



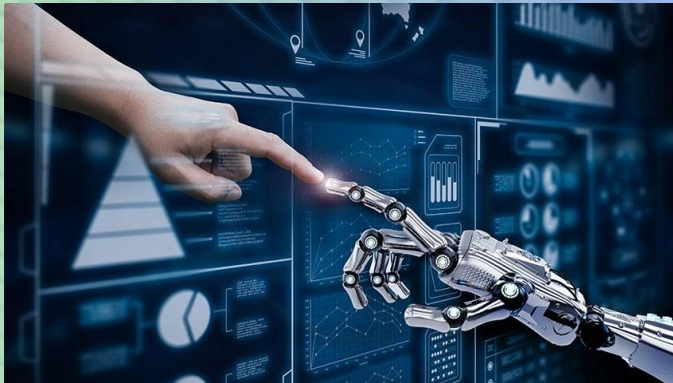
BATCH : **B 150** Data Science
LESSON : **Machine Learning**
DATE : 17.10.2022
SUBJECT : **UNSupervised Learning**



 techproeducation.com

 info@techproeducation.com

 +1 (917) 768-7466



MACHINE LEARNING - 6



Makine Öğrenmesi – 6
Unsupervised Learning



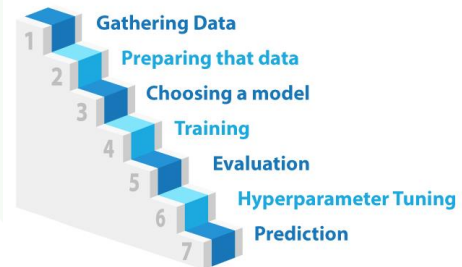
Overall Table of Contents



General Content

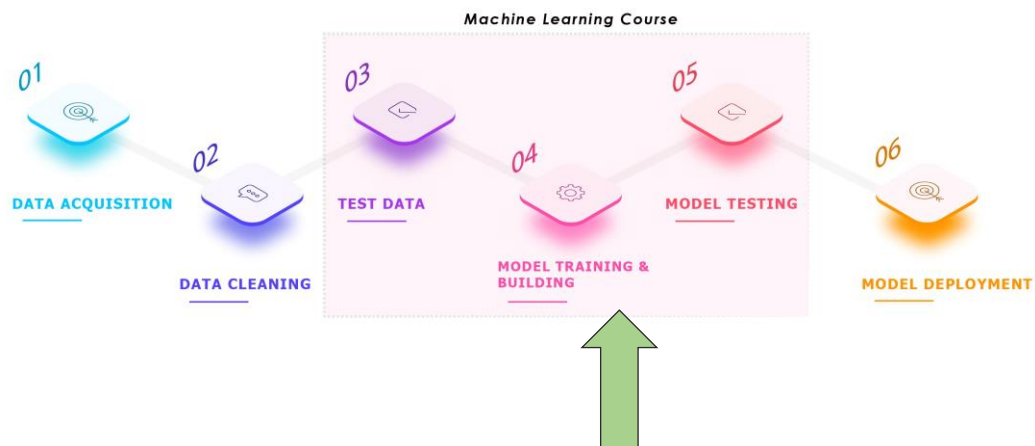
- ✓ UNsupervised Learning–
K-Means - PCA
- ✓ UNsupervised Algorithm
practices Python
application
- ✓ Projects Solutions

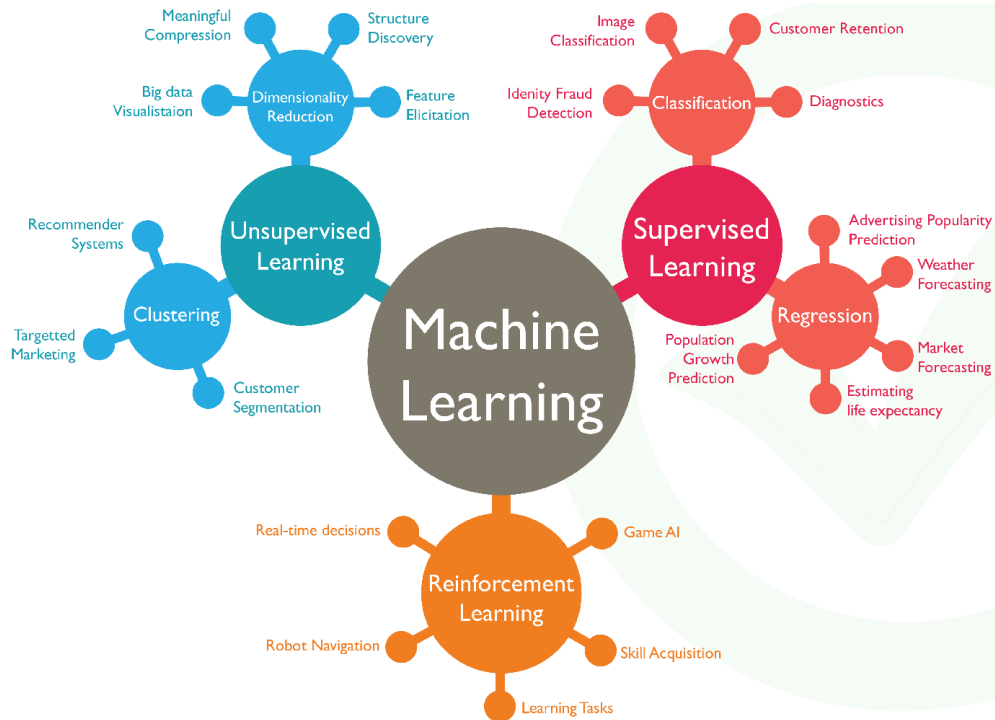
7 steps of Machine Learning



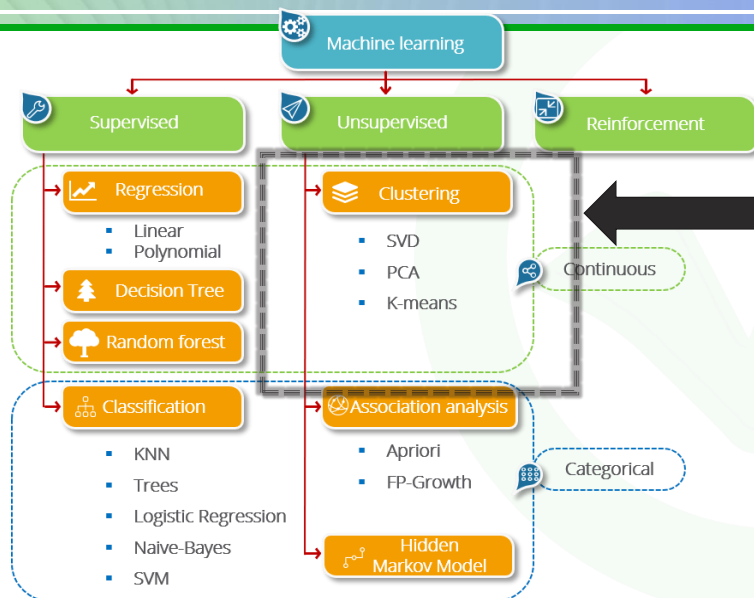
Run Navigation...

