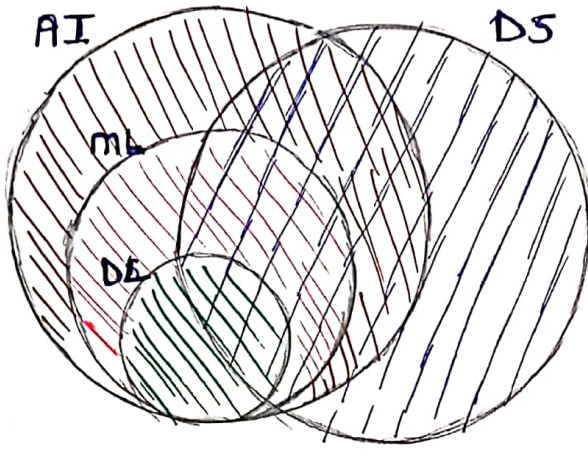
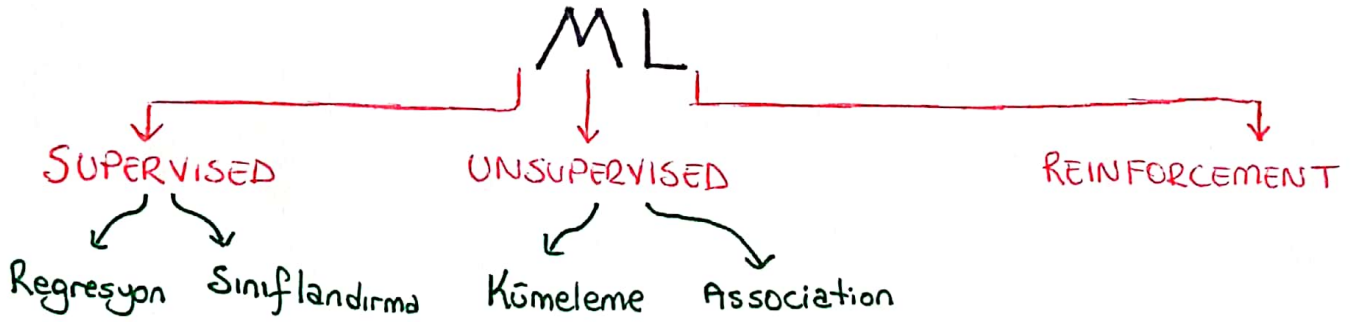


MAKİNE ÖĞRENMESİ



Görüleceği üzere AI ve DS kesişen bir kümededir.

MAKİNE ÖĞRENMESİ SINIFLANDIRMA



• Yukarıdaki kavramlar üzerinde duralım.

• **Supervised Learning** = veri + label \Rightarrow Denetimli Öğrenme

• **Unsupervised Learning** = veri + no label \Rightarrow Denetimsiz Öğrenme

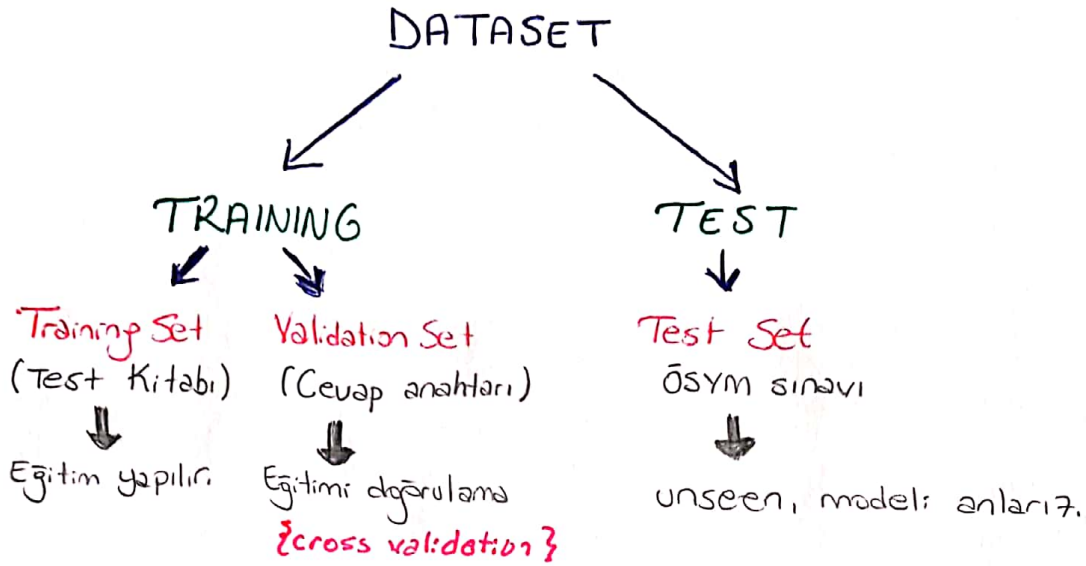
• **Reinforcement Learning** = Deneme - yanılma \Rightarrow Pekiştirmeli Öğrenme

