embarked_new - 3
- Name - 891
- Parch - 7
- Passengerid - 891
- Pclass - 3
- Sex - 2
- SibSp - 7
- Survived - 2
- Ticket - 681
- Validate - 2

Measure

- Fare
- Frequency

Aggregated Measure

- Frequency Percent



Drop a data item or content here

Roles

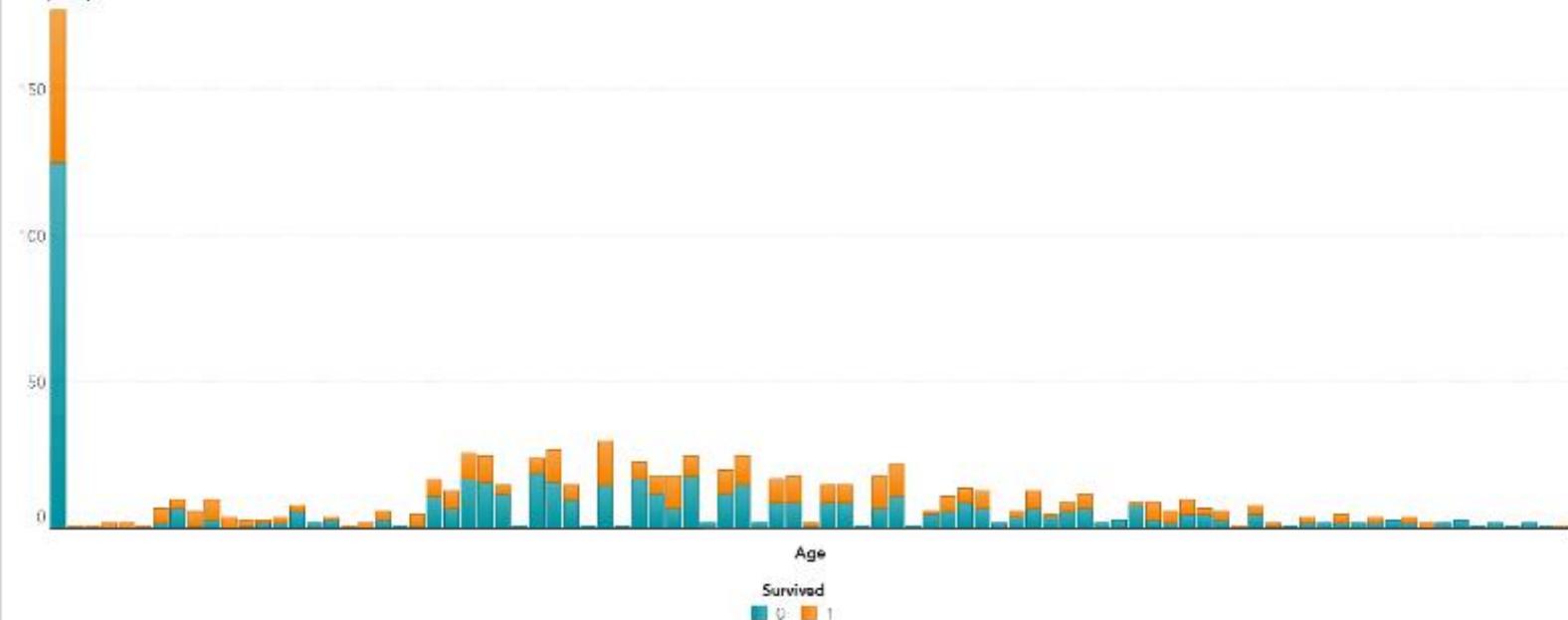
Select an object to see its roles



Explore

PASSANGER AGE

Frequency

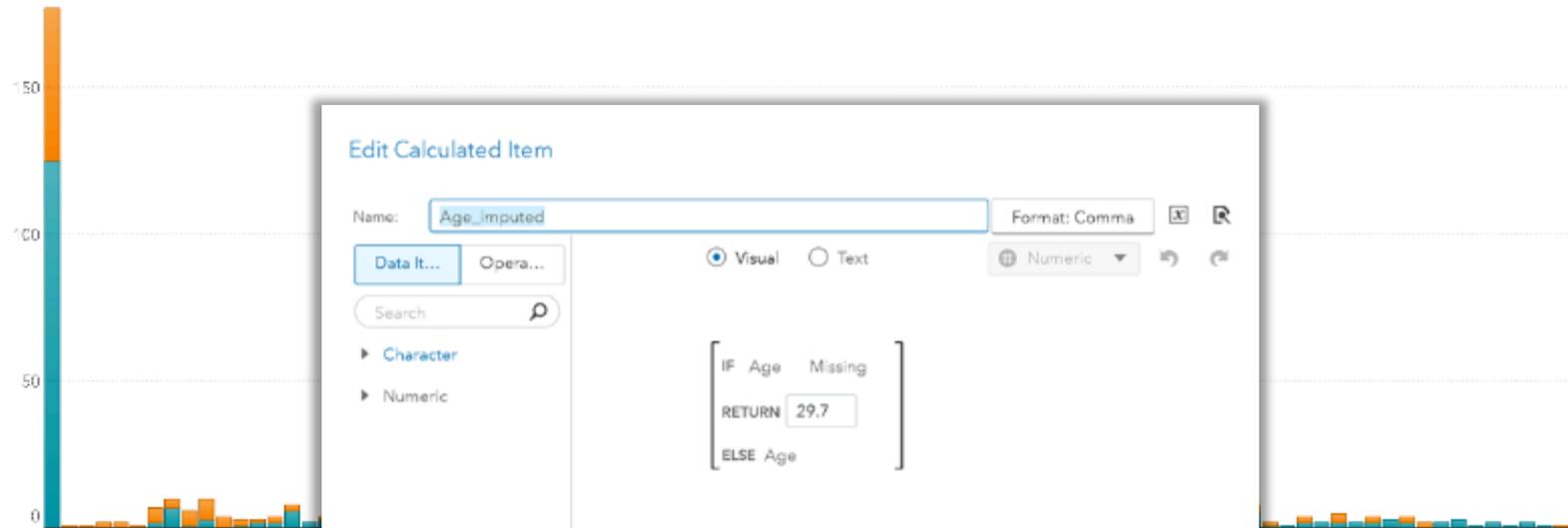


Age	Frequency	Survived
0	125	0
1	50	1
0.2	1	1
0.67	1	1
0.75	2	1
0.83	2	1
0.92	1	1
1.00	2	0

Explore

PASSANGER AGE:

Frequency



Edit Calculated Item

Name:

Format: Comma

Data It...

Opera...