Where are we?





Unsupervised Learning





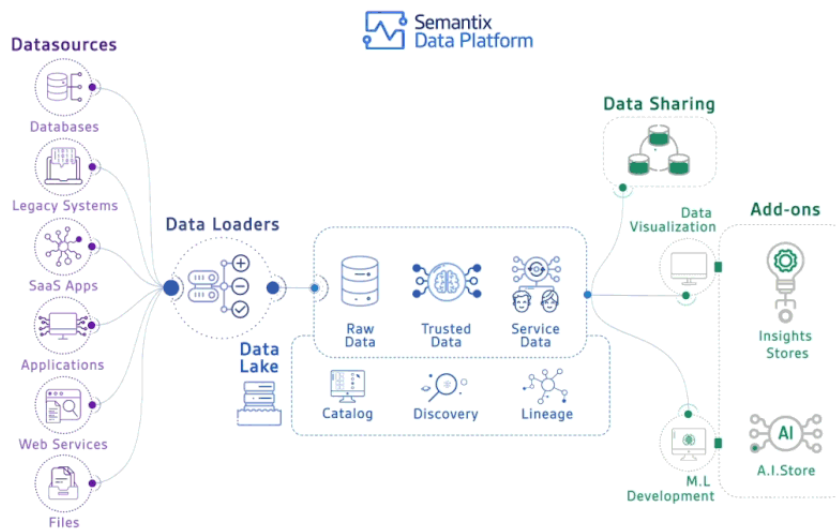
Unsupervised Learning

Overall Table of Contents

- ✓ Unsupervised Learning Algorithm (Denetimsiz Öğrenme)
- ✓ Unsupervised Algorithm practices Python application
- ✓ Projects Solutions