• **Regresyon** = Sürekli, sayısal sonuç verir. Ev fiyatı tahmini.

• **Sınıflandırma** = Sonuç kategoriktir. Spam - spam değil.

• Şimdi de Scikit-learn haritasını inceleyelim.

Bu işlemler yapılırken verimizi iyi şekilde bölmemiz gerekir.

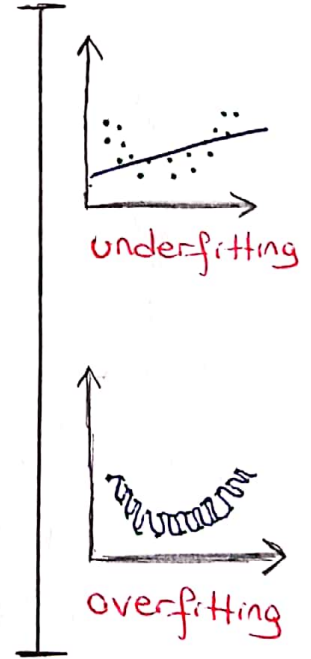
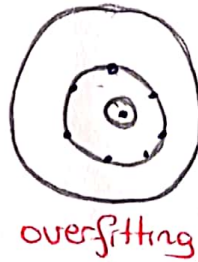
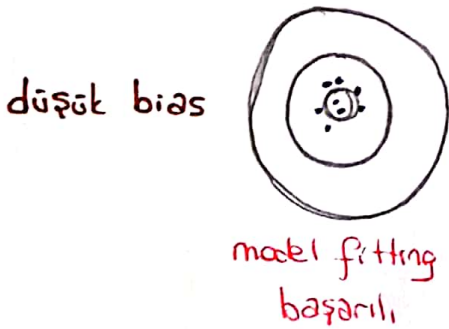
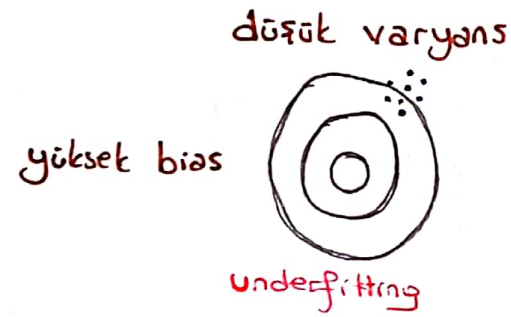


BIAS - VARYANS

bias = merkeze, doğruluğa uzaklık.

varyans = yayılımı verir.

- under-fitting = yüksek bias düşük varyans = Hiç öğrenmeme = doğruluğa çok uzak.
- over-fitting = yüksek varyans düşük bias = ezberleme olmuş, çok iyi dağılmış.



↓ bias ↓ varyans = başarı

REGRESYON

Basit Lineer Reg.

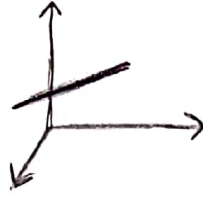
- 1 tane bağımsız değişken



$$y = wx + b$$

Goklu Lineer Reg.

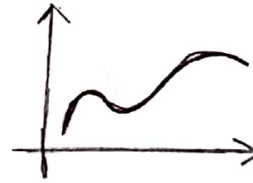
- Birden fazla bağımsız değişken



$$y = b_0 + b_1x_1 + b_2x_2 \dots$$

Polinomial Regression

- bağımsız derece n'inci derecedir.



$$y = x^n \dots n \leq 2$$

① KOD GRAMERİ → BASİT VE GOKLU

from sklearn.linear-model import LinearRegression → kütüphane ekledik

LR = LinearRegression() ⇒ LR objesi oluşt.

