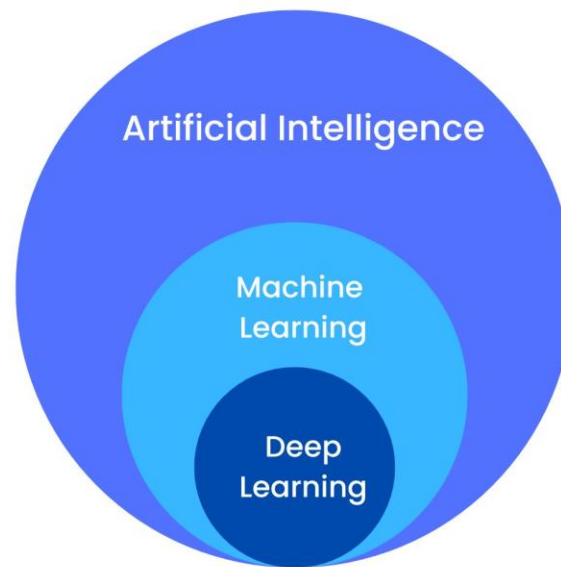


# STROJNO UČENJE

V BIOFIZIKI

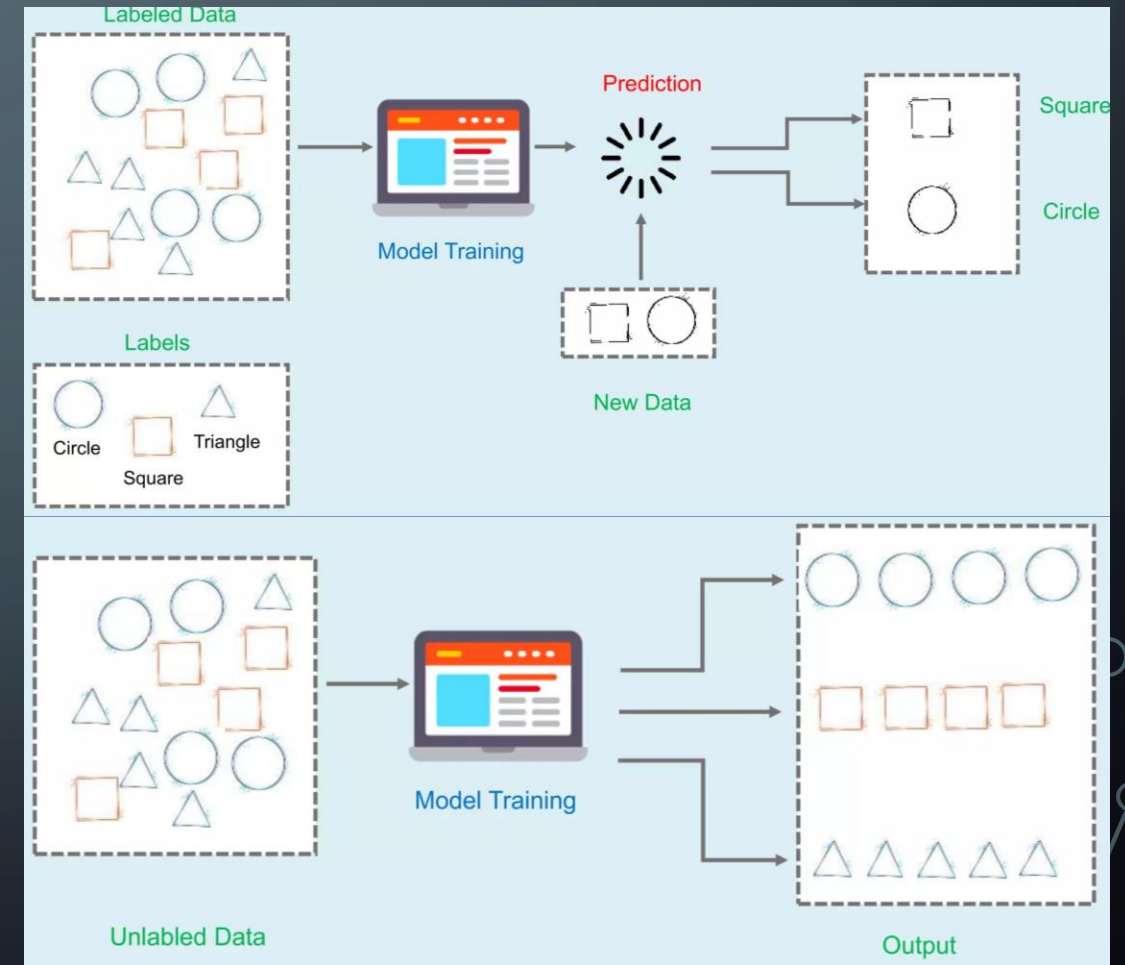
# KAJ JE STROJNO UČENJE?

