

K Nearest Neighbours (kNN)

Jelena Jovanovic

jelena.jovanovic@fon.bg.ac.rs

kNN: Osnovne karakteristike

Jednostavan i vrlo intuitivan algoritam

Može se koristiti za zadatke klasifikacije i regresije

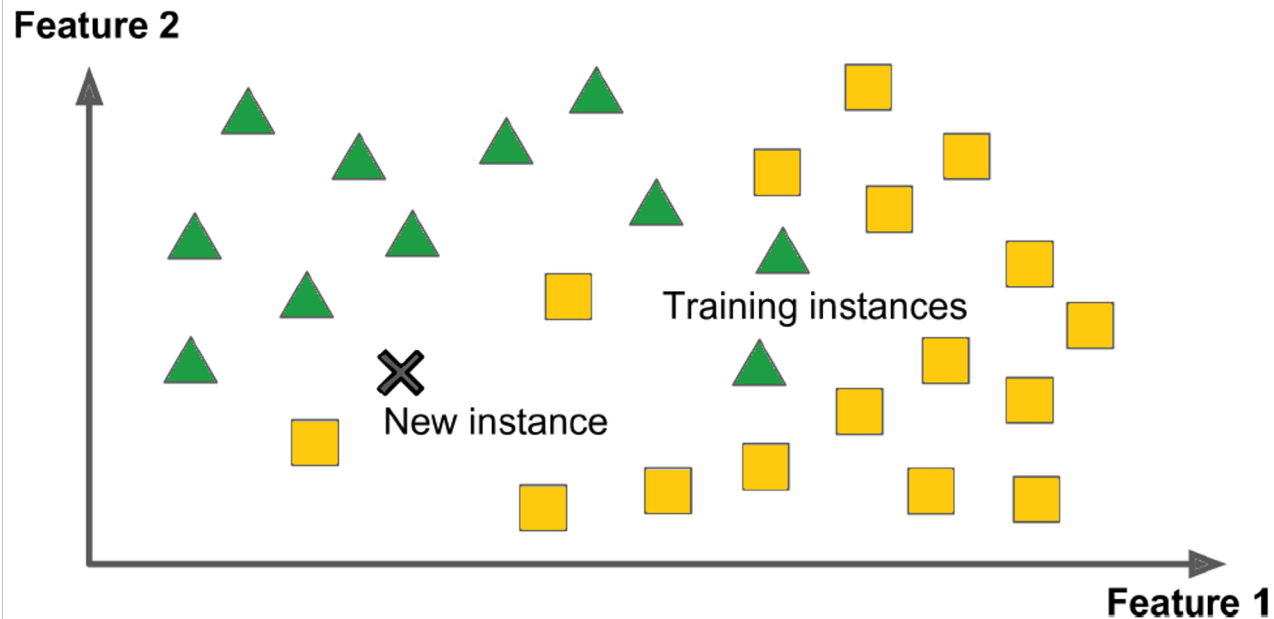
Za razliku od većine ML algoritama, kNN nema trening fazu i ne kreira model

