DATA ANALYSIS

Dimensionality Reduction Techniques



Curse of Dimensionality

... refers to the difficulties that arise when analyzing or modeling data with many dimensions.

Data points become increasingly spread out (Data Sparsity), making it hard to find patterns or relationships.





Three simple techniques to reduce dimensionality

• Missing values ratio

If the percentage of missing values in a variable exceeds the threshold (a pre-specified value), you can drop the variable.

• Low variance filter

All the data columns with variance lower than the threshold value will be eliminated.

• High correlation filter

All the pairs of columns having a correlation coefficient higher than the set threshold will be reduced to 1.



Principal Component Analysis (PCA)

It is a statistical method used to simplify and summarize data by reducing its dimensionality while retaining as much of the original variability as possible.

- 1. Apply PCA to the predictors using the training data. Use the output to determine the number of principal components to be retained. The predictors in the model now use the (reduced number of) principal scores columns. For the validation set, we can use the weights computed from the training data to obtain a set of principal scores by applying the weights to the variables in the validation set. These new variables are then treated as the predictors.
- 2. One **disadvantage** of using a subset of principal components as predictors in a supervised task, is that we might lose predictive information that is nonlinear (e.g., a quadratic effect of a predictor on the outcome variable or an interaction between predictors). This is because PCA produces linear transformations, thereby capturing linear relationships between the original variables.



Cereals

Variable	Description							
mfr	Manufacturer of cereal (American Home Food Products, General Mills, Kellogg, etc.)							
type	Cold or hot							
calories	Calories per serving							
protein	rams of protein							
fat	Frams of fat							
sodium	Milligrams of sodium							
fiber	Grams of dietary fiber							
carbo	Grams of complex carbohydrates							
sugars	Grams of sugars							
potass	Milligrams of potassium							
vitamins	Vitamins and minerals: 0, 25, or 100, indicating the typical percentage of FDA recommended							
shelf	Display shelf (1, 2, or 3, counting from the floor)							
weight	Weight in ounces of one serving							
cups	Number of cups in one serving							
rating	Rating of the cereal calculated by <i>consumer reports</i>							



Cereals

	name	mfr	type	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf	weight	cups	rating
0	100%_Bran	Ν	С	70	4	1	130	10.0	5.0	6.0	280.0	25	3	1.0	0.33	68.402973
1	100%_Natural_Bran	Q	С	120	3	5	15	2.0	8.0	8.0	135.0	0	3	1.0	1.00	33.983679
2	All-Bran	К	С	70	4	1	260	9.0	7.0	5.0	320.0	25	3	1.0	0.33	59.425505
3	All-Bran_with_Extra_Fiber	К	С	50	4	0	140	14.0	8.0	0.0	330.0	25	3	1.0	0.50	93.704912
4	Almond_Delight	R	С	110	2	2	200	1.0	14.0	8.0	NaN	25	3	1.0	0.75	34.384843
4	All-Bran_with_Extra_Fiber	R	С	110	2	2	200	14.0	14.0	8.0	NaN	25	3	1.0	0.75	34.38484





PATRÍCIA XUFRE

SKI FARN

PCA()

Cereals

pcsComponents_df.iloc[:,:5]

	PC1	PC2	PC3	PC4	PC5
calories	0.036594	-0.571142	-0.077927	-0.164032	0.186667
protein	0.330888	0.045367	0.326405	-0.281265	0.372134
fat	0.212549	-0.249143	-0.291274	-0.124324	0.597422
sodium	-0.050894	-0.304684	0.343108	-0.168916	-0.292344
fiber	0.458623	0.201029	0.173912	-0.075295	-0.229653
carbo	-0.268609	-0.175660	0.513002	-0.071305	0.191830
sugars	0.081440	-0.364844	-0.495338	0.016033	-0.360371
potass	0.502603	0.070345	0.130460	-0.112708	-0.124888
vitamins	0.018839	-0.300574	0.327563	0.594798	-0.105582
shelf	0.308929	-0.065114	0.000546	0.667703	0.265485
weight	0.257160	-0.448962	0.122030	-0.160396	-0.232905
cups	-0.371154	-0.106553	0.082695	-0.030683	0.115826