 Visual Text

Numeric

Search 

▶ Character

▶ Numeric

```

[ IF Age Missing
  RETURN 29.7
  ELSE Age
]

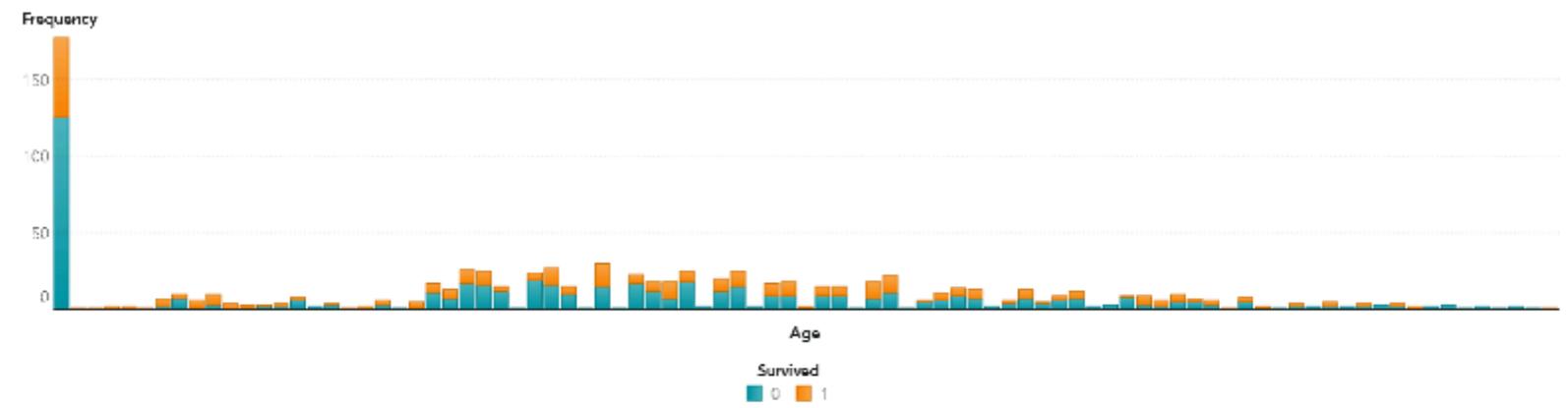
```

OK

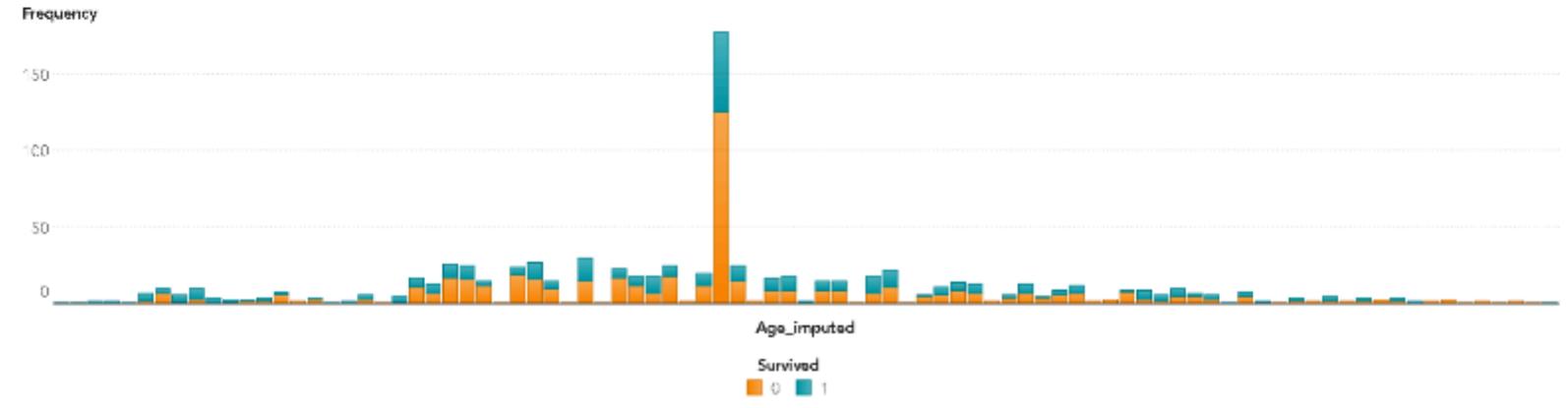
Cancel

Age	Frequency	Survived
.	125	0
.	59	1
0.42	1	1
0.67	1	1
0.75	2	1
0.83	2	1
0.92	1	1
1.00	2	0

PASSANGER AGE



PASSANGER AGE (imputed)

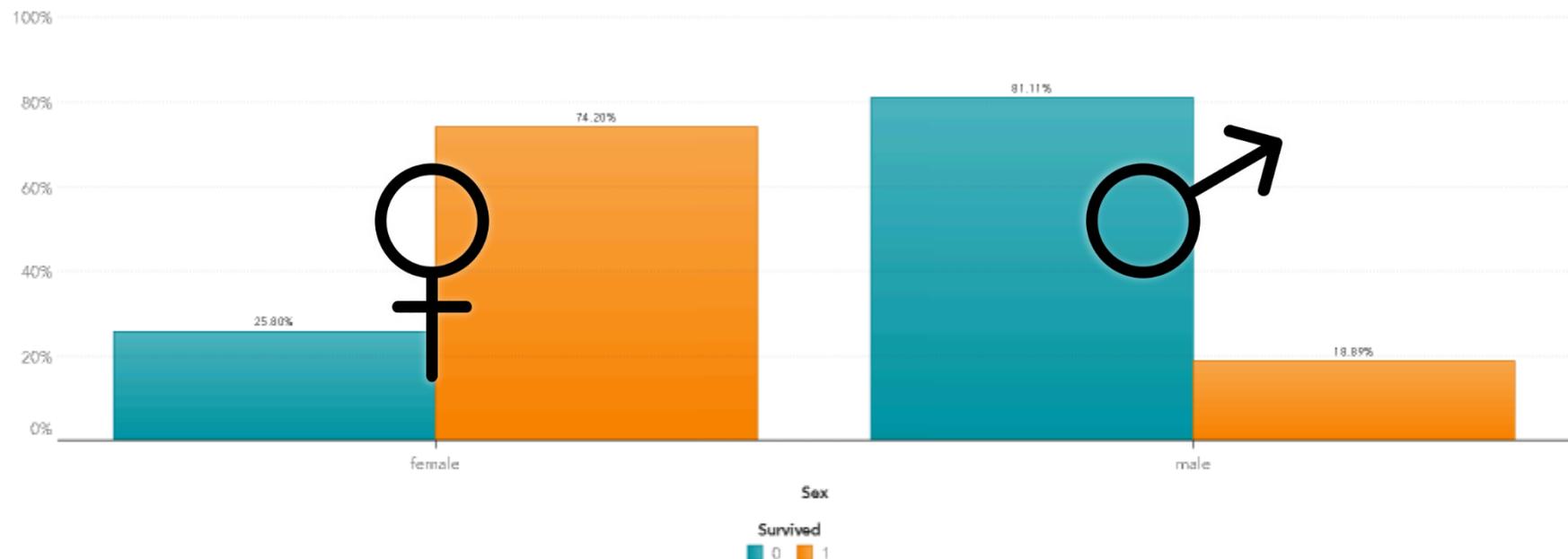


Explore



Survival per Sex (Normalized groups)

Frequency



Sex	Frequency	Survived
fem...	81	0
fem...	233	1
male	468	0
male	109	1



FEATURE ENGINEERING

Variable transformation

Variable / Feature creation

Page 1

Page 2



Data

TRAIN_NEW

Add



Category

- Age - 89
- Cabin - 148
- Embarked - 4
- Embarked_new - 3
- Name - 891
- Parch - 7
- Passengerid - 891
- Pclass - 3
- Sex - 2
- SibSp - 7
- Survived - 2
- Ticket - 681
- Validate - 2