kNN: Osnovne karakteristike

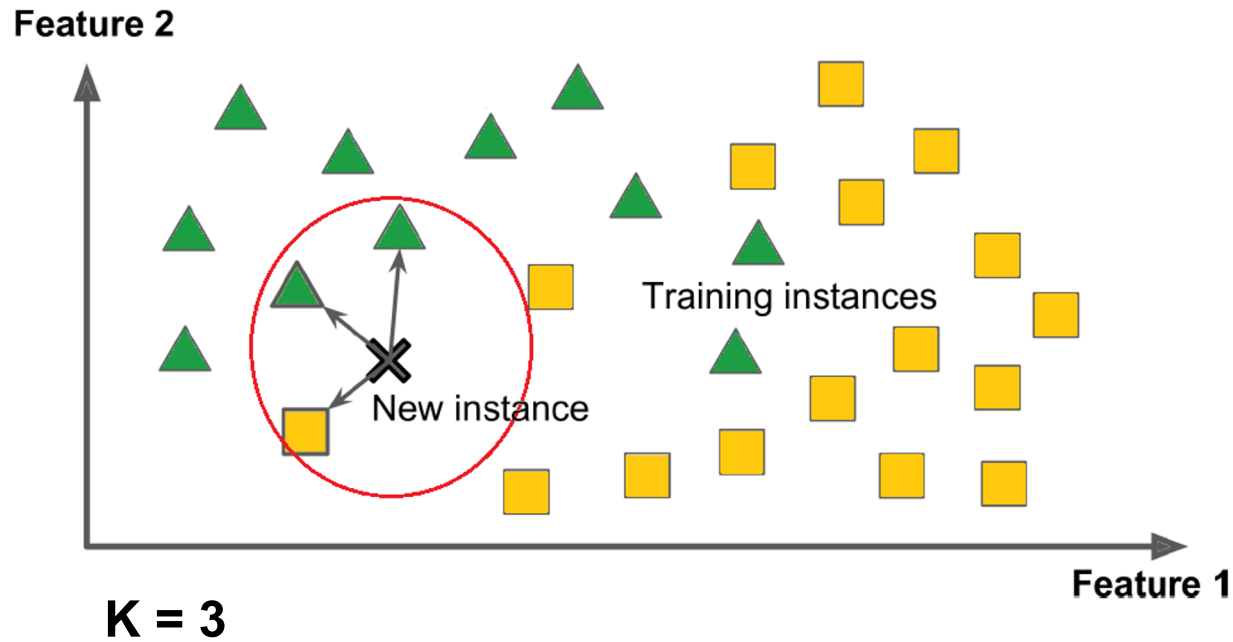
Predviđa izlaznu vrednost za novu instancu u dva koraka:

- 1) Pronađe k najsličnijih (tj., najmanje udaljenih) instanci u skupu za trening - tzv. k najbližih suseda

Kako kNN radi?



Kako kNN radi?

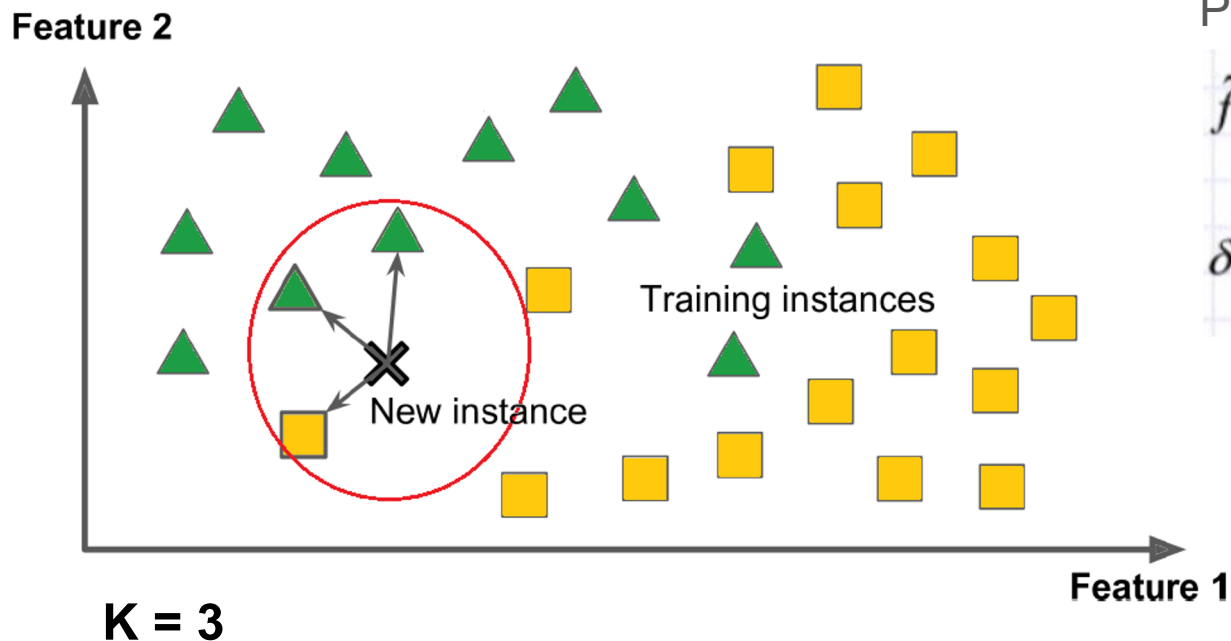


kNN: Osnovne karakteristike

Predviđa izlaznu vrednost za novu instancu u dva koraka:

- 1) Pronađe k najsličnijih (tj., najmanje udaljenih) instanci u skupu za trening - tzv. k najbližih suseda
- 2) Na osnovu izlazne vrednosti k najbližih suseda, određuje izlaznu vrednost nove instance i to:
 - a) Kao dominantnu klasu među k najbližih suseda, u slučaju klasifikacije
 - b) Kao prosečnu vrednost izlazne varijable k najbližih suseda, u slučaju regresije

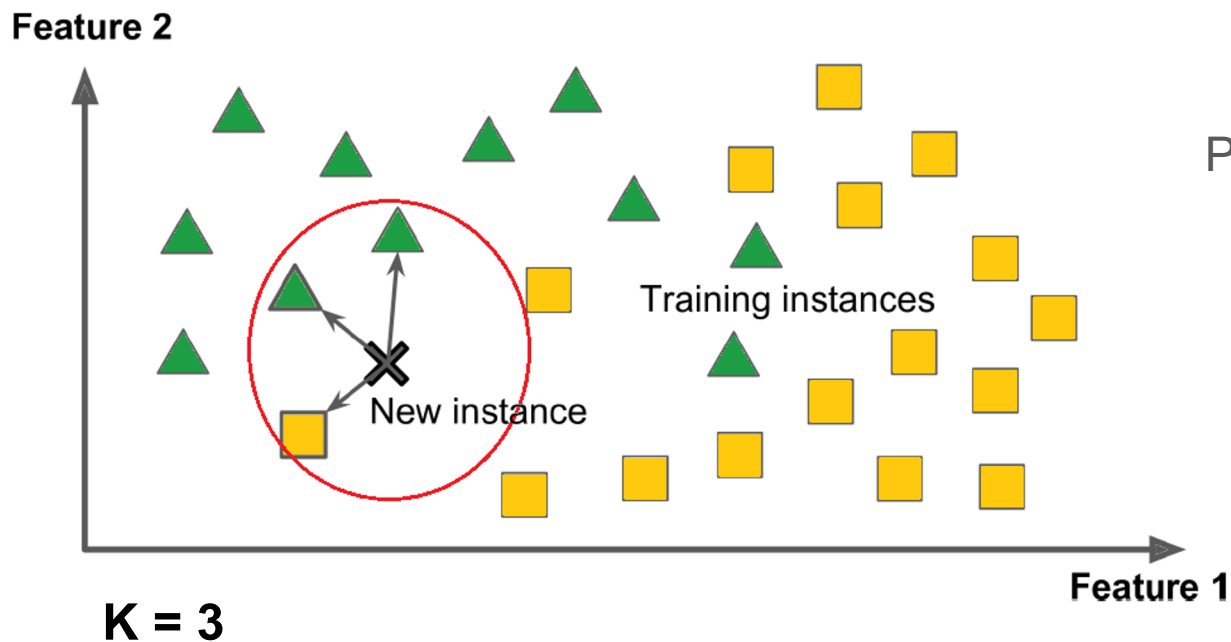
Klasifikacija primenom kNN-a



Predviđena izlazna vrednost:

$$\hat{f}(x_q) \leftarrow \arg \max_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$
$$\delta(a, b) = \begin{cases} 1 & a = b \\ 0 & \text{otherwise} \end{cases}$$