- Učenje je vsak proces pri katerem sistem **izboljša** svoje delovanje preko novih **izkušenj**
- Strojno učenje se ukvarja z računalniškimi programi, ki se avtomatsko izboljšujejo preko novih izkušenj
- Podvrsta umetne inteligence, pri kateri se računalnik "sam" nekaj nauči, ne da bi to **eksplicitno** sprogramirali.
- Glavna ideja: napoved na podobnih podatkih

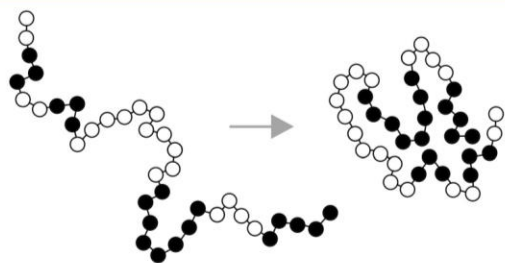


# TIPI STROJNEGA UČENJA

- Nadzorovano učenje
  - Označeni podatki (labeled data)
- Nenadzorovano učenje
  - Neoznačeni podatki
- Spodbujevalno (reinforcement) učenje



reinforcement  
learning

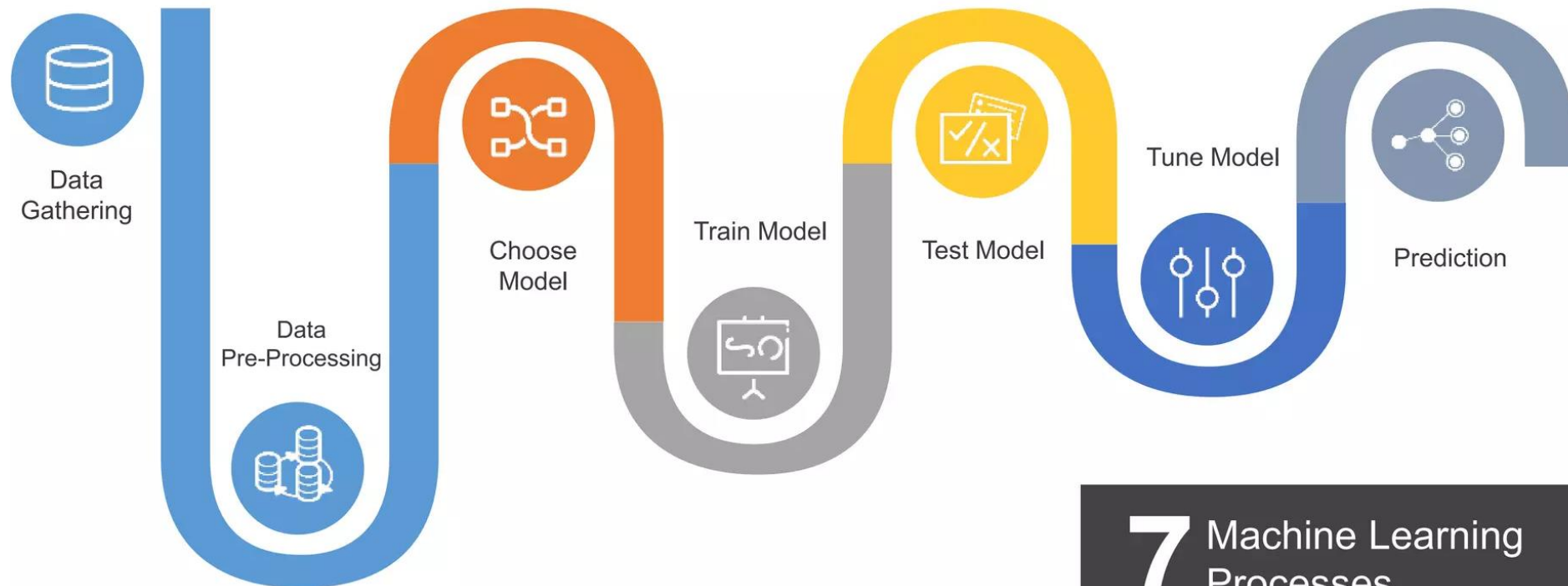


**input:**  
polypeptide chains modeled as a  
sequence of amino acid polarities

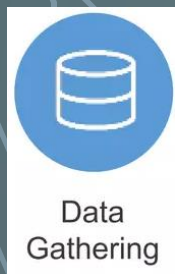
**output:**  
polypeptide conformation with  
minimal free energy

reference: [263]