LR.fit(X_train, Y_train) ⇒ eğitim

y_pred = LR.predict(X_test) ⇒ tahmin

② KOD GRAMERİ → Polinomial

from sklearn.preprocessing import PolynomialFeatures

poly-reg = PolynomialFeatures(degree = n)

X_poly = poly-reg.fit_transform(X)

LR.fit(X_poly, y)

$$y = wx + b$$

• score = R² score

• coef_ = w = ağırlık

• intercept_ = b

$$y = \text{coef_} \cdot x + \text{intercept_}$$

REGRESYON HATA AZALTMA

① Early Stopping

- Gradient Descent gibi yinelenmeli algoritmalarda düzenlilik gerekir. Doğrulama hatası min ulaşıncaya eğitim durur.
- callback

② L1 - L2 Regularizasyon \Rightarrow overfitting önlenir.

• L1 LASSO

Gradient Descent Denklemi + **mutlak değer**

Mutlak değer ile cezalandırılır.

0 olabilir.

$$\frac{\lambda}{m} |\theta_i|$$

from **sklearn.linear-model** import **Lasso**

model = **Lasso** (alpha = n) \Rightarrow fit ve predict \Rightarrow metric

• L2 RIDGE

Gradient Descent Denklemi + **kareli değer**

Karesi ile cezalandırılır.

0 asla olamaz.

$$\frac{\lambda}{m} \theta_i^2$$

from **sklearn.linear-model** import **Ridge**

model = **Ridge** (alpha = n) \Rightarrow fit ve predict \Rightarrow metric

KAVRAMLAR :

α = Learning Rate = Gradient Descent çalışırken yakınsama hızı. :::::

θ = theta = feature ağırlığıdır.

λ = Regularizasyon term = Ağırlığın ne kadar cezalandırılacağını belirler.

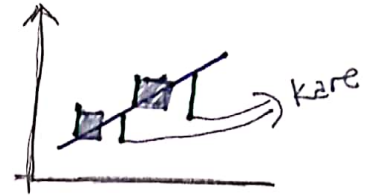
REGRESYON PERFORMANS ÖLÇÜMÜ

{metrics}

Amaç: y ile y -pred arasındaki hatayı ölçmektir.

① R Squared

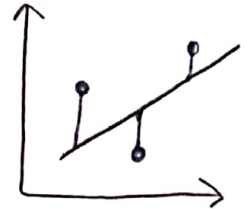
```
from sklearn.metrics import r2_score  
r2_score(y, y-pred)
```



$$R^2 = 1 - \frac{SS_{reg}}{SS_{total}}$$

② Mean Absolute Error = MAE

```
from sklearn.metrics import mean-absolute-error  
mean-absolute-error(y, y-pred)
```



- Fırtınalı değerlerden etkilenmez. \Rightarrow Avantaj
- Minimumda türevlenebilir değil \Rightarrow Dezavantaj = Gradient Descent olmaz.

$$MAE = \frac{1}{n} \cdot \sum |y - \hat{y}| \Rightarrow \text{farkların mutlak değerinin ortalaması}$$

③ Mean Squared Error = MSE

```
from sklearn.metrics import mean-squared-error  
mean-squared-error(y, y-pred)
```

- En sık kullanılan regresyon hata fonksiyonudur.
- Türevlenebilir. \Rightarrow Çünkü üssü var.

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

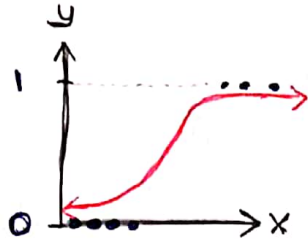
! $y = y_{-test}$

CLASSIFICATION

- 2'li veya 2'den fazla sınıflandırma olabilir.

① LOGISTIC REGRESYON

- Sınıflandırma işlemi yapan regresyon modelidir.



spam / not spam } 0 ile 1 arasında olasılık değeri
yes / no } **threshold** = eşik değeri

Linear Regression → Logistic Regression

Activation Function = * sigmoid, tanh, RELU

→ SOFT MAX = 2'den fazla sınıf.

```
from sklearn.linear_model import LogisticRegression
LR = LogisticRegression(*)
LR.fit(X_train, y_train)
LR.predict(X_test)
```

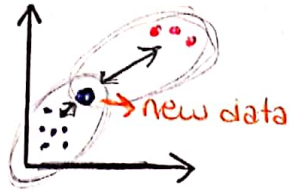
penalty
random_state
⋮

Kısaca özet geçecek olursak:

- Linear fonksiyon + Sigmoid Function = Logistic Regresyon.
- Eşik değeri önemlidir.
- S-shape

! sklearn.metrics → classification_report

② KNN = K-Nearest Neighbors



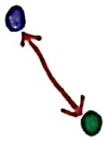
- Çalışma mantığı : en yakın komeye yeni datanın eklenmesidir.

k = bakılması istenen komşu sayısıdır.