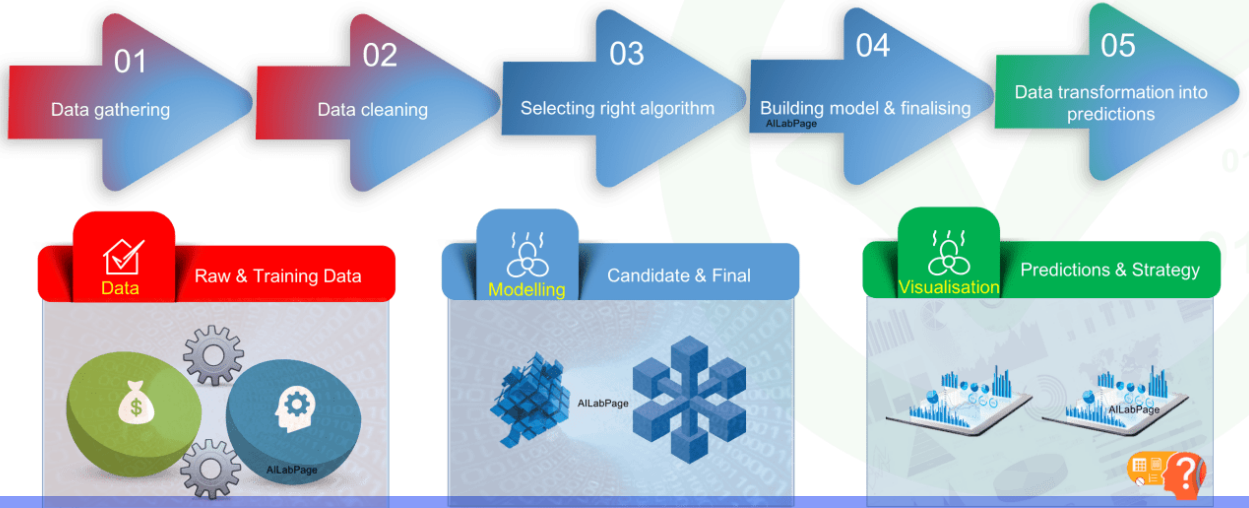
Unsupervised Learning





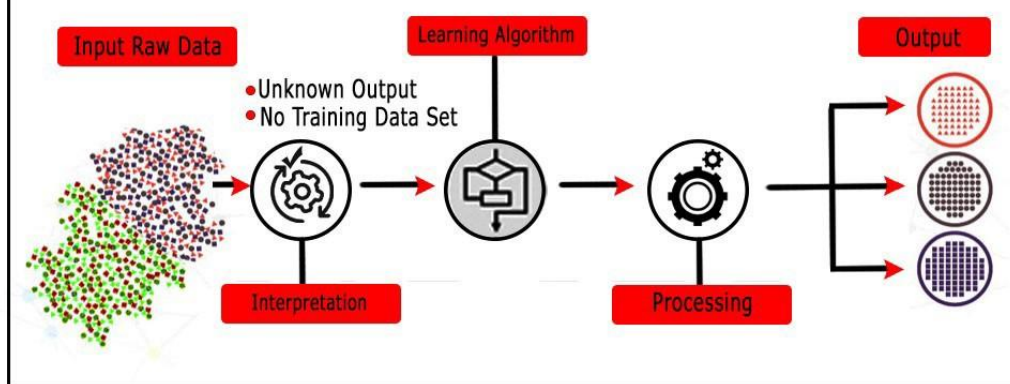
Unsupervised Learning

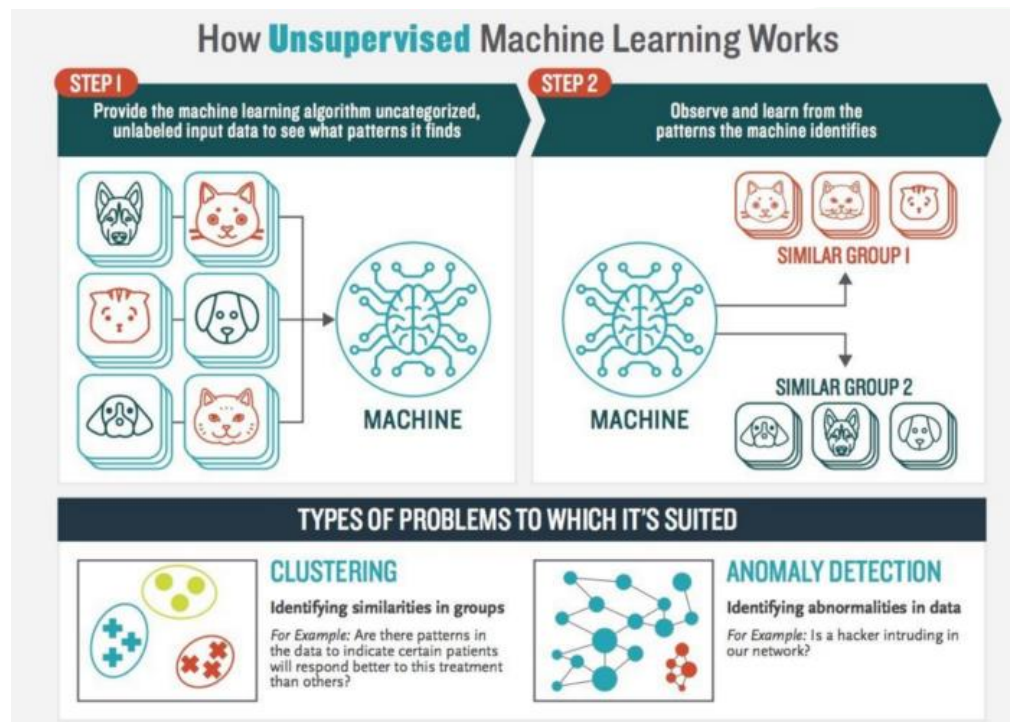
Machine Learning Process



Unsupervised Learning (Denetimsiz Öğrenme)

Unsupervised Learning





➤ Exploratory Data Analysis and Visualization

➤ Machine Learning

❖ Train | Test Split

- `X_train, X_test, y_train, y_test = train_test_split()`

❖ Scalling (if needed)

- `scaler = scaler_name()`
- `scaler.fit_transform(X_train)`
- `scaler.transform(X_test)`

❖ Modelling

- `model = model_name().fit(X_train, y_train)`
- `y_pred = model.predict(X_test)`
- `y_pred_proba = model.predict_proba(X_test)`

❖ Model Performance

- Regression => `r2_score`, `MAE`, `MSE`, `RMSE`
- Classification => `accuracy`, `recall`, `precision`, `f1_score` (`confusion_matrix`, `classification_report`)
- Cross Validate => `cross_val_score`, `cross_validate`

❖ Tunning (if needed)

- `grid_param = {}`
- `GridsearchCV(grid_param)`

❖ Final Model

- `model = model_name().fit(X, y)`

➤ Model Deployment

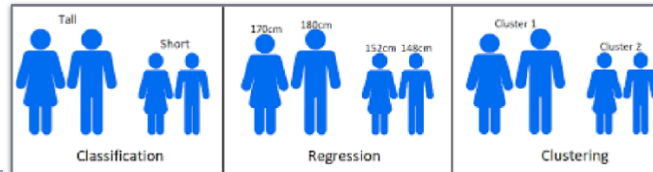
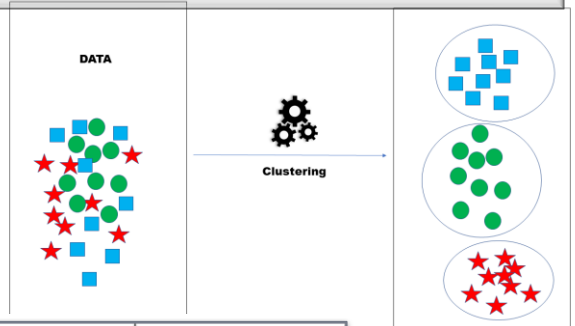
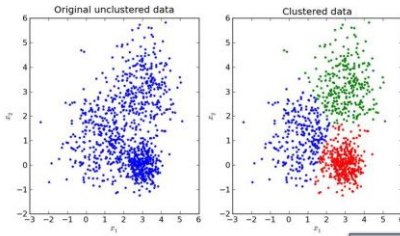


Unsupervised Learning

CLUSTERING' E GİRİŞ (Kümeleme)

✓ Meaningfulness - Usefulness

Unsupervised Learning

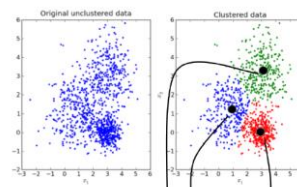
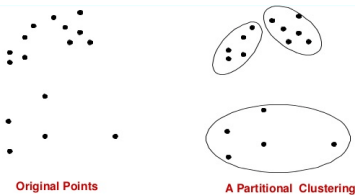


Unsupervised Learning

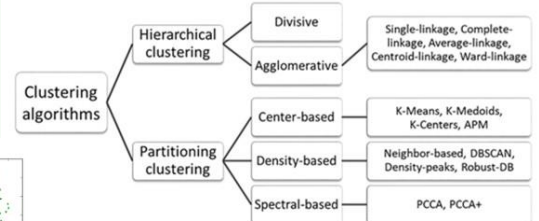
CLUSTERING' E GİRİŞ (Kümeleme)