## Processes involved in Machine Learning

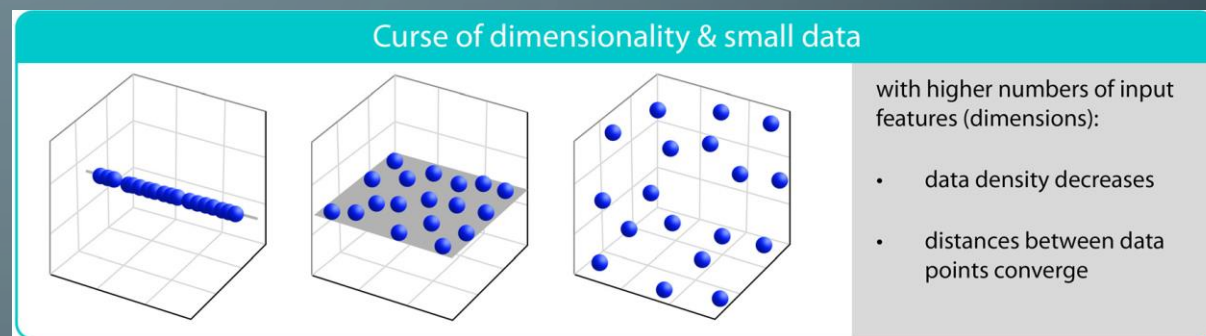


**7** Machine Learning Processes



# PODATKOVNI SETI

- *Garbage in, garbage out*
- Značilke (features) - fizikalne opazljivke
  - *Curse of dimensionality*
  - Korelirane opazljivke
- Potrebujemo velike sete podatkov (statistični modeli)
  - Označevanje podatkov
- Set za trening in set za testiranje



## reducing dimensionality

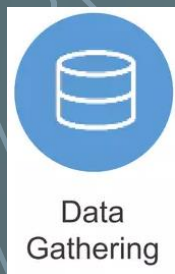
### feature elimination<sup>328,302</sup>

- removing strongly correlating features
- ranking features by importance and removing the weakest

### principal component analysis<sup>329</sup>

- condensing features to combined components





# PODATKOVNI SETI

- Umetno povečevanje setov (*data augmentation*)
- Eksperimentalni podatki ali simulacije
- Vektorizacija podatkov
  - Pogosto potrebna pred-obdelava, npr. skaliranje podatkov, maskiranje, ...

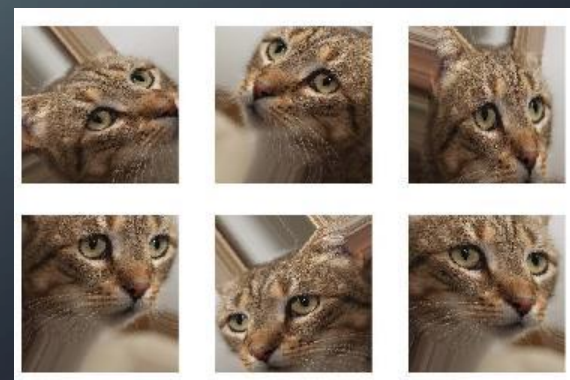
## increasing data size

### **data augmentation**<sup>335</sup>

- adding slightly altered copies of already existing data
- creating new data from existing data

### **data simulation**<sup>326</sup>

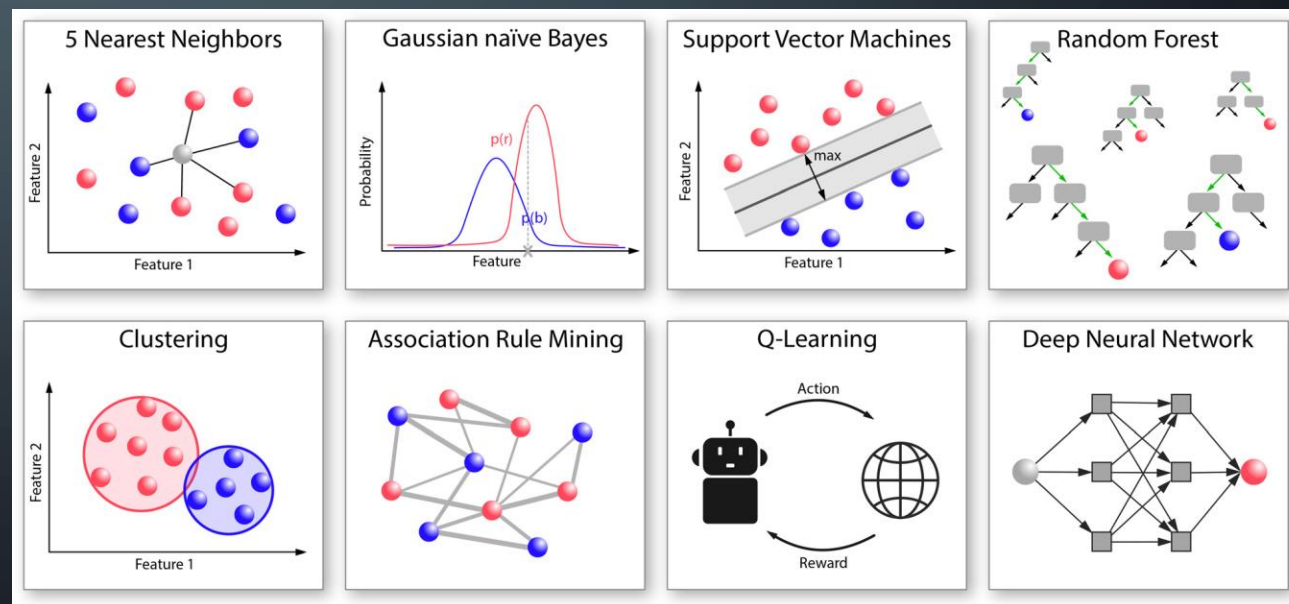
- generating data from real-world simulations