Measure

- Fare
- Frequency

Aggregated Measure

- Frequency Percent



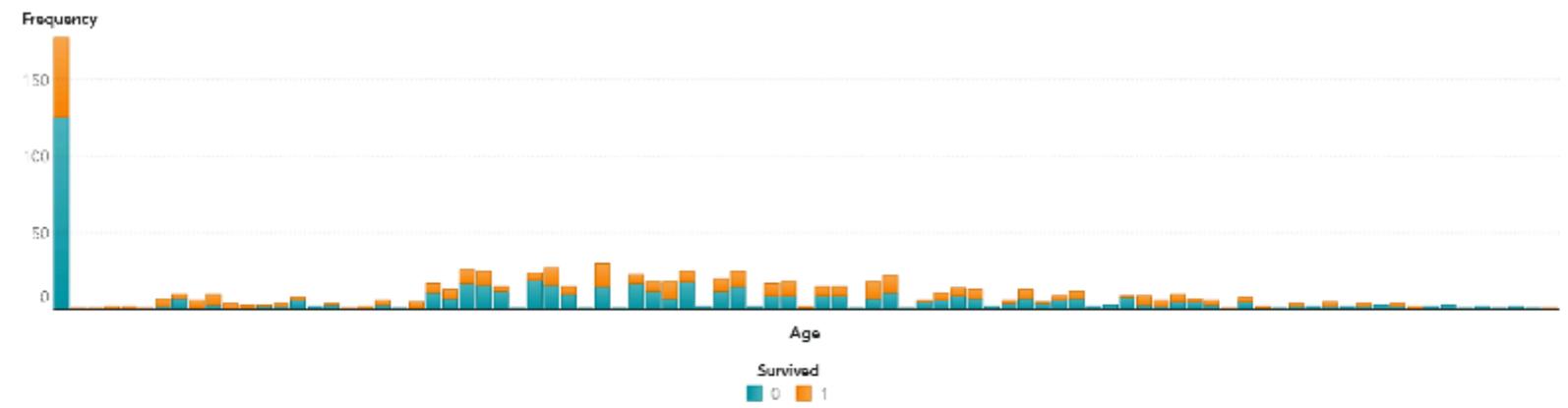
Drop a data item or content here

Roles

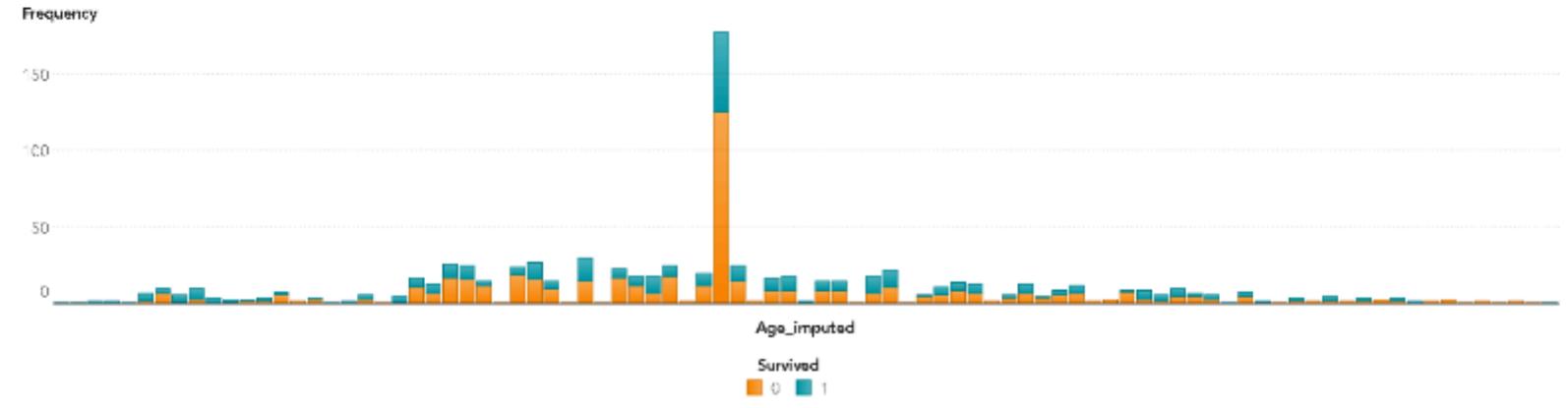
Select an object to see its roles



PASSANGER AGE



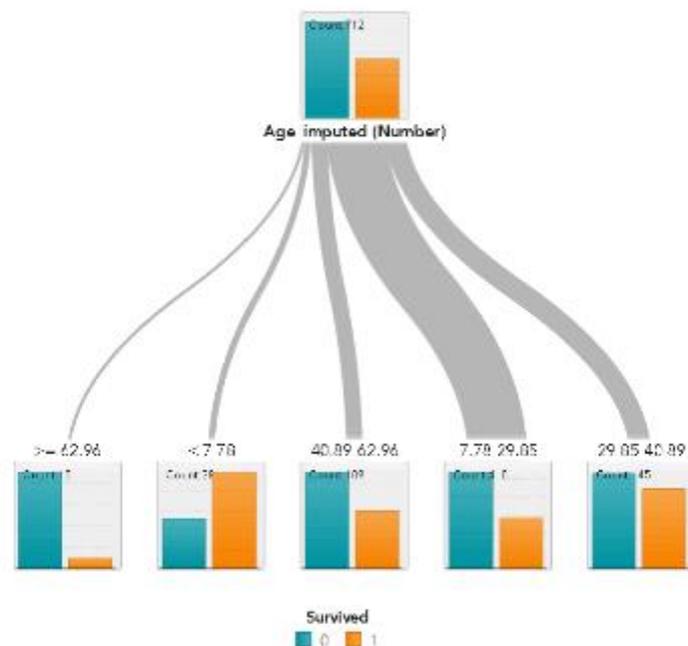
PASSANGER AGE (imputed)



Explore

Decision Tree Survived Observations used 712

Tree

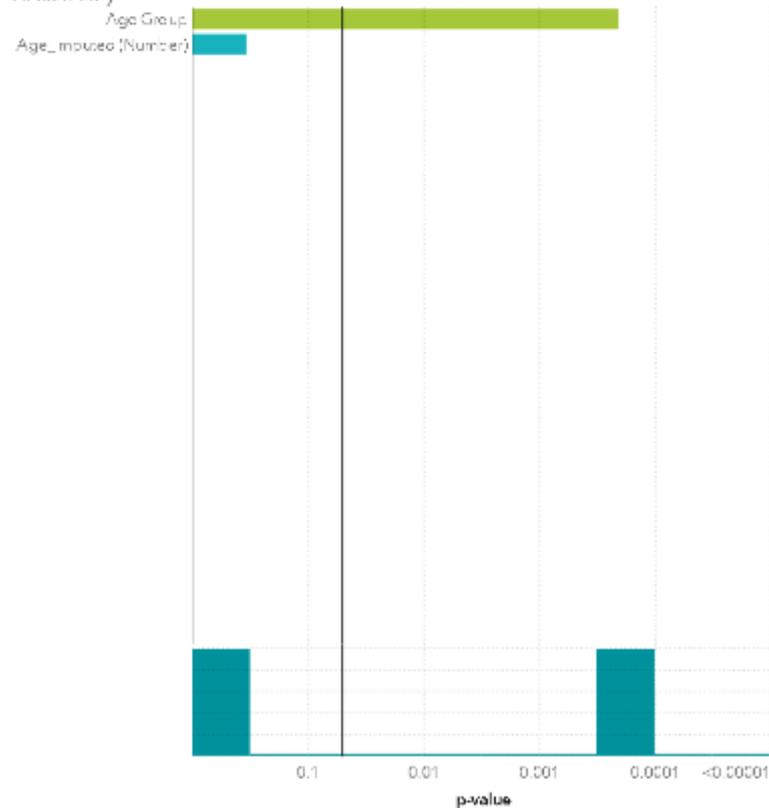


Age imputed (Number)