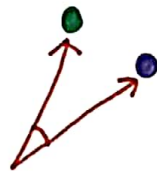
k=5 olsun. \Rightarrow new datanın en yakın 5 komşusuna bakılır.

3 mor , 2 kırmızı varsa \Rightarrow new data = mor olur.

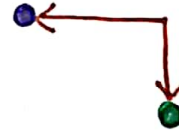
- Mesafe Ölçümleri



euclidean



cosine



manhattan

...

- En iyi k değeri = Decision Boundary



bu nokta en iyi k değeridir. Elbow point

```
from sklearn.neighbors import KNeighborsClassifier
```

```
KNN = KNeighborsClassifier ( n-neighbors = n )
```

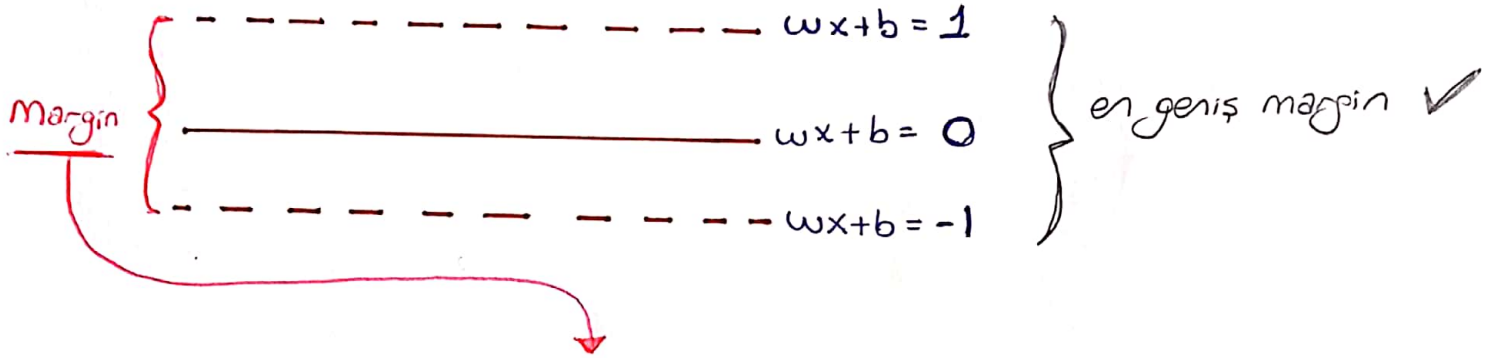
```
KNN.fit ( X-train , y-train )
```

```
KNN.predict ( X-test )
```

③ Support Vector Machine = SVM

SVM : sınıflandırma , regression , aykırı değer sonucu üretir.

- Küçük ve orta ölçekli ^{veriler için} çok uygundur.



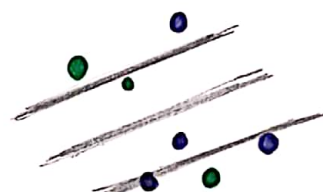
MARGIN

HARD MARGIN



- margin ı ar boş

SOFT MARGIN



- marginde olabilir. Yanlış değerler olabilir.

overfitting
engeller.

```
from sklearn.svm import SVC
```

```
svm_model = SVC(kernel=...) → rbf
```

```
svm_model.fit(X_train, y_train)
```

```
svm_model.predict(X_test)
```

Kernel Library



Gaussian



Polynomial

RBF



Kernel Trick



Boyut artırmak



boyut ↑



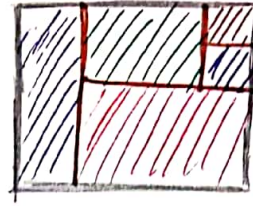
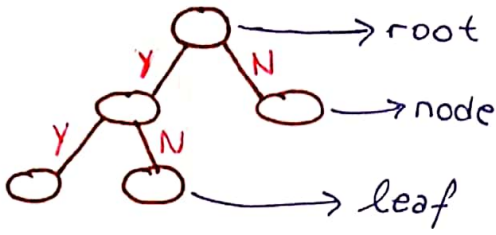
sklearn.kernel_approximation



Nystroem

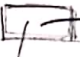
RBFsampler

④ Decision Tree



⇒ overfitting
↳ yaygın
!