# MODELI STROJNEGA UČENJA

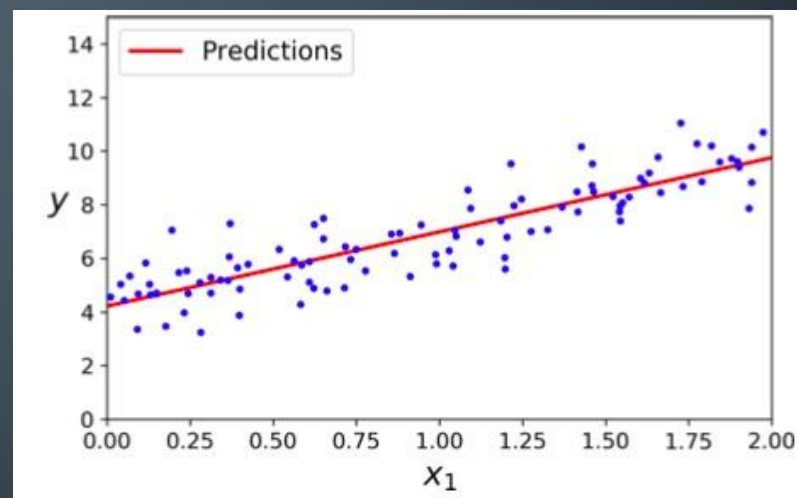
- No free lunch theorem (1996)
  - Če ni nobenih predpostavk (omejitev) na vrsti podatkov, ni formalnega razloga, da bi bil kakšen ML pristop boljši od drugega!
  - Ena metoda ni najboljša za vse probleme
- Regresija, klasifikacija
  - Napoved vrednosti ali razreda





# "FITANJE" (REGRESIJA)

- Model se "uči" z dodajanjem novih podatkov
- Obstaja celo točna rešitev!
- Če imamo veliko količino podatkov - numerika

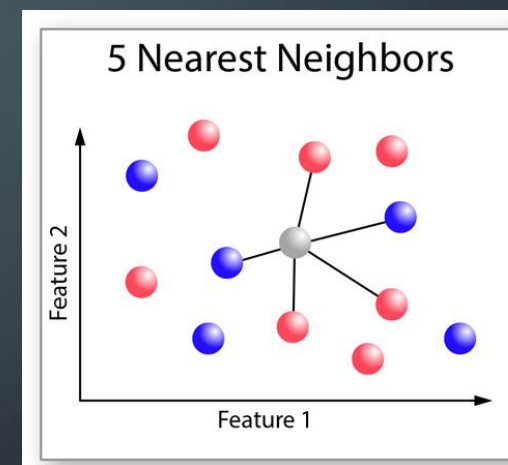






# K NEAREST NEIGHBOUR (KNN)

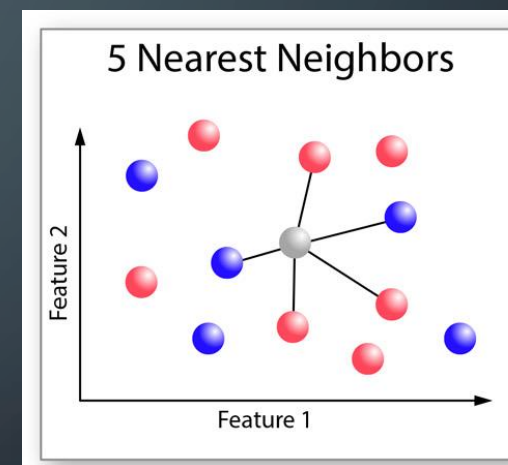
- Klasifikacijski algoritem
- Podatke gledamo v hiperprostoru značilk
- Premisa: Podatki, ki so blizu v hiperprostoru pripadajo istemu razredu
- Glede na že označene podatke za nov podatek:
  - Poiščemo K najbližnjih sosedov
  - Pogledamo katera kategorija je najpogostejša





# K NEAREST NEIGHBOUR (KNN)

- Metrika razdalje
- Problem *outlierjev* (zahtevajo velik  $k$ ), a na ta način izgubimo na natančnosti deljenja (robne točke razredov so problematične)
- Uravnoteženost podatkov (predprocesiranje - skaliranje velikosti)
- Curse of dimensionality





Choose  
Model

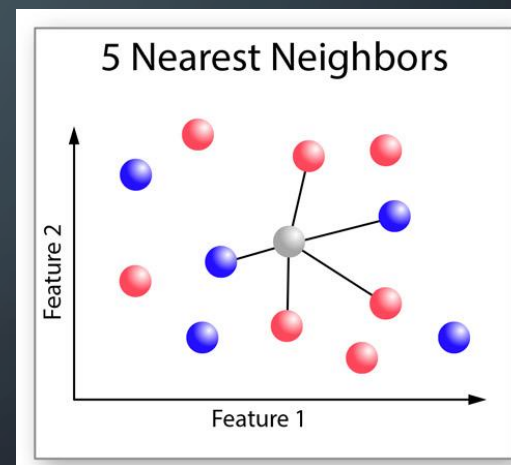
# K NEAREST NEIGHBOUR (KNN)