Raw vs Binning

Logistic Regression Survived (raw=1) R-Square 0.0302 Observations Used 891

Fit Summary



Fit

Overview

Risks

Actions

Filter

Flags

Files

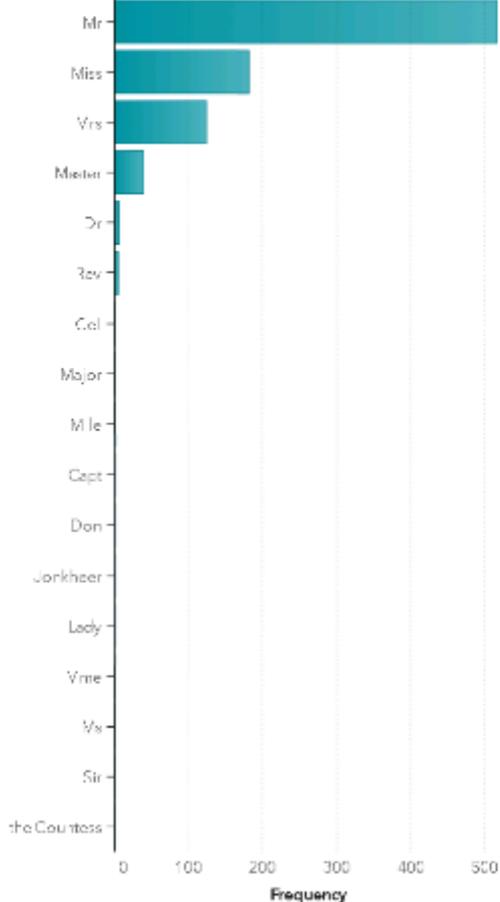
Ranks

Name	Ticket	Cabin
Braund, Mr. Owen Harris	A/5 21171	
Cumings, Mrs. John Bradley (Florence Briggs Thayer)	PC 17599	C85
Heikkinen, Miss. Laina	STON/O2. 3101282	
Futrelle, Mrs. Jacques Heath (Lily May Peel)	113803	C123
Allen, Mr. William Henry	373450	
Moran, Mr. James	330877	
McCarthy, Mr. Timothy J	17463	E46
Palsson, Master. Gosta Leonard	349909	
Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	347742	
Nasser, Mrs. Nicholas (Adele Achem)	237736	
Sandstrom, Miss. Marguerite Rut	PP 9549	G6
Bonnell, Miss. Elizabeth	113783	C103
Saunderscock, Mr. William Henry	A/5. 2151	
Andersson, Mr. Anders Johan	347082	
Vestrom, Miss. Hulda Amanda Adolfina	350406	
Hewlett, Mrs. (Mary D Kingcome)	248706	
Rice, Master. Eugene	382652	

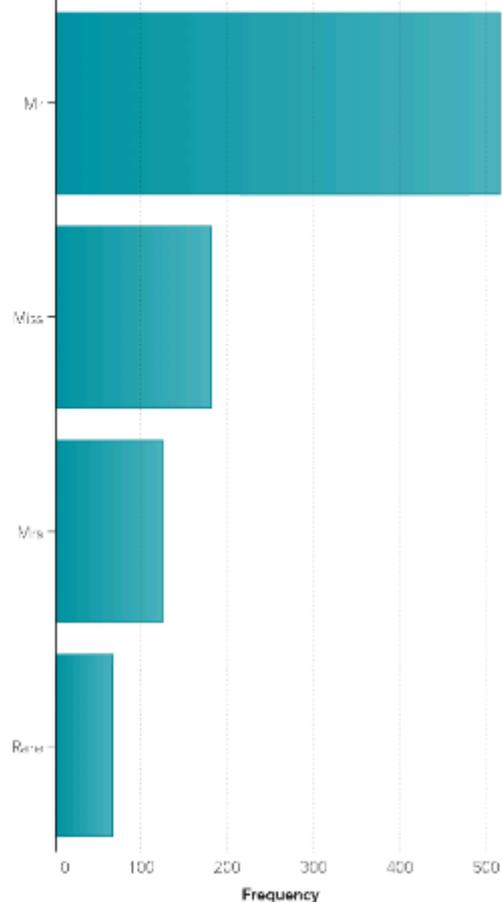
Name

Aikling, Wm. An. Henry
 Abbott, M.: Rosemore Edward
 Abbott, Mrs. Stanton (Rosa Hunt)
 Abelson, M.: Samuel
 Abelson, Mrs. Samuel (J. Hannah Weosky)
 Aasahl, M.: Mauritz Nils Martin
 Adams, Mr. John
 Ahlin, Mrs. ...shen (Johanna Persdotter Larsson)
 Ake, Mrs. Sam (Leah Rosen)
 Albinson, Mr. Nassaf Cassem
 Alexander, Wm. William
 Alho maki, Mr. Ilmari Rudolf
 Ali, Mr. Ahmad
 Ali, Mr. William
 Allen, Miss. Elisabeth Walton
 Allen, Wm. William Henry
 Allison, Master. Hudson Trevor
 Alison, Miss. Helen Loraine
 Alison, Mrs. Hudson J.C. (Bessie Waldo Daniels)
 Allum, Mr. Owen George
 Andersen-Jensen, Miss. Carla Christine Nielsine
 Anderson, Mr. Harry
 Anderson, Master. Sigvard Herald Elias
 Anderson, Miss. Ebba Iris Alfrida
 Anderson, Miss. Ellis Anna Maria
 Anderson, Miss. Erna Alexandra
 Anderson, Miss. Ingeborg Constanza
 Anderson, Miss. Sigrid Elisabeth
 Anderson, Mr. Anders Johan
 Anderson, Mr. August Edvard ("Wennerstrom")
 Anderson, Mrs. Anders Johan (Alfrida Konstantia Brogren)
 Anderson, Mr. Paul Edwin
 Andrew, Wm. Edgardo Samuel
 Andrews, Miss. Kornelia Theodosia
 Andrews, Mr. Thomas Jr.
 Angle, Mrs. William A. (Florence "Mary" Agnes Hughes)

Title



Title Group





Explore



Logistic Regression Survived

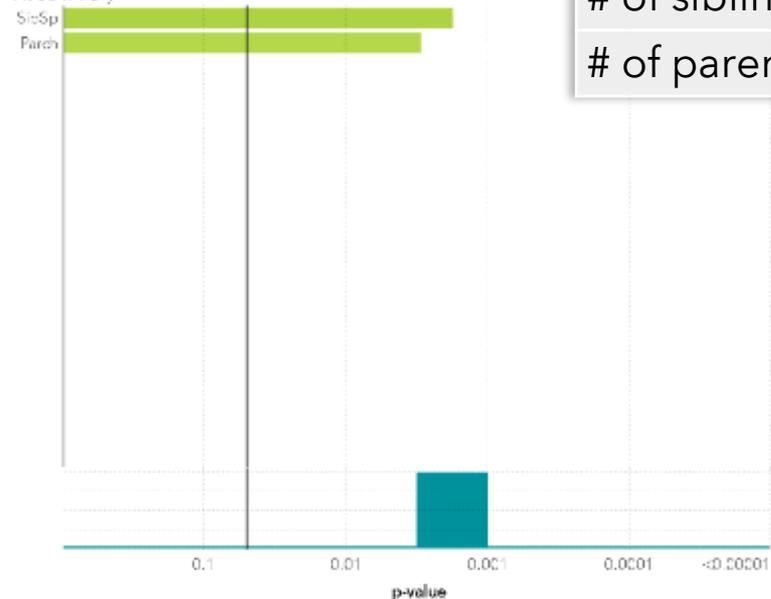
(even=1)

Validation Misclassification

0.2905

Observations Used: 891

Fit Summary

Better 

of siblings & spouses

of parents & children

Fit

Overview

Risks

Actions

Filter

Views

Flags

Ranks

Edit Calculated Item

Name:

Format: Comma



Data It...