	PCA Loadings												
<u>5</u>	0.037	0.33	0.21	-0.051		-0.27	0.081		0.019	0.31	0.26	-0.37	- 0.6
PC2	-0.57	0.045	-0.25	-0.3	0.2	-0.18	-0.36	0.07	-0.3	-0.065		-0.11	
Ŋ.	- 0.078	0.33	-0.29	0.34	0.17			0.13	0.33	0.00055	0.12	0.083	- 0.4
PC4	-0.16	-0.28	-0.12	-0.17	-0.075	-0.071	0.016	-0.11	0.59	0.67	-0.16	-0.031	
PC5	0.19	0.37	0.6	-0.29	-0.23	0.19	-0.36	-0.12	-0.11	0.27	-0.23	0.12	- 0.2
PC6	0.2	0.18	0.3		0.085		0.061	0.068	0.28	-0.13	-0.32		- 0.0
PC7	0.062	0.26	-0.25	-0.61	0.096	-0.027	0.19	0.15	0.17	-0.12	0.19		
BC8	0.0049		0.054	0.17	0.22	0.17	-0.068	0.27	-0.46		0.077		0.
- PC	0.043			0.25	-0.26	-0.14	0.25	-0.19	-0.34		-0.068	0.045	
C10	0.5	-0.02	-0.15	-0.0055	0.29	0.24	0.25	0.19	0.052	-0.047	-0.69	-0.078	0.
CI1	0.19	0.015	-0.0039	-0.0082	0.64	-0.15	-0.13	-0.71	-0.036	0.051	0.11	0.057	
C12	-0.53	0.14	0.26	-0.0088	0.19	0.51	0.55	-0.18	-0.016	0.018	-0.025	-0.029	0.
_	calories	protein	fat	sodium	fiber	carbo	sugars	potass	vitamins	shelf	weight	cups	_



It is difficult to give an interpretation to the Principal Components



77

Cereals





When the data can be reduced to two dimensions, a useful plot is a scatter plot of the first vs. second principal scores with labels for the observations



Performance Evaluation



79

PERFORMANCE EVALUATION

Predictive Power Assessment



Predictive Performance Evaluation



Predicted numerical value

when the outcome variable is numerical

Predicted class membership

when the outcome variable is categorical

Propensity - the probability of class membership when the outcome variable is categorical



Prediction Accuracy Measures

The **prediction error** is defined as the difference between its actual outcome value and its predicted outcome value: $e_i = y_i - \hat{y}_i$.

• Mean Absolute Error (MAE) = $\frac{1}{n} \sum |e_i|$

the magnitude of the average absolute error.

• Mean Absolute Percentage Error (MAPE) = $100 \times \frac{1}{n} \sum \left| \frac{e_i}{v_i} \right|$

a percentage score of how predictions deviate (on average) from the actual values.

• Root Mean Squared Error (RMSE) = $\sqrt{\frac{1}{n}\sum e_i^2}$

is similar to the standard error.





Toyota Car Prices

car_df = pd.read_csv('ToyotaCorolla.csv')

create a list of predictor variables by removing output variables and text columns excludeColumns = ('Price', 'Id', 'Model', 'Fuel_Type', 'Color') predictors = [s for s in car_df.columns if s not in excludeColumns] # the explanatory variables outcome = 'Price' # the dependent variable

partition data

NOVA SCHOOL OF BUSINESS & ECONOMICS

X = car_df[predictors]
y = car_df[outcome]
train_X, valid_X, train_y, valid_y = train_test_split(X, y, test_size=0.4,random_state=1)
train Linear regression modeL
reg = LinearRegression()
reg.fit(train_X, train_y)

	ld	Model	Price	Age_08_04	Mfg_Month	Mfg_Year	KM	Fuel_Type	ΗP	Met_Color	 Powered_Windows	Power_Steering	Radio	Mistlamps	Sp
0	1	TOYOTA Corolla 2.0 D4D HATCHB TERRA 2/3- Doors	13500	23	10	2002	46986	Diesel	90	1	 1	1	0	0	
1	2	TOYOTA Corolla 2.0 D4D HATCHB TERRA 2/3- Doors	13750	23	10	2002	72937	Diesel	90	1	 0	1	0	0	
2	3	TOYOTA Corolla 2.0 D4D HATCHB TERRA 2/3- Doors	13950	24	9	2002	41711	Diesel	90	1	 0	1	0	0	
3	4	TOYOTA Corolla 2.0 D4D HATCHB TERRA 2/3- Doors	14950	26	7	2002	48000	Diesel	90	0	 0	1	0	0	
4	5	TOYOTA Corolla 2.0 D4D HATCHB SOL 2/3- Doors	13750	30	3	2002	38500	Diesel	90	0	 1	1	0	1	