- Primer uporabe:
  - Klasifikacija okuženih s COVID-19: uporaba rentgenskih slik pljuč -> procesiranje slik -> generiranje značilk -> klasifikacije s KNN

## K nearest neighbors<sup>69-72</sup>

No training phase needed  
Intuitive and simple algorithm  
Easily adapts to new training data  
Only one hyperparameter to tune

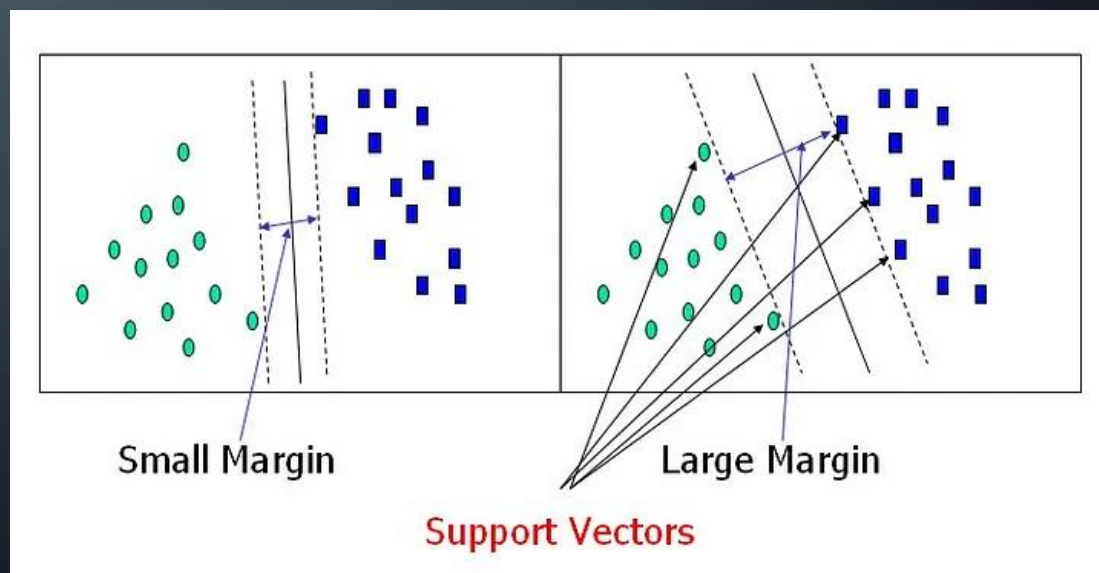
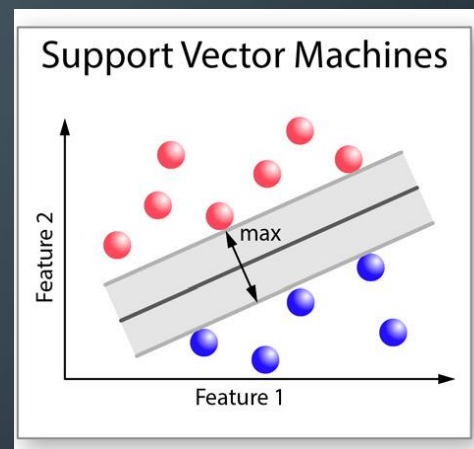
High dimensionality leads to decreased accuracies  
Can become slow for big datasets  
Needs feature scaling  
Has problems with imbalanced datasets  
Missing values are problematic

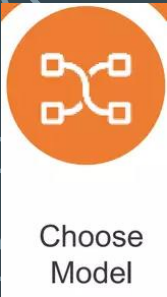




# SUPPORT VECTOR MACHINES (SVM)

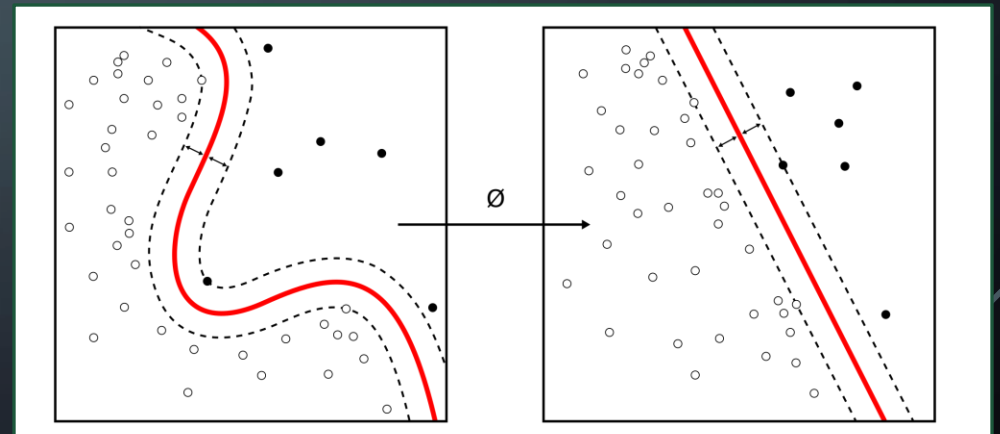
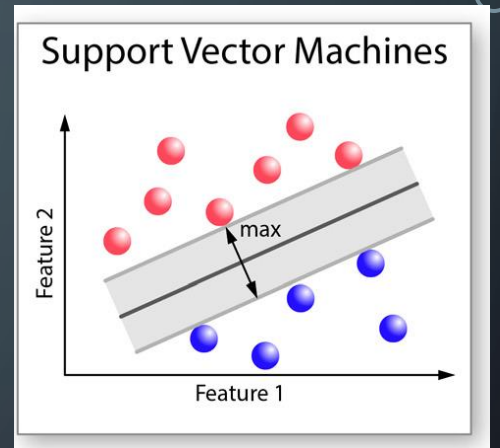
- Binarna klasifikacija
- V hiperprostoru značilk iščemo hiperpovršino, ki razdeli prostor na dva dela – SVM poišče ploskev, kjer je najširši pas nezasedenosti
  - Potem klasifikacija glede na to, na katerem delu pasu se nahaja





