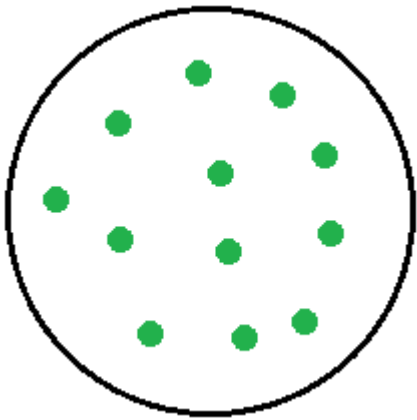


# Bagging

---

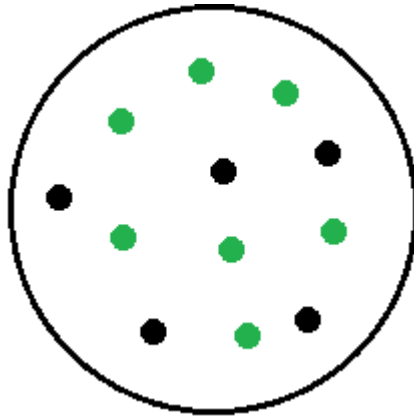
# Manipulacija trening skupom

- Na *različitim* podacima više puta trenirati *isti* model
- Kako da dobijemo različite podatke (u praksi, raspolažemo sa samo jednim skupom podataka)?



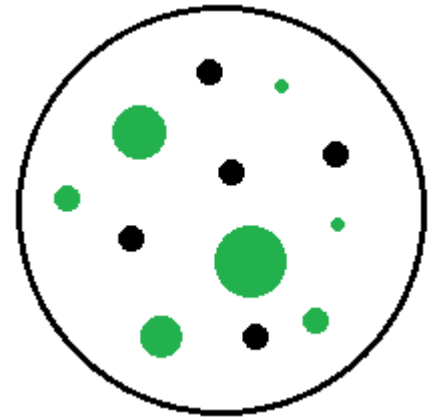
Jedan model

koristiti ceo skup podataka



Bagging

Slučajno uzorkovanje sa  
zamenom



Boosting

Slučajno uzorkovanje sa  
zamenom, pri čemu neki  
primeri imaju veću šansu  
da budu selektovani

\* Zelene tačke su selektovane za treniranje modela

# Bagging (*Bootstrap Aggregation*)

<b>Ulaz</b>	<ul style="list-style-type: none"><li>• <math>T = \{(x^{(i)}, y^{(i)}), i = 1, \dots, N\}</math> – trening skup</li><li>• <math>M</math> – broj slabih prediktora</li><li>• Osnovni model <math>L</math> (korišćen za obuku pojedinačnih članova ansambla)</li></ul>
<b>Postupak</b>	<ul style="list-style-type: none"><li>• for <math>m = 1, \dots, M</math><ul style="list-style-type: none"><li>• Napraviti podskup skupa podataka <math>T_m</math> uzorkovanjem <math>N</math> primera sa zamenom iz <math>T</math></li><li>• Obučiti model <math>h_m = L(T_m)</math></li></ul></li></ul>
<b>Izlaz</b>	$\hat{y}(x) = \frac{1}{M} \sum_{m=1}^M h_m(x)$