Regresija primenom kNN-a

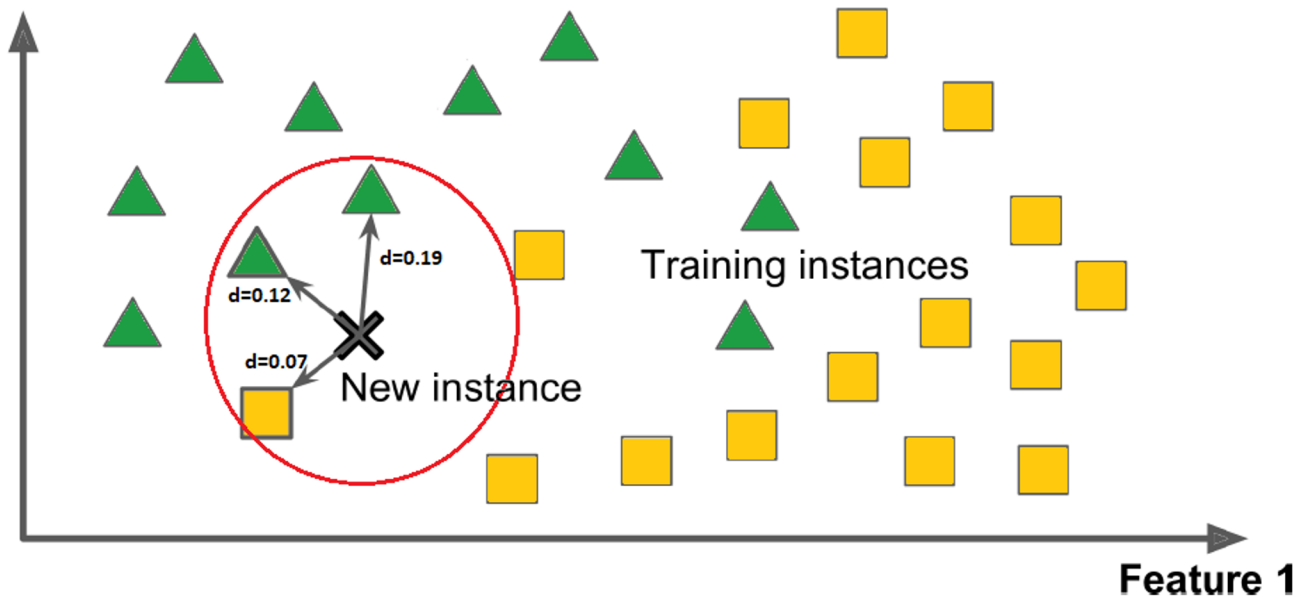


Predviđena izlazna vrednost:

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

Težinski kNN

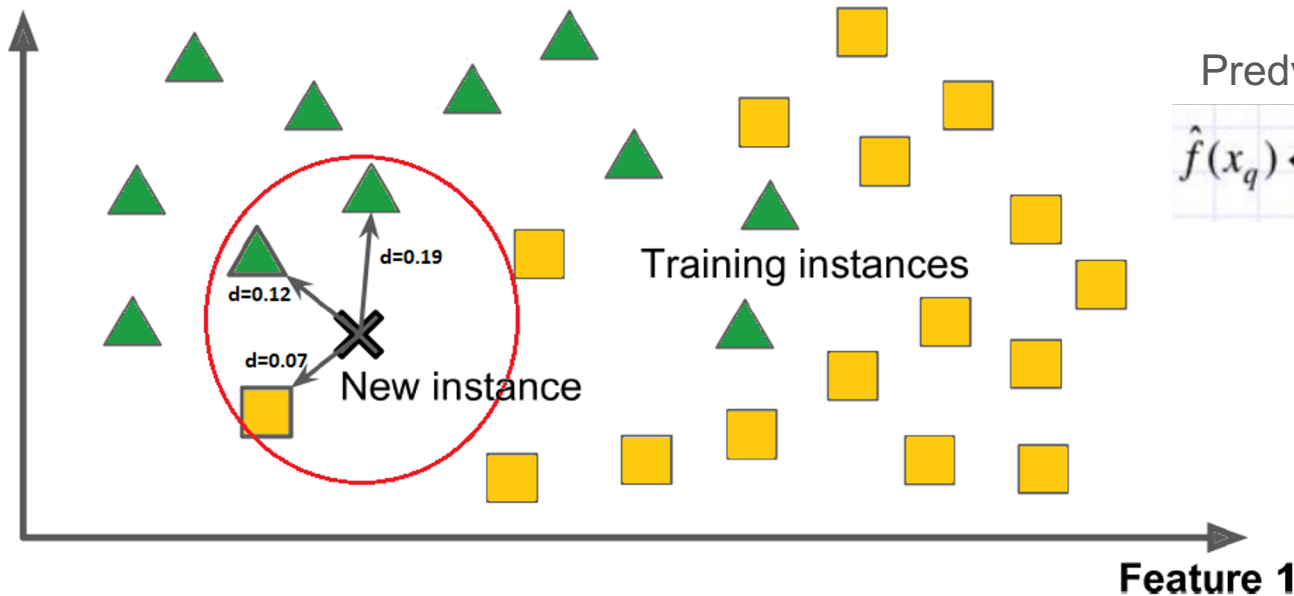
Feature 2



$K = 3$

Težinski kNN: klasifikacija

Feature 2



K = 3

Težina (značaj) “suseda”:

$$w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

Predviđena izlazna vrednost:

$$\hat{f}(x_q) \leftarrow \arg \max_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i))$$

Težinski kNN: regresija

Težina (značaj) “suseda”:

$$w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

Predviđena izlazna vrednost:

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}$$

Feature 2

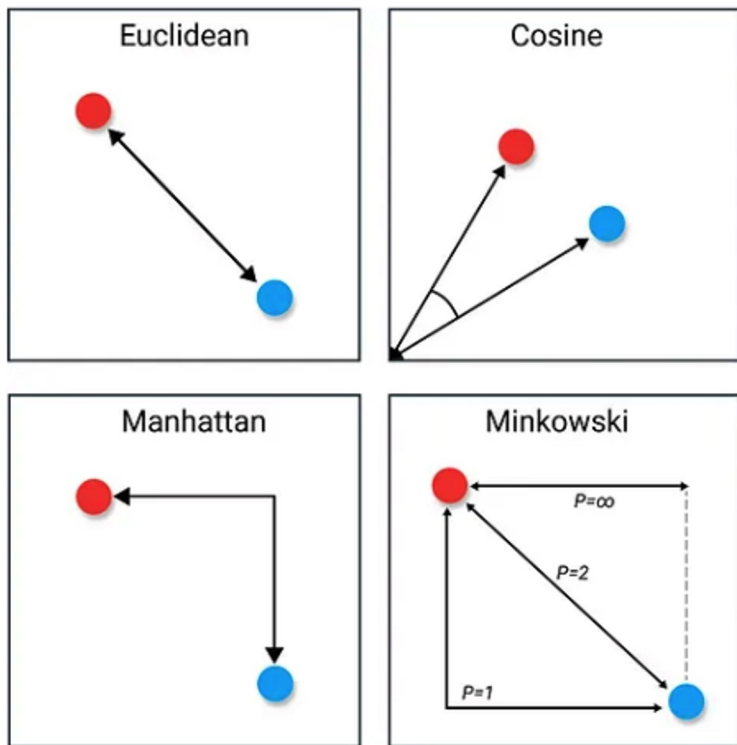


kNN: Hiper-parametri

Rezultati kNN-a zavise od 2 ključne stvari (hiper-parametara modela):

- Broja suseda (k)
- Izabrane mere udaljenosti / sličnosti

Mere udaljenosti / sličnosti



Euclidean:
$$d = \sqrt{\sum_{j=1}^n (x_{sj} - x_{tj})^2}$$

Manhattan:
$$d = \sum_{j=1}^n |x_{sj} - x_{tj}|$$

Minkowski:
$$d = \left(\sum_{j=1}^n |x_{sj} - x_{tj}|^p \right)^{1/p}$$

Cosine:
$$d = 1 - \frac{\sum_{j=1}^n x_{sj} x_{tj}}{\sqrt{\sum_{j=1}^n x_{sj}^2} \sqrt{\sum_{j=1}^n x_{tj}^2}}$$