# SUPPORT VECTOR MACHINES (SVM)

- Problemi, ko se območja prekrivajo (ali ko so podatki zelo zašumljeni)
- Z uporabo jeder lahko transformiramo podatke razdelimo podatke nelinearno



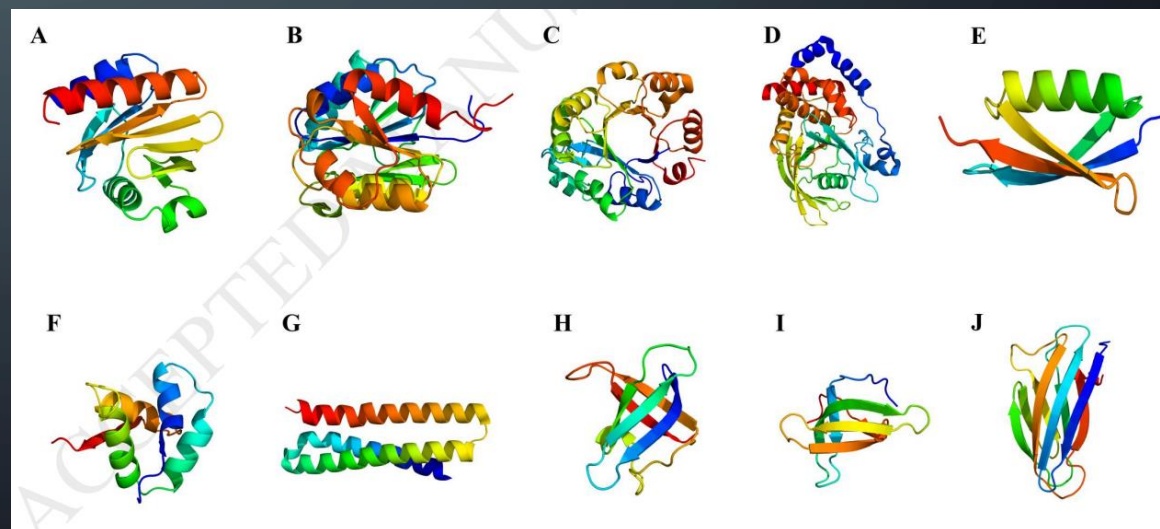
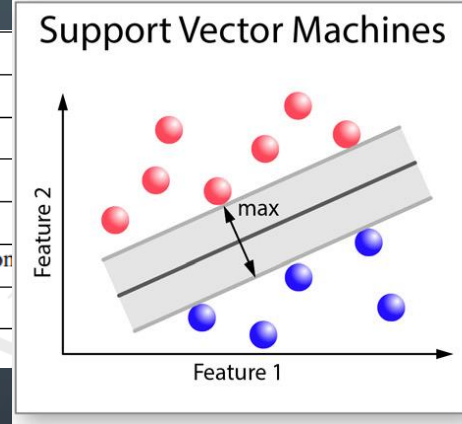




# SUPPORT VECTOR MACHINES (SVM)

- Primer uporabe: Analiza strukture zgibanja proteinov
  - Povezava med biofizikalnimi lastnostmi in obliko proteina – proteini, ki imajo različno sekvenco AK, a podobne povprečne biofizikalne lastnosti, se podobno zgibajo
  - SVM uporabljen za klasifikacijo med razredi tipičnih oblik proteinov