Toyota Car Prices

evaluate performance

training
mae = metrics.mean_absolute_error(train_y, reg.predict(train_X))
mape =metrics.mean_absolute_percentage_error(train_y, reg.predict(train_X))
mse = metrics.mean_squared_error(train_y, reg.predict(train_X))
rmse = np.sqrt(mse)
print('Performance Measures - Training set')
print(f'Mean absolute error: {mae:.2f}')
print(f'Mean absolute percentage error: {mape:.2f}')
print(f'Root mean squared error: {rmse:.2f}')

validation

mae = metrics.mean_absolute_error(valid_y, reg.predict(valid_X))
mape =metrics.mean_absolute_percentage_error(valid_y, reg.predict(valid_X))
mse = metrics.mean_squared_error(valid_y, reg.predict(valid_X))
rmse = np.sqrt(mse)
print('\nPerformance Measures - Validation set')
print(f'Mean absolute error: {mae:.2f}')
print(f'Mean absolute percentage error: {mape:.2f}')
print(f'Root mean squared error: {rmse:.2f}')



Performance Measures - Training set Mean absolute error: 811.68 Mean absolute percentage error: 0.08 Root mean squared error: 1121.06

Performance Measures - Validation set Mean absolute error: 880.14 Mean absolute percentage error: 0.09 Root mean squared error: 1382.04



Toyota Car Prices



pred_error_train = pd.DataFrame({ 'residual': train_y - reg.predict(train_X), 'data set': 'training' }) pred_error_valid = pd.DataFrame({ 'residual': valid_y - reg.predict(valid_X), 'data set': 'validation' }) boxdata_df = pd.concat([pred_error_train, pred_error_valid], ignore_index=True) fig, axes = plt.subplots(nrows=1, ncols=3) fig.set_size_inches(9, 4) common = {'bins': 100, 'range': [-6500, 6500]} pred_error_train.hist(ax=axes[0], **common) pred_error_valid.hist(ax=axes[1], **common) boxdata_df.boxplot(ax=axes[2], by='data set')

axes[0].set_title('training')
axes[1].set_title('validation')
axes[2].set_title(' ')
axes[2].set_ylim(-6500, 6500)
plt.suptitle('Prediction errors')
plt.subplots_adjust(bottom=0.1, top=0.85, wspace=0.35)
plt.show()



PERFORMANCE EVALUATION

Judging Classifier Performance



Confusion Matrix

A **confusion matrix** is a matrix that summarizes the performance of a machine learning model on a set of test data. It is a means of displaying the number of accurate and inaccurate instances based on the model's predictions. It is often used to measure the performance of classification models, which aim to predict a categorical label for each input instance.

		Actual Value					
		POSITIVE	NEGATIVE				
Predicted	POSITIVE	TP	FP				
Value	NEGATIVE	FN	TN				



Accuracy Measures



• Accuracy = $\frac{TP+TN}{TP+FT+FN+TN}$ fraction of predictions our model got right.

• **Precision** =
$$\frac{TP}{TP+FP}$$

how many of the positive class samples present in the dataset were correctly identified by the model.

• Recall =
$$\frac{TP}{TP+FN}$$

how many of the "positive" predictions made by the model were correct.



Accuracy Measures

• $F1 - score = \frac{2 \times Precision \times Recall}{Precision + Recall}$

The F1 score combines precision and recall using their harmonic mean, and maximizing the F1 score implies simultaneously maximizing both precision and recall.



Performance in Case of Unequal Importance of Classes

The **sensitivity** of a classifier is its ability to detect the important class members correctly. This is the recall and is the percentage of C1 members classified correctly.

The **specificity** of a classifier is its ability to rule out C2 members correctly. This is measured by the recall of C2 and is the percentage of C2 members classified correctly.

		Actual	Value
		CLASS 1	CLASS 2
Predicted	CLASS 1	<i>n</i> _{1,1}	n _{2,1}
value	CLASS 2	n _{1,2}	n _{2,2}



ROC Curve

A ROC (Receiver Operating Characteristic) curve is a plot of the true positive rate (Sensitivity) in function of the false positive rate (100 - Specificity) for different cut-off points of a parameter. Each point on the ROC curve represents a sensitivity/specificity pair corresponding to a particular decision threshold. The Area Under the ROC curve (AUC) is a measure of how well a parameter can distinguish between two classes.