Uticaj različitih vrednosti k na predviđanje

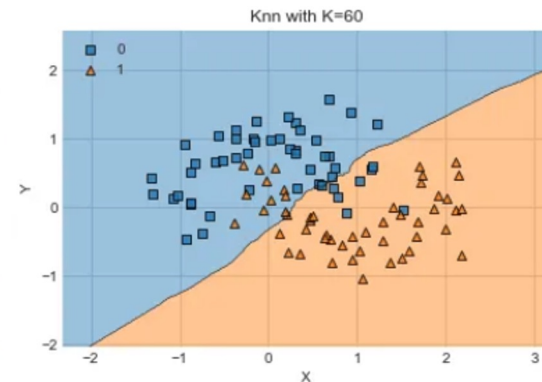
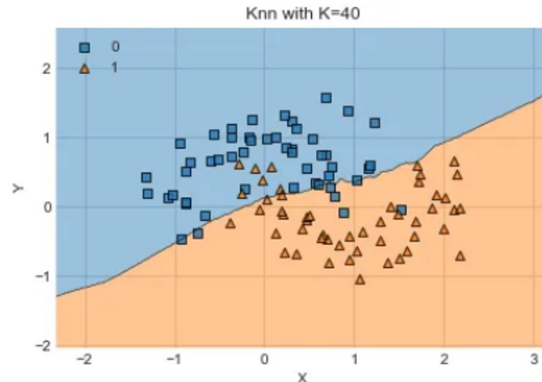
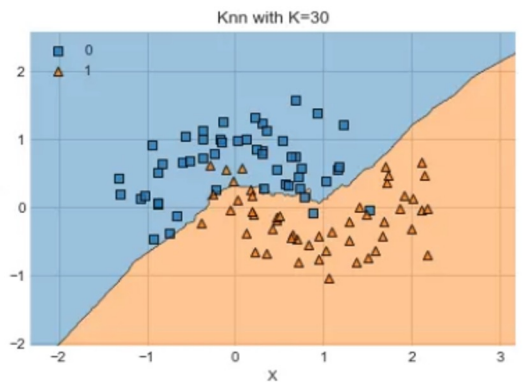
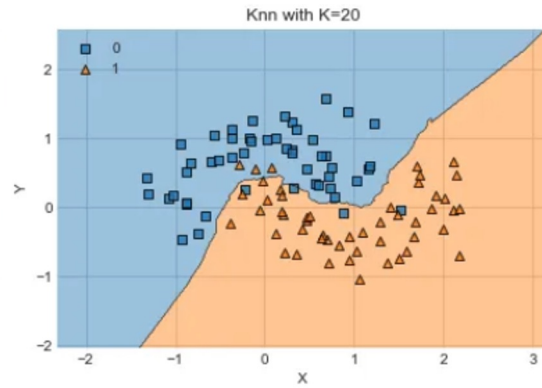
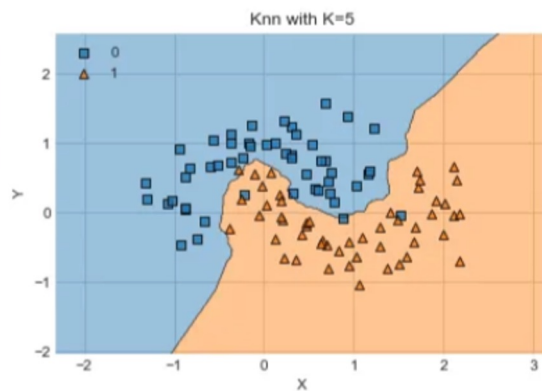
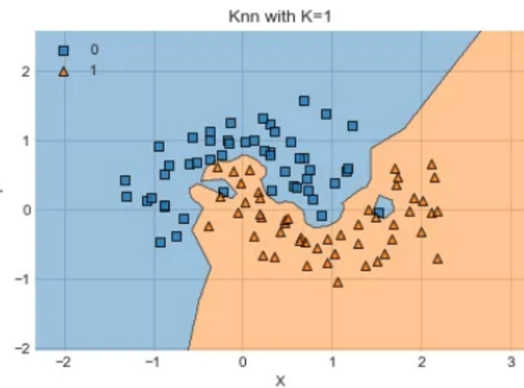
Veoma mala vrednost za K