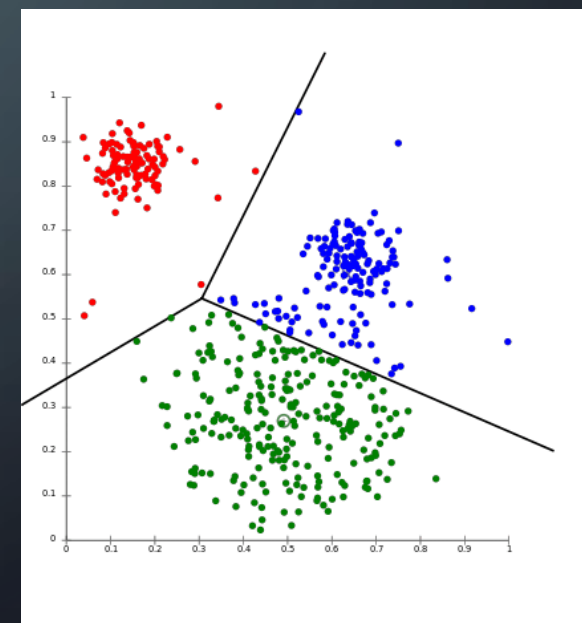
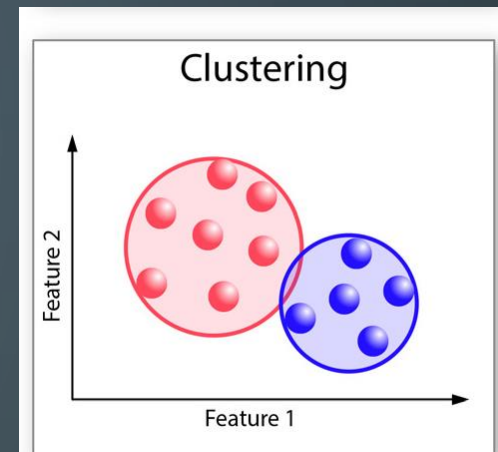
1	<b>K0</b>	Compressibility
2	<b>Ht</b>	Thermodynamic transfer hydrophobicity
3	<b>Hp</b>	Surrounding hydrophobicity
4	<b>P</b>	Polarity
5	<b>pHi</b>	Isoelectric point
6	<b>pK'</b>	Equilibrium constant with reference to the ionization
7	<b>Mw</b>	Molecular weight
8	<b>BI</b>	Bulkiness

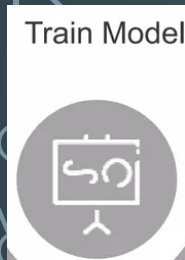




# CLUSTERING

- Nenadzorovana klasifikacija
- K-means gručenje
  - Razdelitev v K gruč preko iterativnega postopka
  - Težave pri različnih velikostih gruč (težko pravilno najti manjše gruče)
  - Kakšen K je pravi?
- Problem računanja razdalj v veliko dimenzijah - drago





# TRENIRANJE MODELA

- Matematičen postopek – navadno neke vrste minimizacija
- Definira se Cost/Loss funkcijo
- Iskanje pravih parametrov modela, tako da se minimizira loss funkcijo
  - Linearna regresija – iskanje parametrov linearne funkcije, tako da je npr. MSE minimalen
  - SVM – iskanje parametrov hiperploskve, tako da je pas, ki ločuje razrede najširši
  - Nevronske mreže - iskanje vrednosti uteži, da je output kar se da podoben predvideni oznaki

Loss function

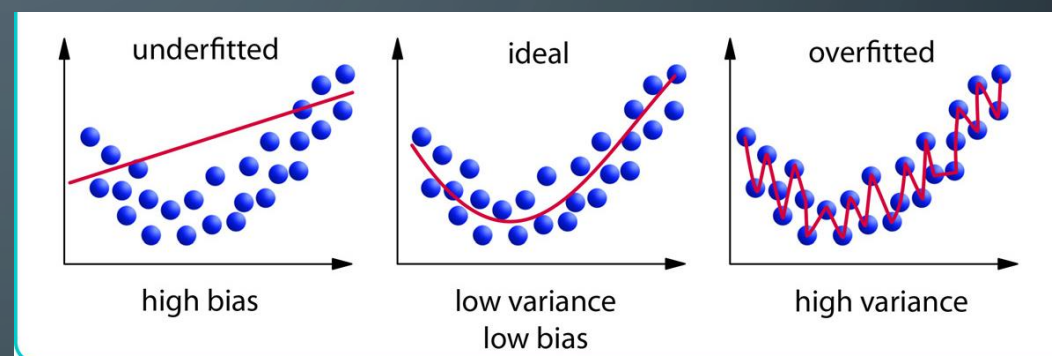
$$\mathcal{L}(\mathbf{W}) = \sum_{i=0}^{N_L} (y_i - \hat{y}_i)^2$$

prediction      label



# TRENIRANJE MODELA

- Over-fitting, under-fitting
  - Problem: slaba napovedna moč
  - V več dimenzijah lahko ni očitno, da smo pretirano prilagodili model
- Rešitve:
  - Uporabimo prav red polinoma za fitanje
  - Early stopping



## avoiding overfitting

### regularization<sup>330</sup>

- constraining model coefficients

### cross-validation<sup>331</sup>

- holding back data to test the model on truly unseen data

### dropout<sup>332</sup>

- ignoring a subset of neurons with a set probability

### ensembling<sup>333</sup>

- combining predictions from multiple separate models

- randomly creating subsets of samples by bootstrap aggregation

### early stopping<sup>334</sup>

- reducing the number of iterations performed to train a model



Test Model

# TESTIRANJE MODELA

- Napovedujemo na še ne-videnih podatkih in gledamo natančnost
  - Več različnih mer za natančnost, odvisno od konteksta (npr. pri iskanju bolnih oseb, si pred vsem želimo, da ne zgrešimo koga, ki je bolan)
- Testni set podatkov (npr. 20% celotnega seta)
- *Confusion matrix*
- Primerjava z drugimi modeli