✓ Clustering tipleri

1. Partitional clustering (K-means)



Centroids

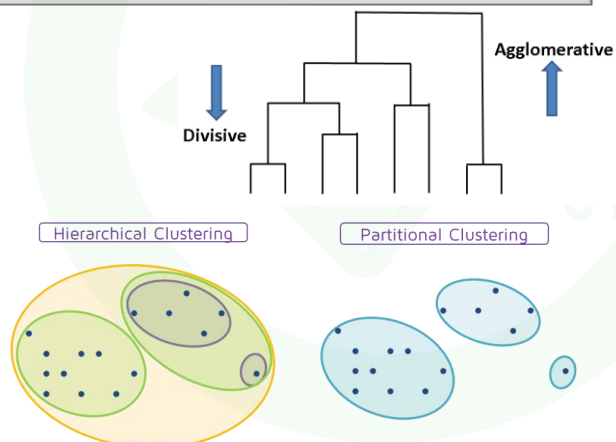
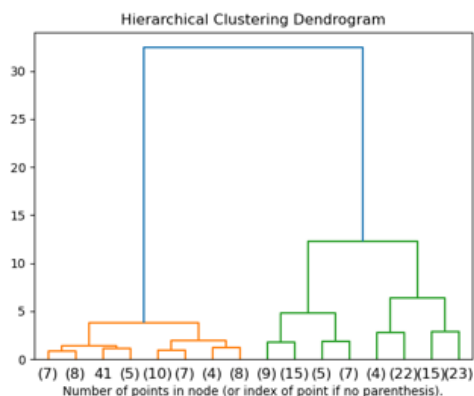




Unsupervised Learning

CLUSTERING' E GİRİŞ (Kümeleme)

- ✓ Clustering tipleri
- ✓ 2. Hierarchical clustering (PCA)



Unsupervised Learning

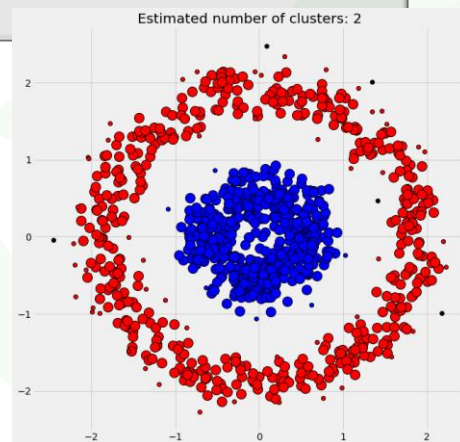
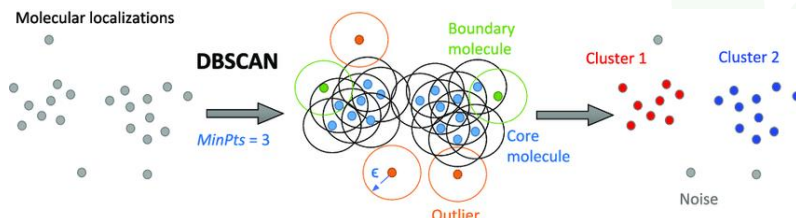
CLUSTERING' E GİRİŞ (Kümeleme)

- ✓ Clustering tipleri
- ✓ 3. Density-based clustering

Molecular localizations

DBSCAN

MinPts = 3

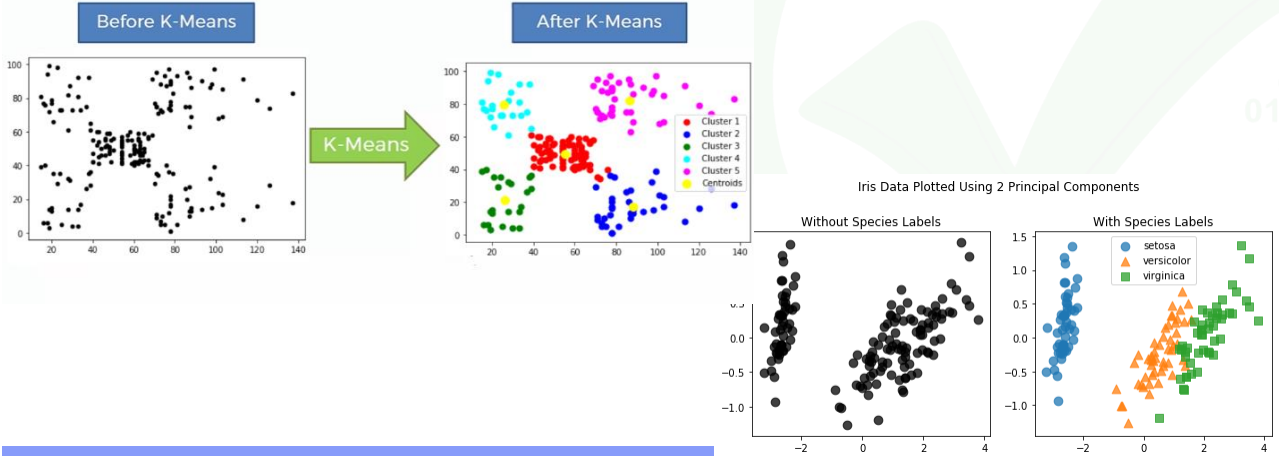




Unsupervised Learning

K-MEANS (K-Ortalamalar) CLUSTERING

✓ Hyperparameters

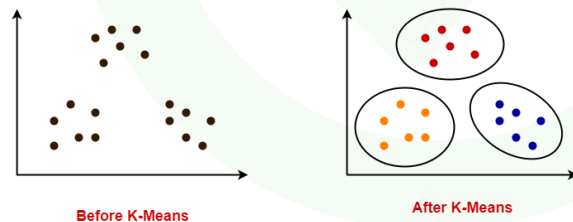
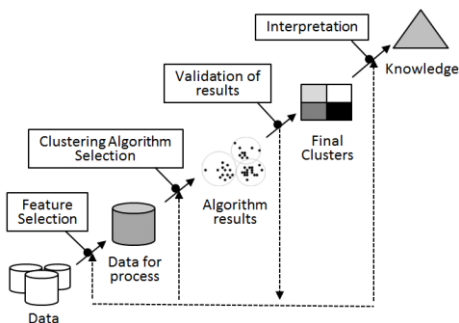


Unsupervised Learning

K-MEANS (K-Ortalamalar) CLUSTERING

✓ Cluster Modelleri için PERFORMANS DEĞERLENDİRME METRİKLERİ

- ✓ Cluster'lar ne kadar anlamlı
- ✓ Optimum sayıda cluster var mı
- ✓ Cluster ları doğrulamaya ihtiyaç var mı