- *over-fitting*-a i visoke osetljivosti na promene u dataset-u (*high variance*)
- osetljivost na outliers

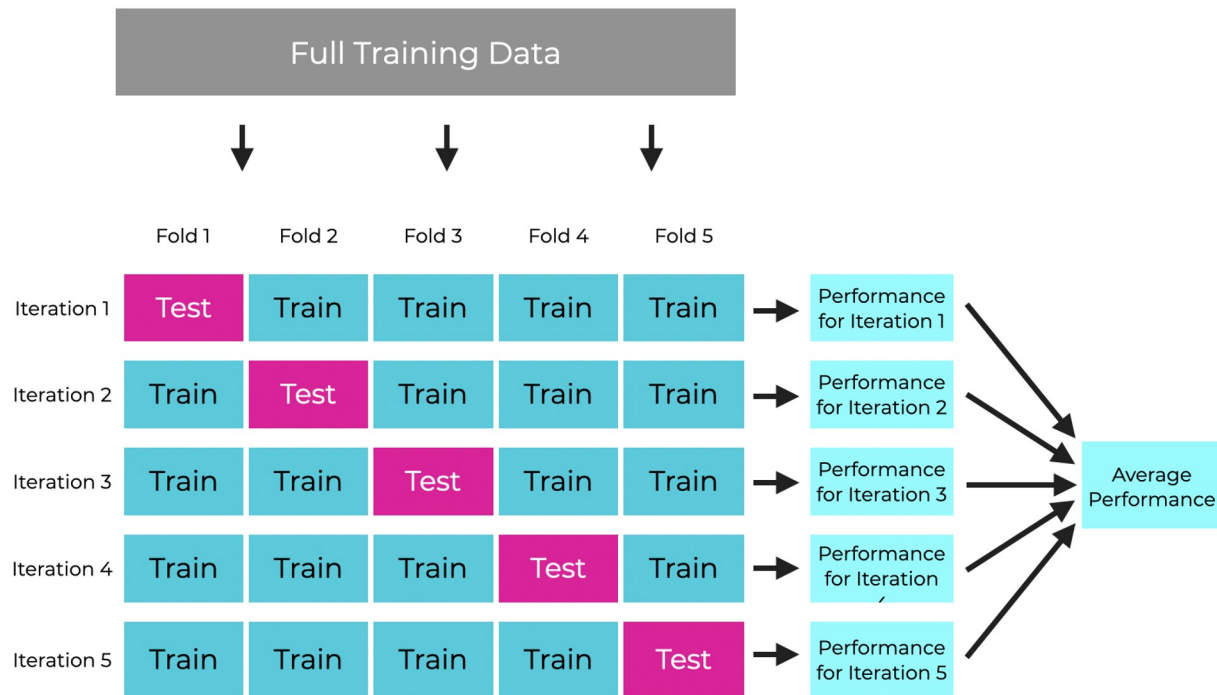
Suviše velika vrednost za K

- *under-fitting*
- pristrasnost prema većinskoj klasi

Uticaj različitih vrednosti k na predviđanje



Određivanje k primenom kros-validacije



kNN: Prednosti

Jednostavnost: *kNN ima malo pretpostavki i zahteva minimalno podešavanje hiper-parametara (svega dva)*

Neparametrijski: *ne zahteva neku posebnu distribuciju podataka*

Fleksibilnost: *može se koristiti za zadatke klasifikacije i regresije*

Interpretabilnost: *predviđanja se mogu lako protumačiti*

Bez faze treninga: *skup podataka za obuku se koristi direktno tokom faze predviđanja*

kNN: Mane / Izazovi

Slabo skaliranje: kako veličina skupa podataka raste, vreme predviđanja se povećava

Osetljivost na opseg vrednosti atributa: atributi sa većim opsegom vrednosti mogu da dominiraju u proračunu udaljenosti

Neadekvatan u slučaju nebalansiranih podataka: većinska klasa ima tendenciju da dominira predviđanjima, što čini izazovom identifikovanje slučajeva manjinske klase

Osetljivost na vrednost k