# Različitost trening skupova

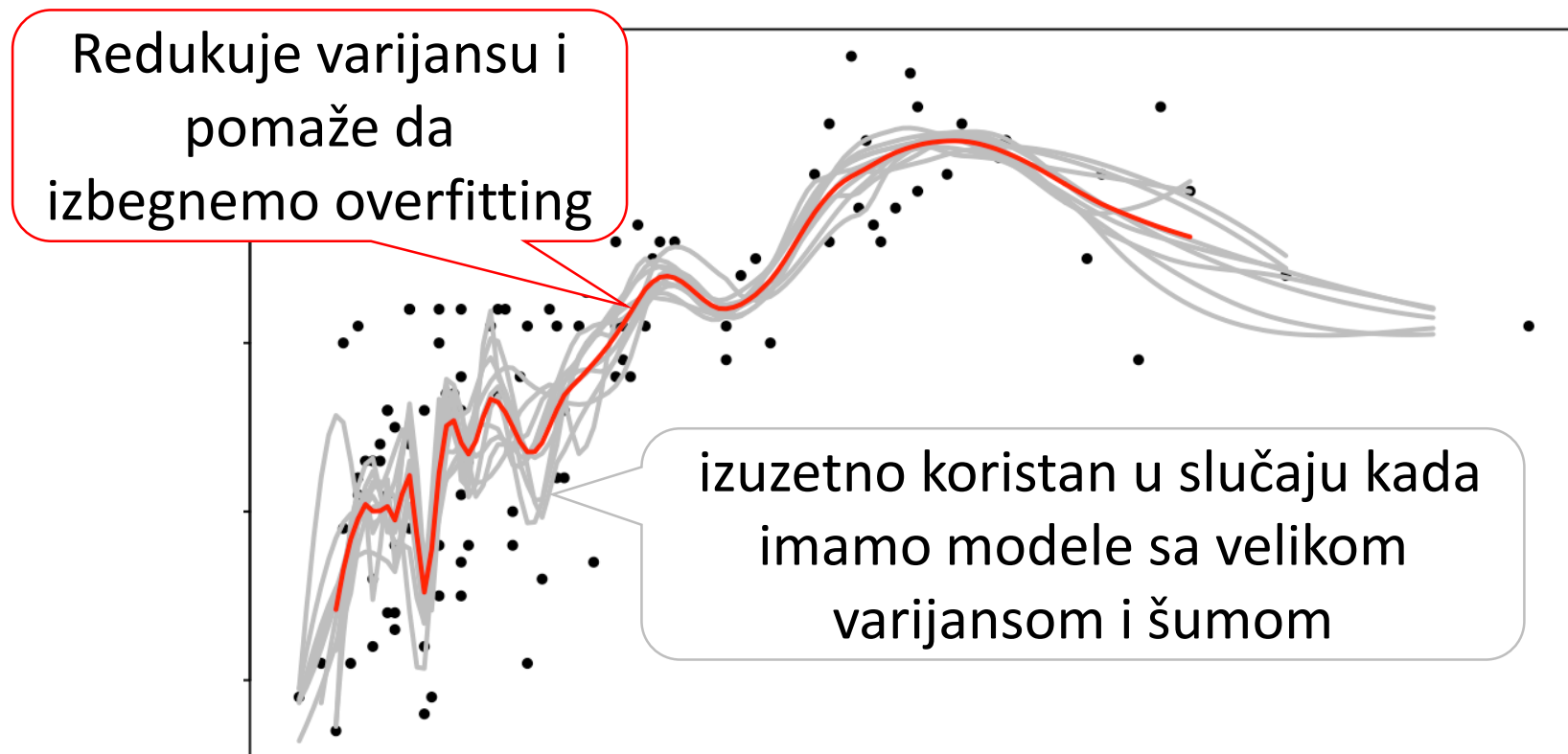
- Svaki podskup  $T_m$  ima isti broj primera  $N$  kao i originalni skup podataka  $T$
- Neki primeri iz  $T$  se ponavljaju u  $T_m$ , a neki su izostavljeni
- Podskupovi se u razumnoj meri razlikuju
  - Instanca ima verovatnoću  $\left(1 - \frac{1}{N}\right)^N \approx e^{-1} \approx 0.368$  da ne bude odabrana za trening skup
  - To znači da će u  $T_m$  biti oko 63.2% originalnih/jedinstvenih instanci

# Random forest

- Najpopularniji algoritam je *Random forests*
- Kao osnovni model se koristi stablo odluke (*decision tree*)
- Umesto podskupova primera, koriste se podskupovi obeležja (svaki model je zasnovan na  $\sqrt{\text{no. features}}$  obeležja)

# Bagging prednosti

- Iznenadjujuće kompetitivne performanse



- Manje osetljiv na nerelevantna obeležja