Unsupervised Learning

K-MEANS (K-Ortalamlar) CLUSTERING

✓ Performans Ölçütleri

- ✓ Dışsal Kümeleme Doğrulaması
- ✓ True-Matching (TM)
- ✓ Adjusted Rand Index (ARI)
- ✓ Mutual Information Score (MIS)
- ✓ V-measure:
- ✓ Eksiksizlik (Completeness)
- ✓ İçsel Kümeleme Doğrulaması



Unsupervised Learning

K-MEANS (K-Ortalamlar) CLUSTERING

✓ Performans Ölçütleri

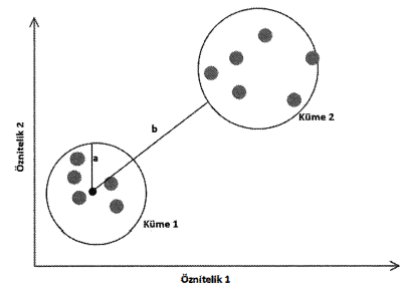
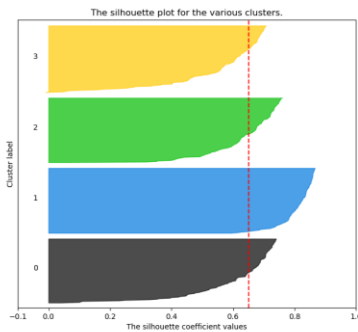
✓ Silhouette Metodu

$$\text{Silhouette Katsayısı} = (b - a) / \max(a, b)$$

a= küme içi ortalama mesafe (bir kümenin içindeki tüm veri noktalarının birbirleri arasındaki ortalama mesafe)

b= kümeler arası ortalama mesafe (tüm kümelerin arasındaki mesafelerin ortalaması)

Silhouette analysis for KMeans clustering on sample data with n_clusters = 4



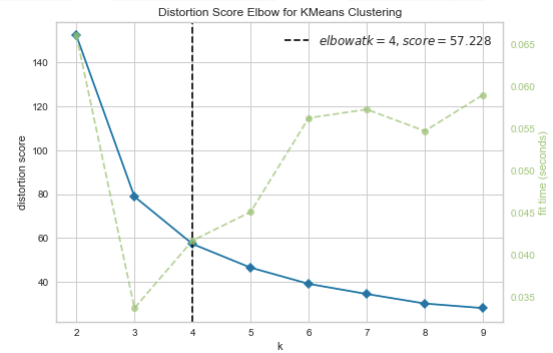
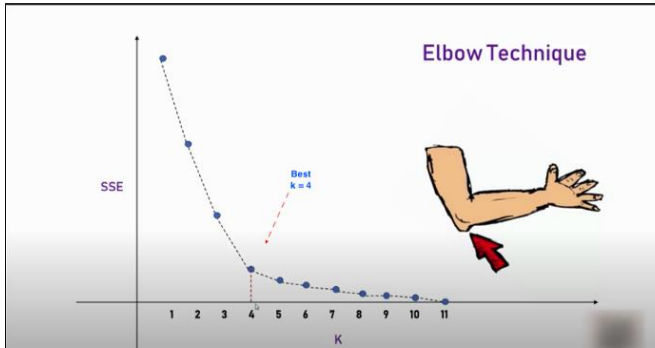


Unsupervised Learning

K-MEANS (K-Ortalamlar) CLUSTERING

✓ Performans Ölçütleri

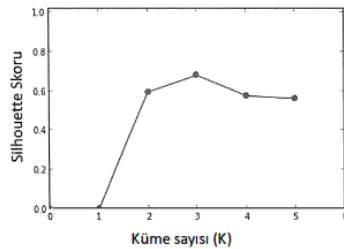
✓ Elbow Metodu



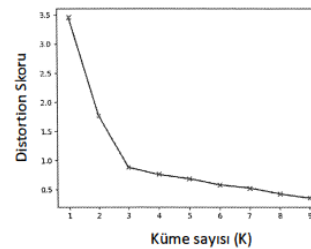
Unsupervised Learning

✓ K-MEANS (K-Ortalamlar) CLUSTERING

Silhouette ve Elbow Metotlarına Göre Optimum Küme Seçimi



Silhouette Metodunda, optimum küme sayısı (K) farklı küme sayıları için hesaplanan skorlardan en büyüğüne karşılık gelen değerdir. Bu örnekte K=3 optimum olarak görülmektedir.



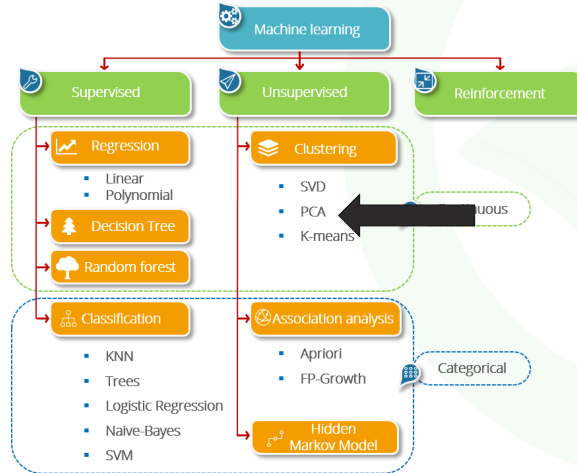
Dirsek Metodunda, optimum küme sayısı (K) grafiğin dirsek şeklinde karıldığı/büküldüğü noktaya karşılık gelen değerdir. Bu örnekte K=3 optimum olarak görülmektedir.

Son tahlilde, bir değerlendirme ölçütü olmakla birlikte denetimli makine öğrenmesinde kullanılan değerlendirme ölçütleri kadar kesinlik atfetmeden bir miktar ihtiyatla yaklaşmak gerekir.