Opera...

Visual

Text

Numeric



Search



▶ Character

▼ Numeric

- _comma
- _period
- Age
- Age (Num..
- Age_imo...
- Age_imo...
- Fare
- Leaf ID (A...

(Parch + SibSo)

Message (0)

OK

Cancel

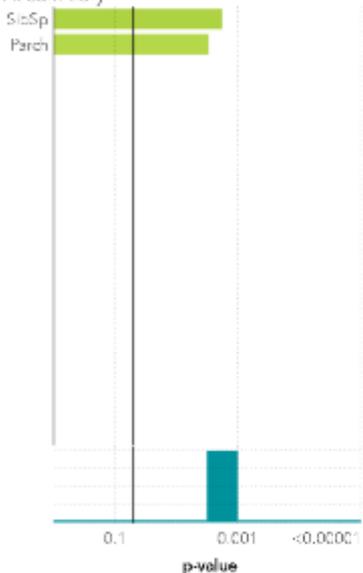
Explore

Logistic Regression **Survived**

(event=1)

Validation Misclassification

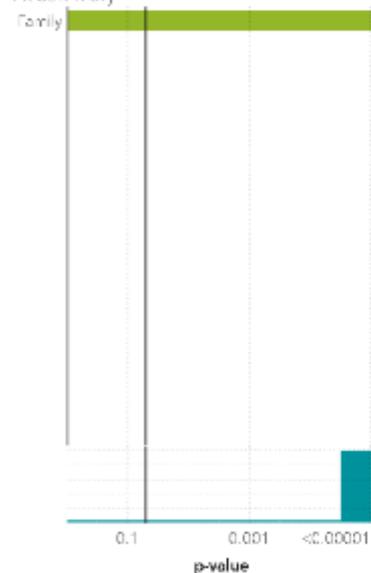
0.2905 Observations Used 891

Fit Summary **Better** Logistic Regression **Survived**

(event=1)

Validation Misclassification

0.2961 Observations Used 891

Fit Summary **Better** 

- Overview
- Roles
- Actions
- Filter
- Flags
- Rank



MODEL

Logistic Regression

Decision tree

Random Forest

Neural Network

Gradient Boosting

Support Vector Machines

Factorization Machines



Drop a data item or constant here

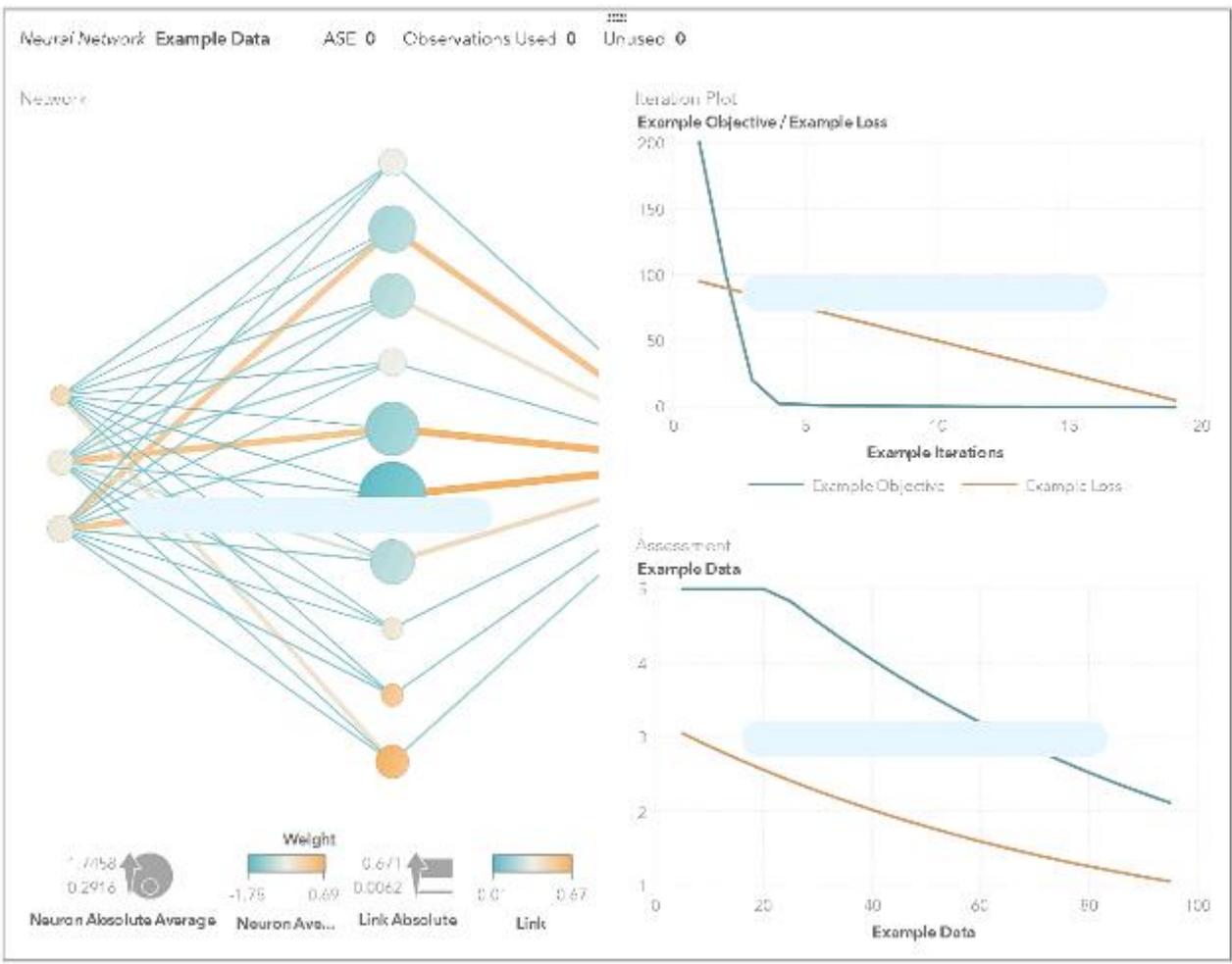
Page 1 Page 2 Name Age Logistic Regression ANN Decision Tree CBoost Random Forest

Objects

- Slider
- Text Input
- Analysis
 - Network Analysis
 - Text Topics
- Other
 - Container
 - Image
 - Form Container
 - Text
 - Web Content
- SAS Visual Statistics
 - Cluster
 - Decision Tree
 - Generalized Linear Model
 - Linear Regression
 - Logistic Regression
 - Model Comparison
- SAS Visual Data Mining and Machine Learning
 - Federated Machine Learning
 - Ensemble
 - Gradient Boosting
 - Neural Network
 - Support Vector Machine