ID3


• pure =  → tüm veriler aynı sınıfa aittir.

• impure =  → tüm veriler aynı sınıfa ait değildir.

• Gini impurity = verinin ne sıklıkla yanlış tanımlandığını anlatır.
↳ Gini INDEX

! Gini Index ↓ ⇒ Tercih Edilme ↑

ENTROPY = Düzensizlik


Entropi: low
Bilgi: high


medium
medium


high
low

! graphviz

! pruning ⇒ optimize

INFORMATION GAIN

• 0 koto - 1 iyi

• $IG = \text{entropy}(\text{parent}) - \frac{1}{n} [\text{entropy}(c_1) + \dots + \text{entropy}(c_n)]$ ↳ dal sayısı

! Gini = [0, 0.5]

Entropy = [0, 1]

from sklearn.tree import DecisionTreeClassifier

DTC = DecisionTreeClassifier()

DTC. ~~predict~~ (X_train, y_train) ⇒ DTC. predict (X_test)

} aynı

fit

ENSEMBLE LEARNING

- Birden fazla algoritma geliştirilmesi ile sonuç verir.
 - Maximum oylama
 - Ortalama
 - Ağırlıklı ortalama
 - Yığma {Stacking}
 - Blending

BAGGING

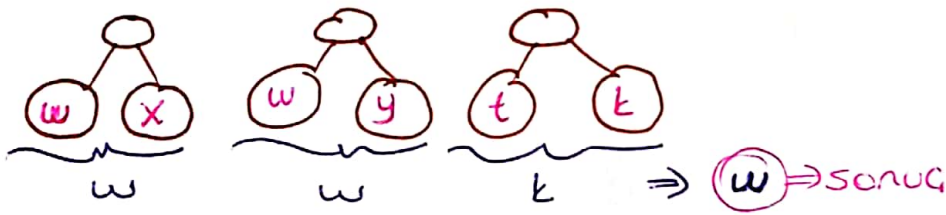
- Random Forest

BOOSTING

- Ada Boost
- Gradient Boosting
- XG Boost

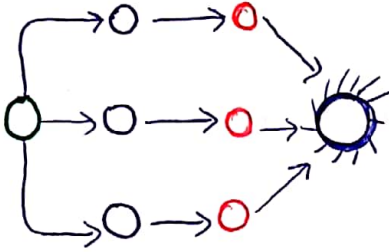
⑤ Random Forest

- Decision tree'deki overfitting sorunu için alternatiftir.
- Kolonlar rastgele birleşir ve birden çok ağacı oluştur.
 - ↳ Sonucunda en iyisi seçilir.



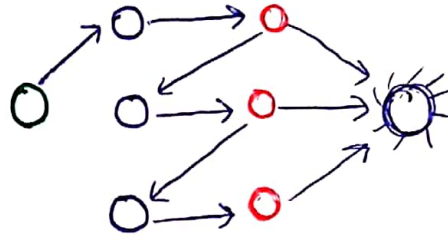
```
from sklearn.ensemble import RandomForestClassifier  
RF = RandomForestClassifier ( )  
RF.fit(X_train, y_train)  
RF.predict(X_test)
```

BAGGING - BOOSTING



parallel

Her algoritma ayrı ayrı çalışır ve sonunda birleştirilir.



sequential

Her algoritmanın en iyi olduğu kısım göz önüne alınır. Sıralı olarak birleştirilir.

○ ⇒ weak learner
⊙ ⇒ strong learner

⑥ Adaboost

doğruluk



weight

LARGE
+

medium



Large
-

modelin ağırlığına göre
birleştirme yapılır.

$$\begin{array}{c} + \\ \downarrow \\ \text{weight} \uparrow \end{array} \quad \begin{array}{c} - \\ \downarrow \\ \text{weight} \downarrow \end{array} = + \text{ olur.}$$

base-estimator = weak learner ⇒ DecisionTree ... vs.

n-estimator = kaç tane learner olacak?