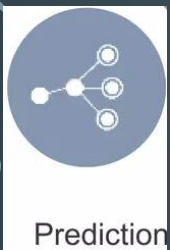
		Expected			
		1	2	3	4
Predicted	1	52	3	7	2
	2	2	28	2	0
	3	5	2	25	12
	4	1	1	9	40





# PRILAGANJE MODELA

- Spreminjanje hiperparametrov
  - Število gruč, arhitektura nevronske mreže, število iteracij, metrika razdalje, inicializacija modela
  - Grid-search
- Spreminjanje vhodnih podatkov - dodatna obdelava, drugačna reprezentacija



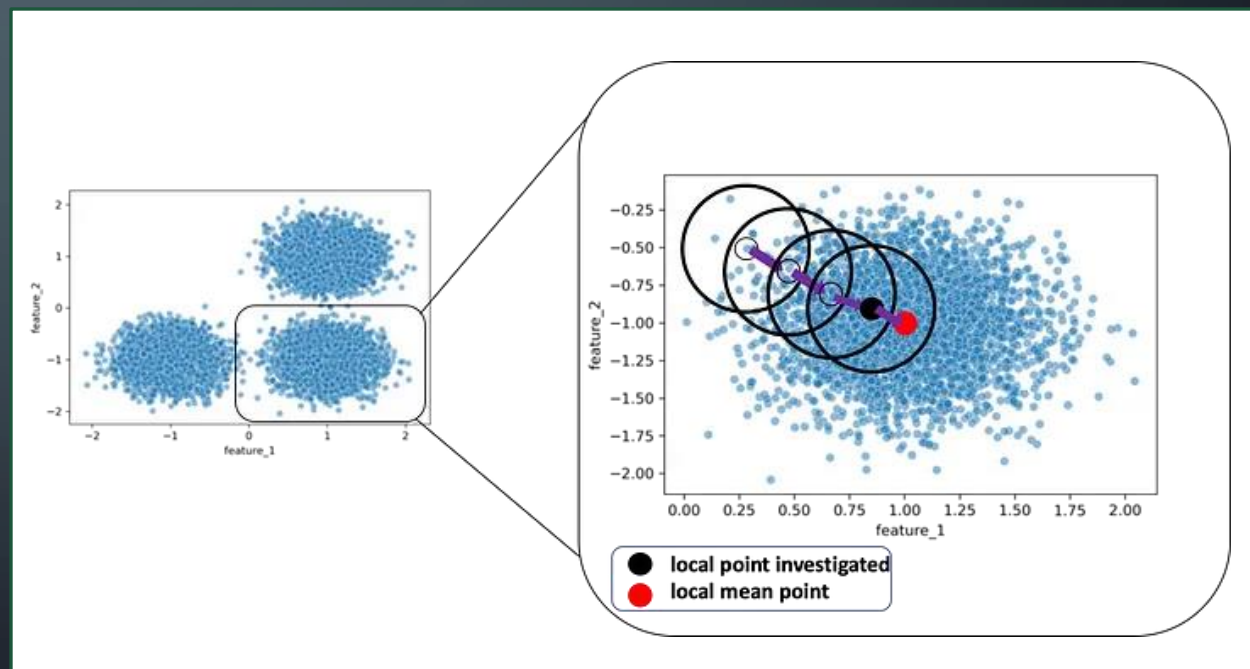
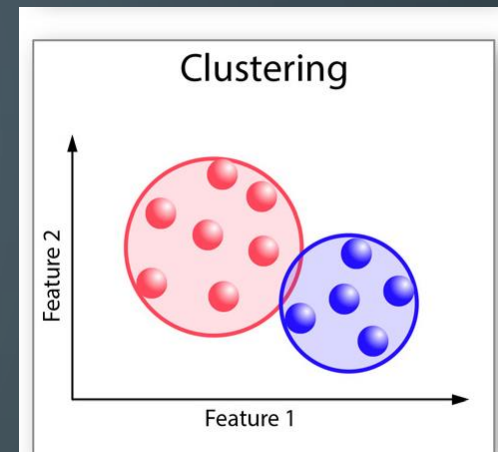
# NAPOVEDOVANJE

- Napoved bo smiselna le na podatkih, ki so podobni vhodnim
- Pogosto problem z interpretacijo – model ne pojasnjuje, na kakšen način je klasificiral
- Dodaten plus, če zna model povedati s kakšno verjetnostjo poda napoved



# CLUSTERING

- Mean-shift gručenje
  - Ne definiramo števila razredov
- Težave z interpretacijo delitve
- Kdaj končati algoritem? Ni trivialno





# TRENIRANJE MODELA

- Stochastic gradient descent
  - Spreminjanje parametrov modela glede na lokalni gradient (v hiperprostoru parametrov)
  - Stohastičnost - ne pregledamo celotnega prostora, pač pa začnemo v naključni točki.