Drop a data item or content here



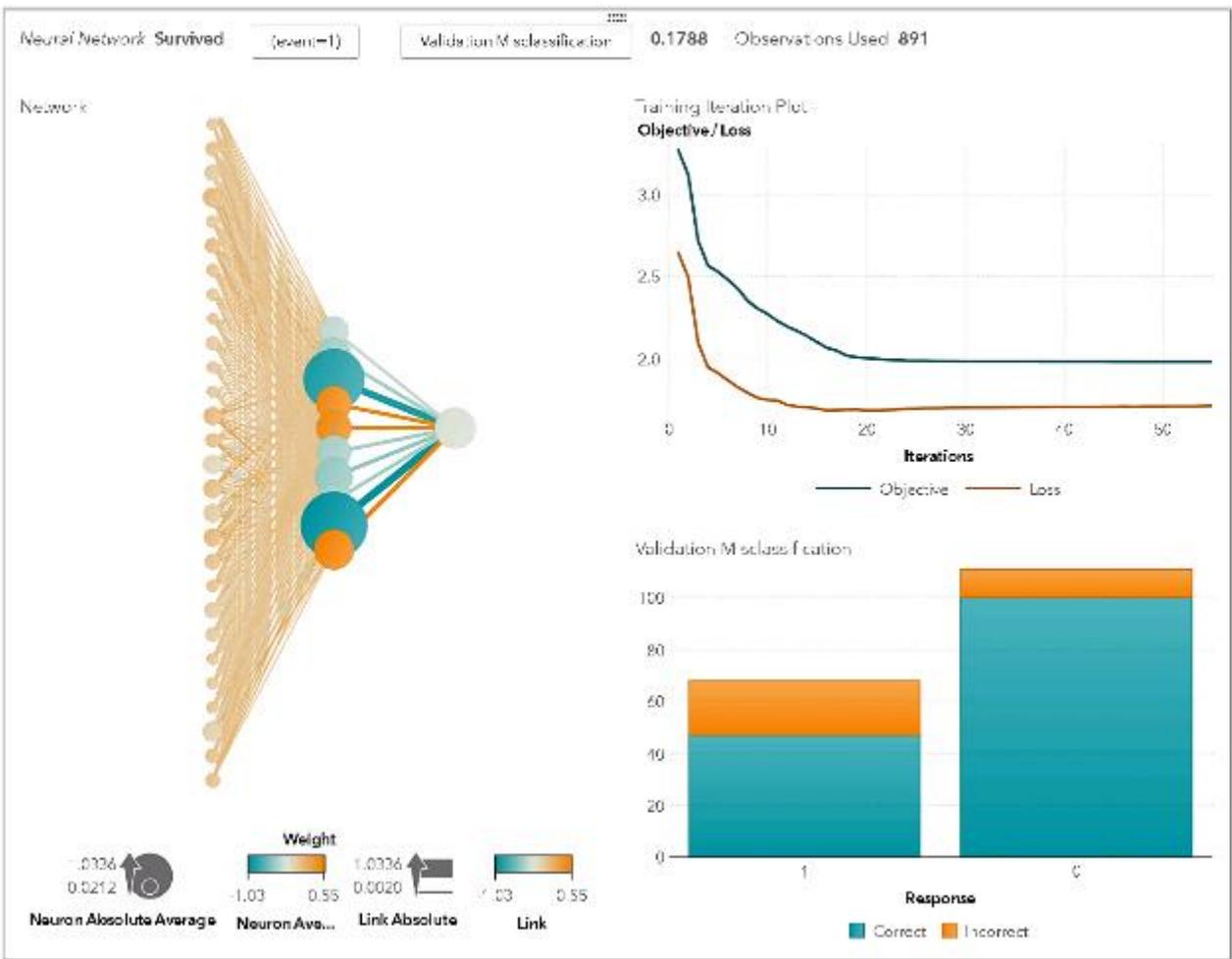


Roles

- Response*
 - + Add
- Predictors*
 - + Add
- Weight
 - + Add



Control Panel with icons for Roles, Actions, Filter, Flow, and Tables.



Roles

- ▼ Response
 - Survived
- ▼ Predictors
 - Leaf D(Age)
 - Sex
 - Title_new
 - Sex
 - Pclass
 - Parth
 - SibSp
 - Embarked_new
 - + Add
- ▼ Weight
 - + Add

Navigation sidebar with icons for Home, Roles, Actions, Filter, Flow, and Help.



WHICH Hyperparameters?

Decision Trees

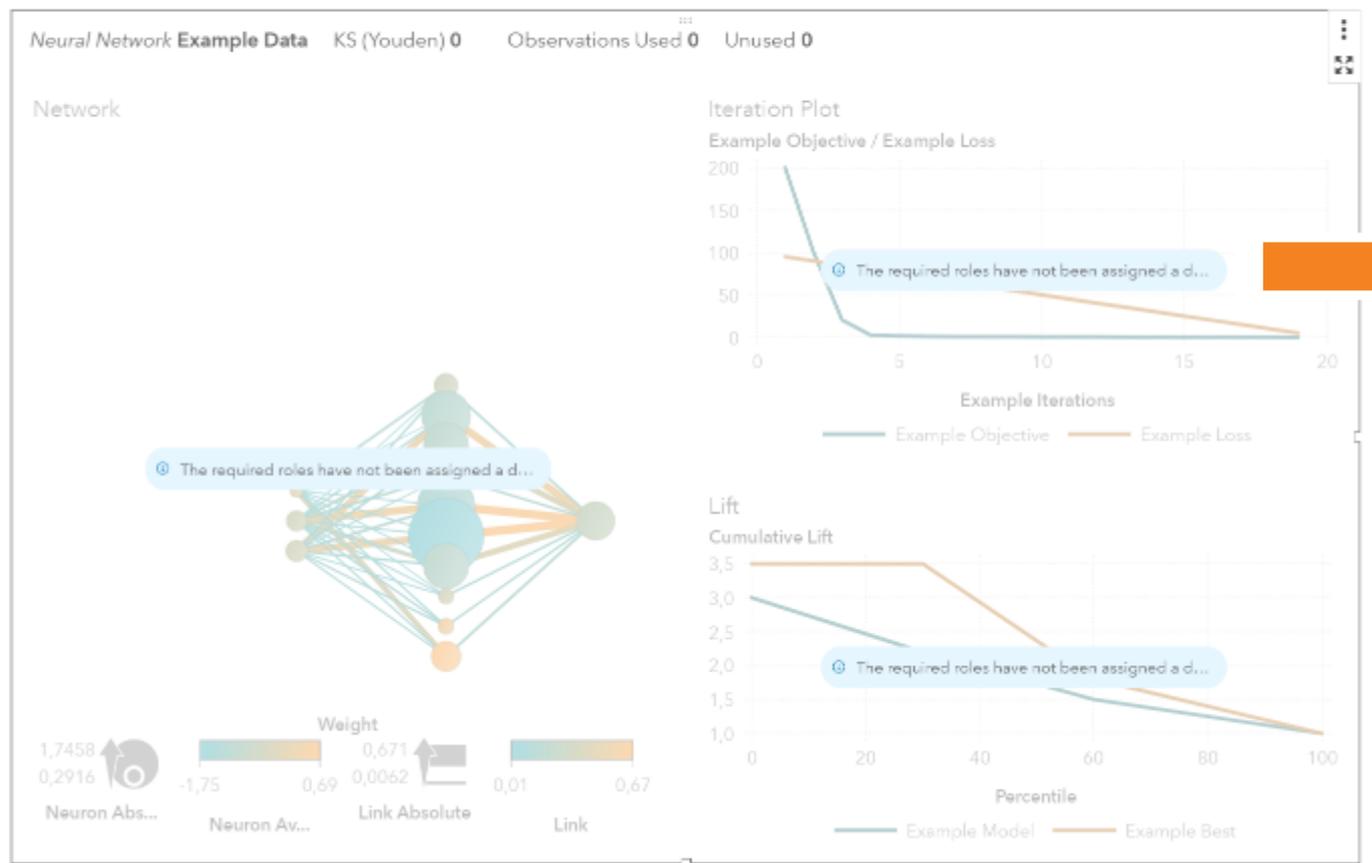
Max tree depth, splitting criterion

Neural Networks

Network architecture, solver options

Support Vector Machines

Kernel, penalty



Options

Neural Network 1

Neural Network

General

Autotune:

Include missing

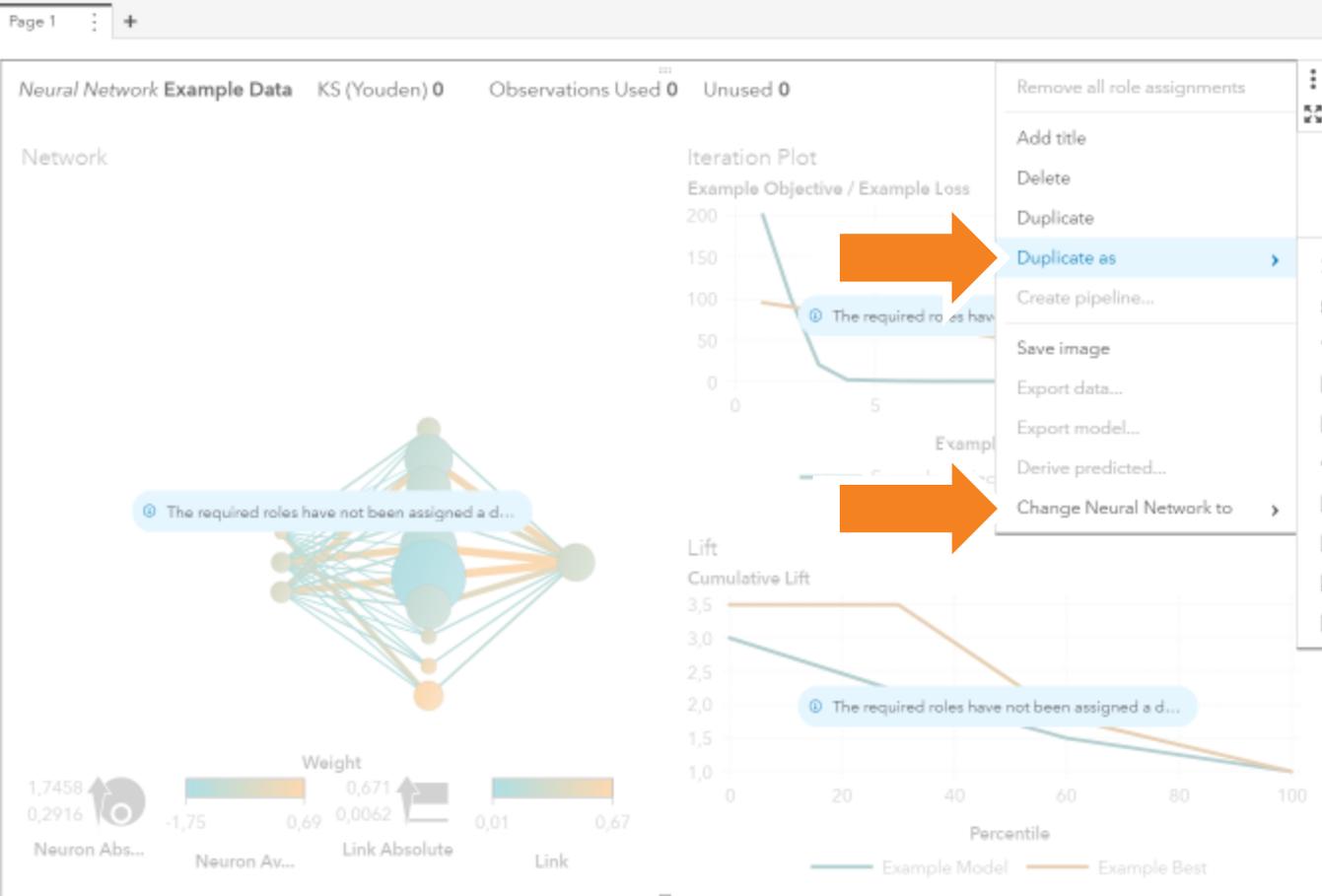
Standardization: Midrange
Maximum iterations: 250
Maximum time(sec): 270

Optimization method: LBFGS
L1: 0
L2: 0,1

Hidden Layers
Number of hidden layers: 1



Report 1



- Remove all role assignments
- Add title
- Delete
- Duplicate
- Duplicate as**
- Create pipeline...
- Save image
- Export data...
- Export model...
- Derive predicted...
- Change Neural Network to

Options

Neural Network 1

Neural Network

General

- Cluster
- Decision Tree
- Forest
- Generalized Additive Model
- Generalized Linear Model
- Gradient Boosting
- Linear Regression
- Logistic Regression
- Nonparametric Logistic Regression
- Support Vector Machine

LBFGS

L1: 0

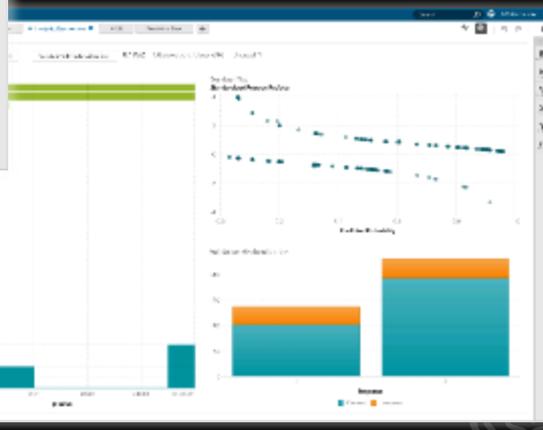
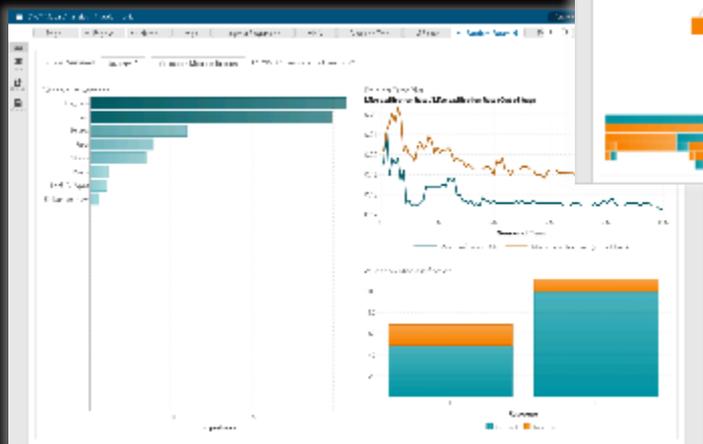
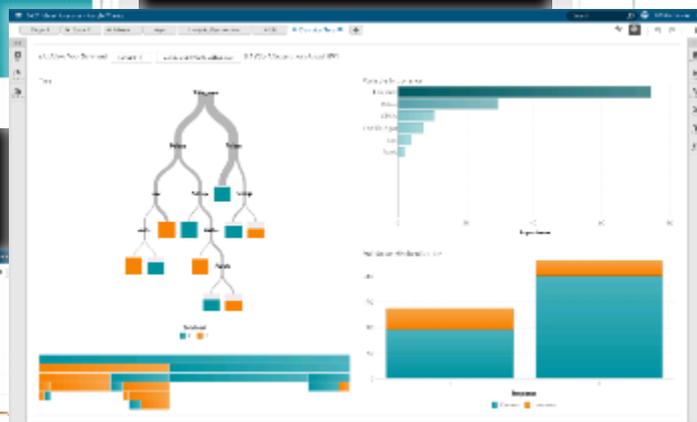
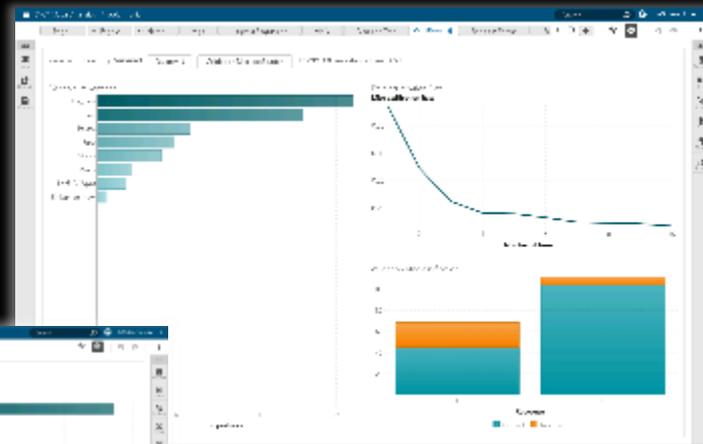
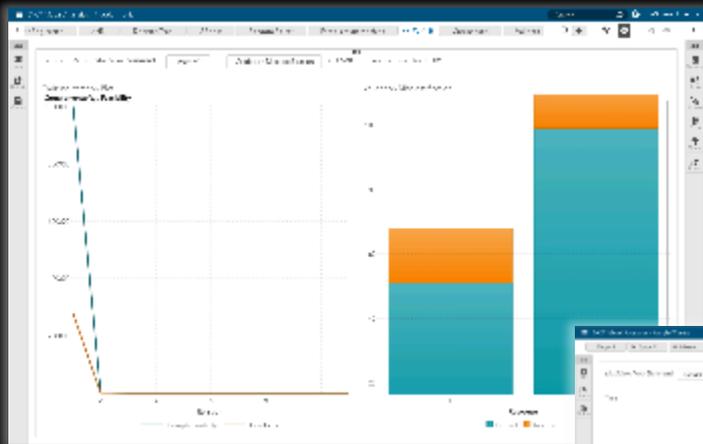
L2: 0,1