! Kod grameri aynıdır. DecisionTree vs. import etmeyi unutma.

⑦ Gradient Boosting

⑧ XGBoost

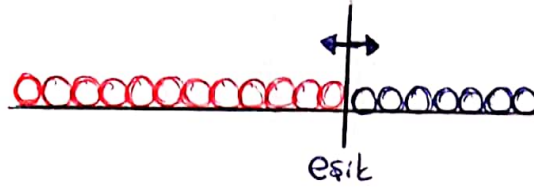
⑨ Naive Bayes

SINIFLANDIRMA PERFORMANS ÖLGÜMÜ

{metrics}

- Bu bölümde ölçüm ve hata azaltma üzerine duracağız.

KAVRAM Thresholding = Eşik değerimiz değiştirilirse performans değişir.



KAVRAM TRUE - FALSE & Positive - Negative

- Ölçüm yapabilmek için TN, TP, FN, FP gereklidir.

Positive = YES } olsun. \Rightarrow NO tahmin ettik. YES tahmin ettik.
Negative = NO } DOĞRU ise Yanlış ise DOĞRU ise Yanlış ise
TN FN TP FP

① CONFUSION MATRIX

TP	FP
FN	TN

from sklearn.metrics import confusion_matrix
confusion_matrix(y-test, y-pred)

② Accuracy

- Doğru tahmin edilenlerin tüm tahminlere bölümü ile bulunur.

from sklearn.metrics import accuracy_score
accuracy_score(y-test, y-pred)

$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

\uparrow = yanlışlar yakalanmıyor.
yani kötü

③ Precision & Recall

Precision = "Pozitif tahminlerin ne kadarı doğru?" cevabını verir.

Recall = "Gerçekle pozitif olanların kaçı doğru?" cevabını verir.

④ F1 SCORE = F_β Score

• Precision ve Recall arasında seçim yapmak zordur.

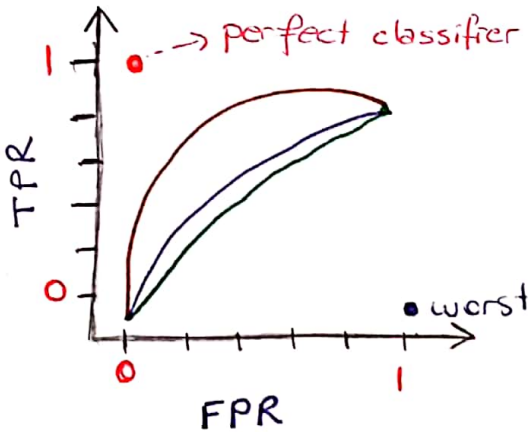
F1 Score = Precision ve Recall harmonik ortalamasıdır.

→ Aynı eşit katkı ile doğru sonuç.

⑤ ROC - AUC

ROC = Her threshold adımı için metrikleri gösterir.

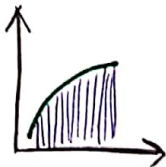
$$TPR = \frac{TP}{TP + FN}$$



kahverengi > mor > yeşil

↓
daha iyi sınıflandırma •'a yakın.

AUC = ROC eğrisinin integrali {altında kalan alan}



! sklearn.metrics'den kod kısmını yazabilirsin.

! Aşağıda formüllere göt geçirebilirsin.

PYTHON LISTS []

İçerisinde farklı veri tipleri barındırır

fam = [1.73 , "dogukan" , True , False , "mom" , 32 , "sewal"]

fam2 = [["dogu", 174],
["arif", 172]] } liste iacide liste de olabirin

`type(fam)` `type(fam2)` \Rightarrow list

 \propto Erizim

fam[0] = 1.73 fam[-1] = ~~True~~ "several"
 ↘ genislethin. ↗

\propto Slicing

$\text{fam} [x : y]$
 ↙ ↘
 başlangıç bitiş index
 dahil dahil değil

$\text{fam} [: 2] = 2.$ indexe kadar olan
 $\text{fam} [1 :] = 1$ dahil ve sonrası

\propto Manipulating Lists

- Changing list elements

- `fam[4] = 1.86` \Rightarrow 4. indexteki değer "mom" \rightarrow 1.86 olur.
- `fam[0:2] = ["lisa", 1.74]` \Rightarrow 0 ve 1. index değeri değişir.