Unsupervised Learning

PRINCIPAL COMPONENT ANALYSIS - PCA (Temel Bileşenler Analizi)



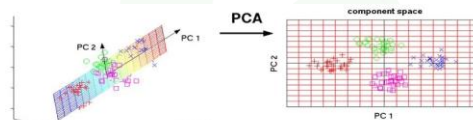
Unsupervised Learning

PRINCIPAL COMPONENT ANALYSIS –PCA

- ✓ Feature Selection
- ✓ Dimensionality Reduction



Set A

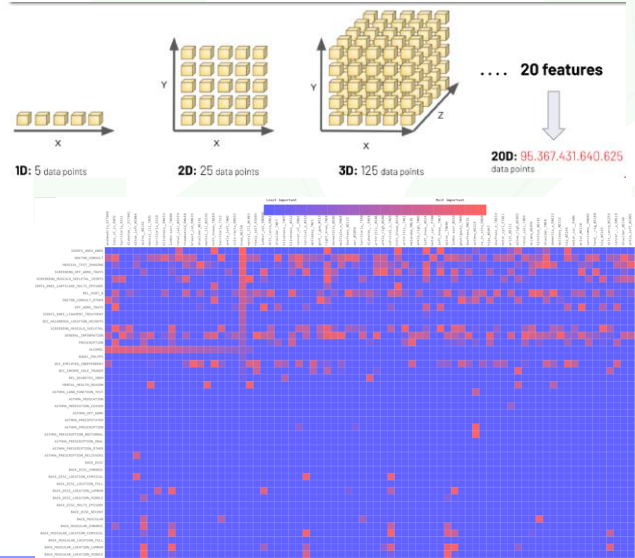
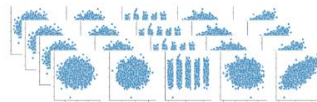
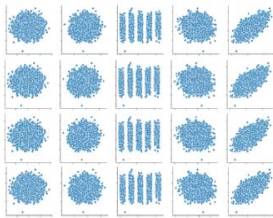
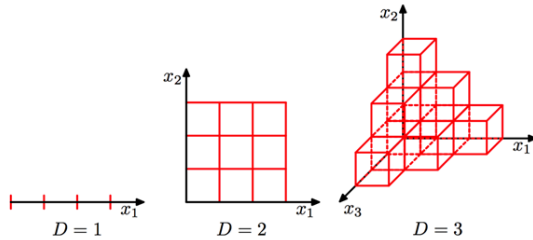


Principal Components Analysis (PCA)

PCA, değişkenlerdeki bilgilerin çoğunu temsil eden daha küçük bir bileşen kümesine indirmeye çalışan bir karmaşıklık azaltma tekniğidir.



PRINCIPAL COMPONENT ANALYSIS –PCA

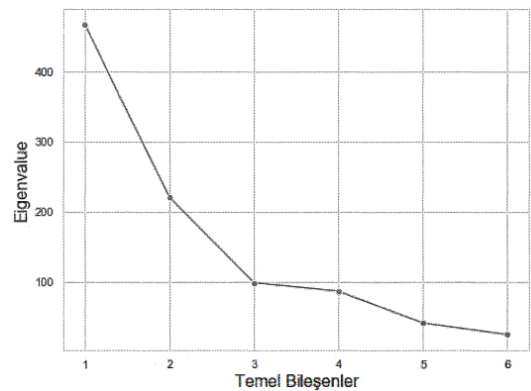


Unsupervised Learning

PRINCIPAL COMPONENT ANALYSIS –PCA

```
87 # PCA modeli için hiperparametrelere değerleri girelim:
88 import prince
89 pca = prince.PCA(n_components=6, n_iter=10, rescale_with_mean=False,
90                 rescale_with_std=False, copy=True, check_input=True,
91                 engine='sklearn', random_state=42)
92 # PCA modelini pca adı altında oluşturalım:
93 pca = pca.fit(df2[features])
94
```

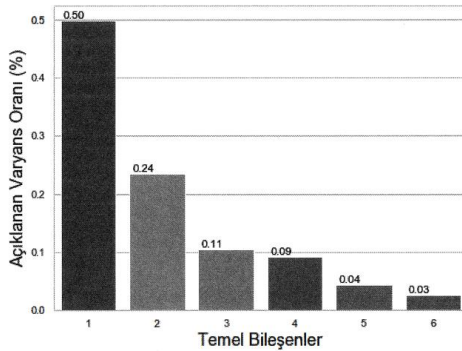
```
99 # Grafik: Elbow metoduyla optimum temel bileşen sayısı
100 ev = pd.DataFrame()
101 ev['pca'] = range(1,7)
102 ev['eigenvalue'] = pd.DataFrame(pca.eigenvalues_)
103 plt.figure(figsize=(8,6))
104 sns.lineplot(x='pca', y='eigenvalue', marker='o', data=ev)
105 plt.ylabel('Eigenvalue', fontsize=16)
106 plt.xlabel('Temel Bileşenler', fontsize=16)
107 plt.show()
108
```





Unsupervised Learning

PRINCIPAL COMPONENT ANALYSIS –PCA



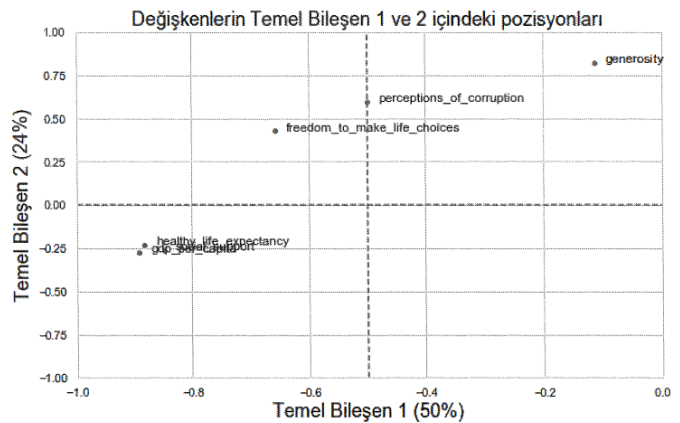
```
109 # Açıklanan Varyans Oranları:
110 # Hangi temel bileşen incelenen olayı ne kadar açıklıyor?
111 explained_variance=pca.explained_inertia_
112 print(explained_variance)
113
```

```
[0.4984116865967306, 0.2355096143382977,
0.10512215087762719, 0.09197939016712874,
0.0434179988625223, 0.025559159157693058]
```