Hidden Layers

Number of hidden layers: 1

Validation Misclassification







ASSESS MODELS

Model comparison

[Overview](#)[Data](#)[Kernels](#)[Discussion](#)[Leaderboard](#)[More](#)[My Submissions](#)[Submit Predictions](#)**SUBMIT_NNEUR.csv**a month ago by [PetriROI](#)

Neural Network

0.75120

**SUBMIT_LREGR.csv**a month ago by [PetriROI](#)

Logistic Regression

0.79904

**SUBMIT_DTREE.csv**a month ago by [PetriROI](#)

Decision Tree

0.74641



**k**[Overview](#)[Data](#)[Kernels](#)[Discussion](#)[Leaderboard](#)[More](#)[My Submissions](#)[Submit Predictions](#)

1110	▼ 111	SergeyKuntsev		0.79904	48	23d
1111	▼ 111	Elizabeth M. Cruz		0.79904	25	20d
1112	▼ 111	Nilson de Lima Jr		0.79904	3	1mo
1113	▼ 111	PetriROI		0.79904	10	1mo
1114	▼ 111	StevenJin		0.79904	10	1mo
1115	▼ 111	lunge		0.79904	21	1mo
1116	▼ 111	SaloniKakkar		0.79904	31	1mo
1117	▼ 111	Naoki Narimatsu		0.79904	16	1mo
1118	▼ 111	srishtideepani		0.79904	3	1mo
1119	▼ 111	Dmitry Kovalev		0.79904	12	10d
1120	▼ 111	Md Sadik Hussain		0.79904	4	1mo
1121	▼ 111	Koffi		0.79904	32	24d

7000
6000
5000
4000
3000
2000
1000
0

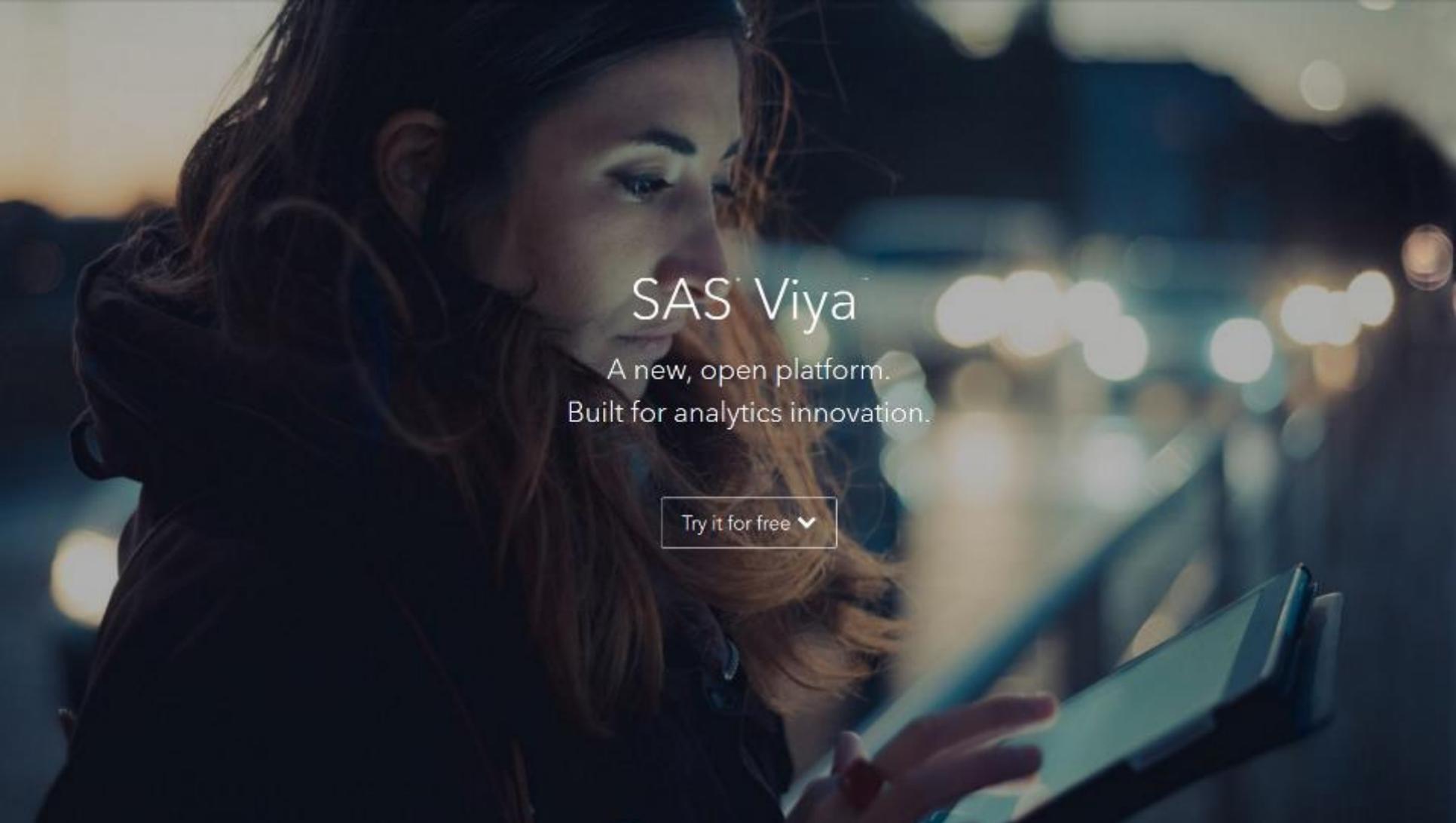


Titanic

600000
500000
400000
300000
200000
100000
0



Kaggle

A woman with long brown hair is looking at a tablet in her hands. The background is a blurred city street at night with warm lights.

SAS Viya

A new, open platform.
Built for analytics innovation.

Try it for free 