- Add list elements

$fam = fam + ["me", 1.79] \Rightarrow fam$ 'in sonuna iki deger ekledik.

- Remove list elements

`del(fam[2])` \Rightarrow True değeri silinir.

list \rightarrow Numpy

$\text{np-array} = \text{np.array}(\text{list})$

$$\max(f_{am})$$

fam. index ("mom")

fam. count (1.73)

fam. append (x)

DICTIONARIES { }

- Listeleri incelemiştik. Eğer listeler arasında bağlantı kuruyorsak dictionaries bizim için bir çözümdür.

pop = [30, 40, 50]
country = ['TR', 'FR', 'BE'] } 2 bağlantılı liste

world = { 30: 'TR', 40: 'FR', 50: 'BE' } - dictionary
 key value

işlemler için = []

[word [key] = value]

α YENİ DEĞER EKLEME

world [60] = 'PL' ⇒ sözlüğe yeni bir değer ekledik.

α DEĞER SİLME

del(world[key])

world.keys() - world.values()

pd.DataFrame(dict)

MATPLOTLIB

```
import matplotlib.pyplot as plt
```

```
year = [1950, 1970, 1990, 2010]
```

```
pop = [2, 3, 5, 7]
```

```
plt.plot(year, pop) = line plot
```

```
plt.scatter(year, pop) = scatter plot
```

```
plt.hist(year, bins=3) = hist plot
```

```
plt.xlabel('Year')
```

```
plt.ylabel('Pop')
```

```
plt.title('Title')
```

```
plt.yticks([0, 2, 4, 6, 8, 10])
```

```
plt.yticks([0, 2, 4, 6, 8, 10], ['0B', '2B', '4B', '6B', '8B', '10B'])
```

```
plt.xscale('log')
```

```
plt.grid(True) #
```

import **NUMPY** as np

- Listeleri daha iyi ve matematiksel olarak manipüle etmemize yarar.

list1 = [1.73, 1.86, 1.55, 1.90, 2.00]

list2 = [60, 70, 80.2, 95, 105.4]

list1 / list2 ** 2 ⇒ **ERROR X**

SOLUTION

np-list1 = np.array(list1) } listeleri → numpy array yaptık.
np-list2 = np.array(list2) }

np-list1 / np-list2 ** 2 ⇒ **SONUÇ ✓**

- list1 + list2 ⇒ [1.73, 1.86, ..., 2.00, 60, 70, ..., 105.4] ⇒ birleştirir
- np-list1 + np-list2 ⇒ [61.73, 71.86, ..., 107.4] ⇒ toplar

numpy subsetting

np-list > 1 ⇒ [True, False, True, True, ...]

np-list[np-list > 1] ⇒ True ise değeri verir ⇒ filtrleme

2D numpy array

np.array([[1.73, 1.68, 1.71, 1.89, 1.79],
[65, 59, 63, 88, 68]]) } liste

x[row][col] = x[row, col] ⇒ x[1, 2] = 63

x[:, :] ⇒ tüm satır ve kolonu getirir ⇒ filtrleme ⇒ ayrıntılı var. ↓

np.mean(np-city[:, 0]), median, std, sum, sort, round

np.corrcoef(np1, np2) x.shape logical_and() - or - not

PANDAS

```
import pandas as pd
```

```
brics = pd.DataFrame(dict)
```

```
brics = pd.read_csv("----")
```

} DataFrame oluşturduk.

Filtreleme - ACCESS

column {
brics["column"] = pandas.Series
brics[["column"]] = pandas.DataFrame
brics[["col1", "col2"]] ⇒ 2 kolonlu df olabilir.

row {
brics[1:4] ⇒ 1, 2, 3. satırı verir.

- **loc** ⇒ Kategorik verilerde işlem yapılır.

brics.loc[["RU"]] = RU satırını verir.

brics.loc[["RU", "AL"]] = RU ve AL satırını verir.

brics.loc[["RU", "AL"], ["col1", "col2"]] = row-col loc

brics.loc[["RU", "AL"], :] = tüm kolonlar, ru-al satırı

- **iloc** ⇒ integer {index} ile erişim yapılır.

- Karşılaştırmalı filtreleme

brics[brics["col"] > 8] ⇒ !np.logical kullan.!

loop ⇒ datacamp / Intermediate Python / Loops

random ⇒ datacamp / Intermediate Python / Case Study