Unsupervised Learning - PCA

```
130 # Grafik: Sütun Korelasyonu
131 pca.column_correlations(df2[features])
132
133 scatter = pd.DataFrame(pca.column_correlations(df2[features])).reset_index()
134 plt.figure(figsize=(10,6))
135 ax = sns.scatterplot(x=0, y=1, data=scatter)
136 ax.set(ylim=(-1, 1), xlim=(-1,0))
137 def label_point(x, y, val, ax):
138     a = pd.concat({'x': x, 'y': y, 'val': val}, axis=1)
139     for i, point in a.iterrows():
140         ax.text(point['x']+0.02, point['y'], str(point['val']))
141 label_point(scatter[0], scatter[1], scatter['index'], plt.gca())
142 plt.axvline(-0.5, ls='--')
143 plt.axhline(0, ls='--')
144 plt.title('Değişkenlerin Temel Bileşen 1 ve 2 içindeki pozisyonları',
145           fontsize=18)
146 plt.xlabel('Temel Bileşen 1 (50%)', fontsize=18)
147 plt.ylabel('Temel Bileşen 2 (24%)', fontsize=18)
148 plt.show()
149
```



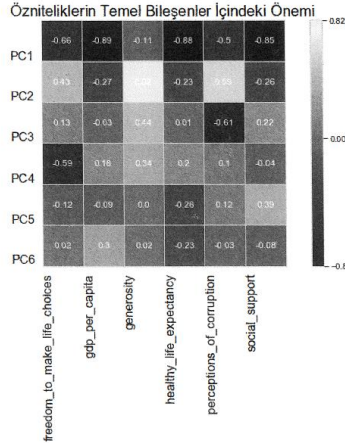


Unsupervised Learning - PCA

```

194 # Özniteliklerin her bir temel bileşen içinde açıklanan varyansa yaptığı
195 # katkıyı (Eigenvectors) tabiileştirilerek inceleyelim:
196 tpecs.column_correlations(df2[features]).transpose()
197 features_names=pd.Series(t.columns)
198 s= pd.Series(["PC1", "PC2", "PC3", "PC4", "PC5", "PC6"])
199 t=t.set_index(s)
200
201 # Grafiğe hazır hale getirmek için t dataframeini numpy array formatına
202 # dönüştürelim ve 2 basamaklı ondalık değerine yuvarlayalım:
203 t=t.to_numpy()
204 t=t.round(2)
205
206 # Graf: Özniteliklerin Temel Bileşenlerdeki Önemi:
207 fig, ax = plt.subplots(figsize=(12, 6))
208 plt.imshow(t, interpolation = 'none', cmap = 'plasma')
209 feature_names = features_names
210 # Temel bileşenler ve özniteliklerin grafikteki konumları...
211 ax.set_xticks(np.arange(-.5, len(feature_names)));
212 ax.set_yticks(np.arange(0.5, 6));
213 ax.set_xticklabels(feature_names, rotation=90, ha='left', fontsize=16);
214 ax.set_yticklabels(s, va='bottom', fontsize=16);
215 # Kutucuklarda eigenvalues değerlerini göstermek için...
216 for i in range(len(s)):
217     for j in range(len(feature_names)):
218         text = ax.text(j, i, t[i, j],
219                        ha="center", va="center", color="w")
220
221 plt.title('Özniteliklerin Temel Bileşenler İçindeki Önemi', fontsize=20)
222 plt.colorbar(orientation='vertical', ticks=[t.min(), 0, t.max()]);
223

```



Unsupervised Learning

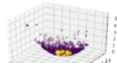
DIMENSION REDUCTION (Boyut Azaltma)

Giriş

Why Do We Need Dimension Reduction?

Visualization is one of the easiest ways to grasp data.

But Visualize > 3D?

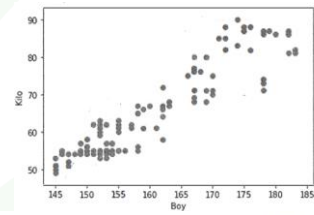


Some machine learning algorithms are affected by the multidimensionality curse

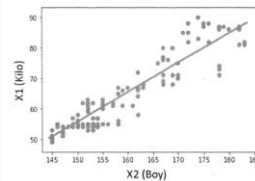


The more features we have, the more data we need.

Additional data collection is expensive and often not possible.



Boyut Azaltma



Burada boy (X1) ve kilo (X2) bir ilkök sınıfındaki öğrencilerin cinsiyetini tanımlayan iki öznitelik olsun. Her iki değişken de bir ölçüde tek başlarına öğrencilerin cinsiyetlerini tanımladığını varsayalım (ilkokul çağında erkekler kızlardan daha ağır; daha önce boy atıyor). Bu değişkenlerin her ikisini de kullanırsak, benzer bilgileri aktarırlar. Bu nedenle, yalnızca bir değişken kullanmak mantıklı olacaktır. Verileri 2D'den (X1 ve X2) 1D (Y1) 'e aşağıda gösterildiği gibi dönüştürebiliriz. Benzer şekilde, verilerin 4 sayıdaki boyutunu aynı bilgiyi içermek kaydıyla bir alt kümesine indirgeyebiliriz. Buna boyut azaltma denir.

y1 vektörü



Unsupervised Learning

✓ DIMENSION REDUCTION (Boyut Azaltma)

- ✓ Feature Selection
- ✓ Feature Extraction



Unsupervised Learning

✓ DIMENSION REDUCTION (Boyut Azaltma)

- ✓ PCA ile Dimension Reduction



- 2 human
- children
- long hair
- put their hands together
- height etc.

The area marked in Set A contains approximately **70% of the information** contained in the pink, blue and yellow circles.



Set A



2D (shadow of the children) *represent most of the information of real picture (3D)*

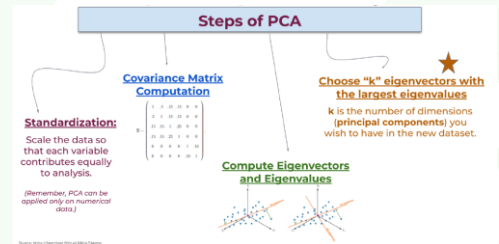
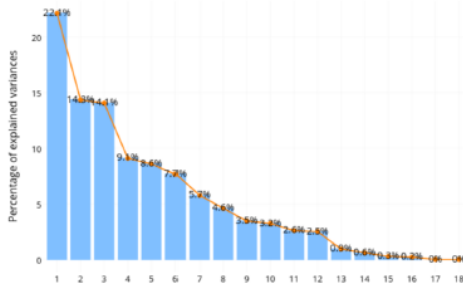




Unsupervised Learning

DIMENSION REDUCTION (Boyut Azaltma)

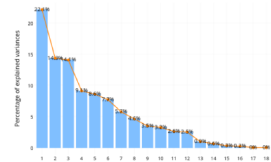
PCA ile Dimension Reduction



```
sklearn.decomposition.PCA
class sklearn.decomposition.PCA(n_components=None, *, copy=True, whiten=False, svd_solver='auto', tol=0.0,
                                iterated_power='auto', random_state=None)
```

Choose "k" eigenvectors with the largest eigenvalues

k is the number of dimensions you wish to have in the new dataset.



Unsupervised Learning

DIMENSION REDUCTION (Boyut Azaltma)

PCA ile Dimension Reduction

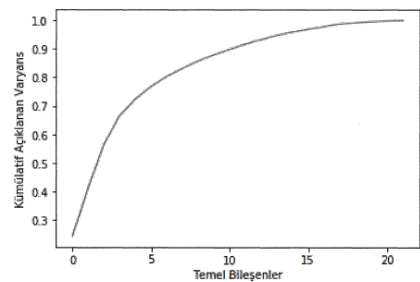
```
26 # KNN Sınıflandırıcısı:
27 n_neighbors=3 # en yakın komşu sayısını 3 seçelim.
28 KNN = neighbors.KNeighborsClassifier(n_neighbors, weights='distance')
29 KNN.fit(X_train, y_train)
30
```

```
31 # Eğitim Seti performansı
32 pred_invest_grade_train = KNN.predict(X_train)
33 print('Doğruluk (Accuracy): %0.2f' % accuracy_score(y_train, pred_invest_grade_train))
34 # Test Seti Performansı
35 pred_invest_grade_test = KNN.predict(X_test)
36 print('Doğruluk (Accuracy): %0.2f' % accuracy_score(y_test, pred_invest_grade_test))
37
```

Doğruluk (Accuracy): 1.00
Doğruluk (Accuracy): 0.92

```
49 # top. varyansın %80'ini açıklayan temel bileşenlere
50 # göre modeli tekrar oluşturalım:
51 pca = PCA(n_components = 0.80)
52 X_pca = pca.fit_transform(X_train) # eğitim setini dönüştürüp boyutunu azaltır.
53 print(pca.n_components_) # top. varyansın %80'i kaç temel bileşen ile açıklanıyor?
54
```

```
38 # sklearn kütüphanesi: PCA ile Boyut Azaltma
39 df.iloc[:, 1:23].columns # özellikler
40
41 pca = PCA()
42 X_pca = pca.fit(X_train)
43
44 # Grafik: Kümülatif Açıklanan Varyans (KAV)
45 plt.plot(np.cumsum(pca.explained_variance_ratio_))
46 plt.xlabel('Temel Bileşenler')
47 plt.ylabel('Kümülatif Açıklanan Varyans');
48
```





Unsupervised Learning

DIMENSION REDUCTION (Boyut Azaltma)

✓ PCA ile Dimension Reduction

✓ Avantajları

✓ Dezavantajları

PC1	0.0	0.06	-0.03	0.08	-0.06	0.01	0.0	-0.01	-0.01	0.04	-0.02	0.01	0.01	-0.09	-0.03	-0.99	0.03	-0.01	-0.02	0.0	0.01	0.02
PC2	0.01	-0.06	0.06	-0.05	0.04	-0.01	0.01	-0.02	-0.04	-0.11	0.06	-0.01	-0.08	-0.02	-0.02	0.04	-0.05	0.03	0.0	-0.06	0.19	0.02
PC3	0.02	0.01	-0.01	-0.09	-0.04	0.01	0.0	-0.01	-0.04	-0.02	0.22	0.12	-0.02	-0.15	0.06	0.04	0.05	0.06	0.11	0.01	-0.04	0.02
PC4	-0.02	-0.04	0.01	-0.02	0.02	-0.04	-0.01	0.11	0.05	0.1	-0.05	0.09	-0.01	-0.93	0.03	0.08	-0.03	-0.2	0.09	0.0	-0.07	-0.2
PC5	0.08	0.02	0.08	0.11	0.25	0.07	0.01	0.04	0.25	0.21	0.03	-0.63	0.02	-0.14	0.08	0.01	0.01	0.2	0.14	-0.01	0.01	0.18
PC6	-0.04	0.19	0.09	0.07	0.17	0.32	0.02	0.44	0.12	-0.65	-0.05	-0.04	0.0	-0.03	0.1	-0.01	0.13	-0.09	-0.21	0.01	0.11	0.0
PC7	0.04	0.08	-0.04	0.12	-0.09	0.12	0.02	0.40	0.17	-0.25	0.08	0.05	0.01	0.19	0.02	0.04	0.25	-0.43	0.06	0.03	0.04	0.04
v1																						
v2																						
v3																						
v4																						
v5																						
v6																						
v7																						
v8																						
v9																						
v10																						
v11																						
v12																						
v13																						
v14																						
v15																						
v16																						
v17																						
v18																						
v19																						
v20																						
v21																						
v22																						

```

104 # Eğitim Seti Performansı
105 pred_invest_grade_pca_train = clf.predict(PC_Df)
106 print('Doğruluk (Accuracy):%0.2f'%
107       accuracy_score(y_train,pred_invest_grade_pca_train))
108 # Test Seti Performansı
109 pred_invest_grade_pca_train = clf.predict(PC_Df)
110 print('Doğruluk (Accuracy):%0.2f'%
111       accuracy_score(y_train,pred_invest_grade_pca_train))
112

```

Doğruluk (Accuracy):1.00
Doğruluk (Accuracy):1.00

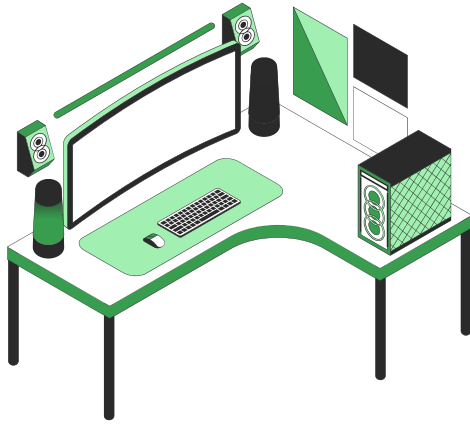


Bu dersi anladım..



Everything is clear ?





Do you have any questions?

Send it to us! We hope you learned something new.

