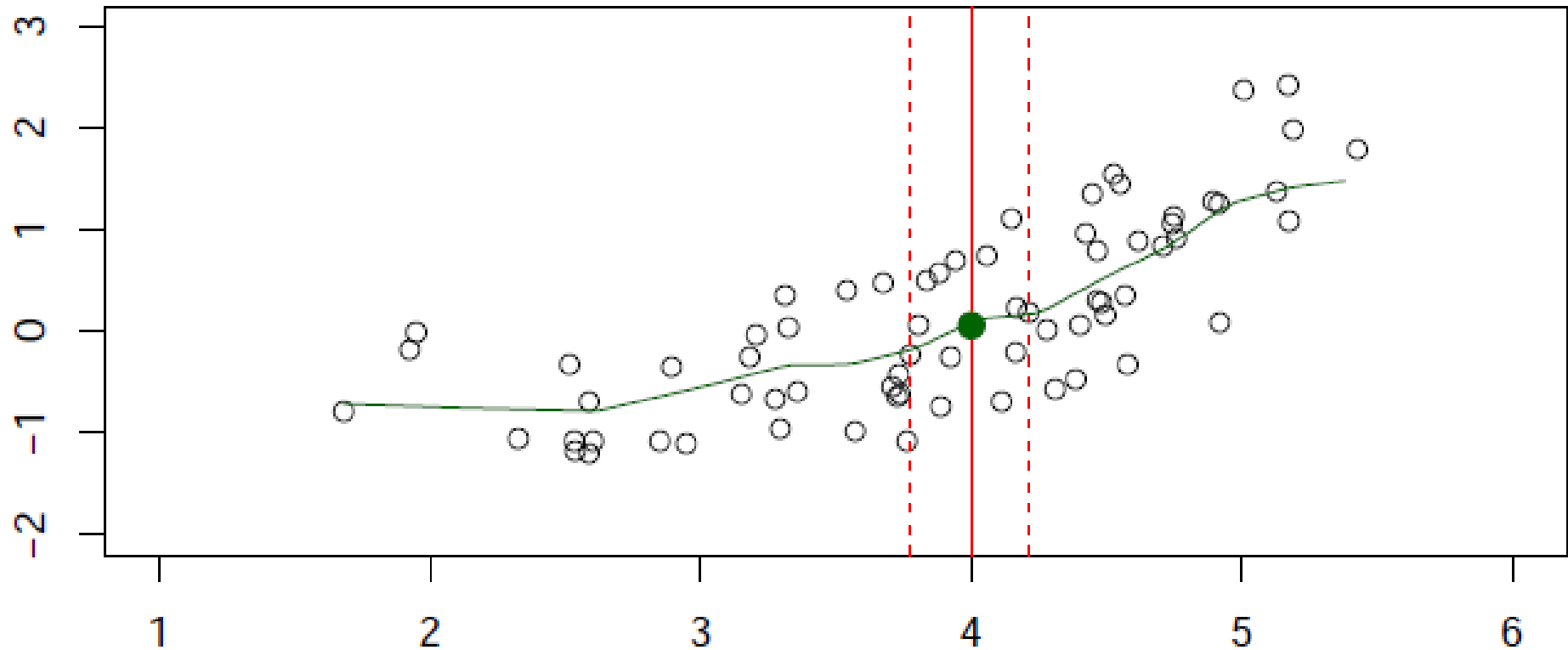
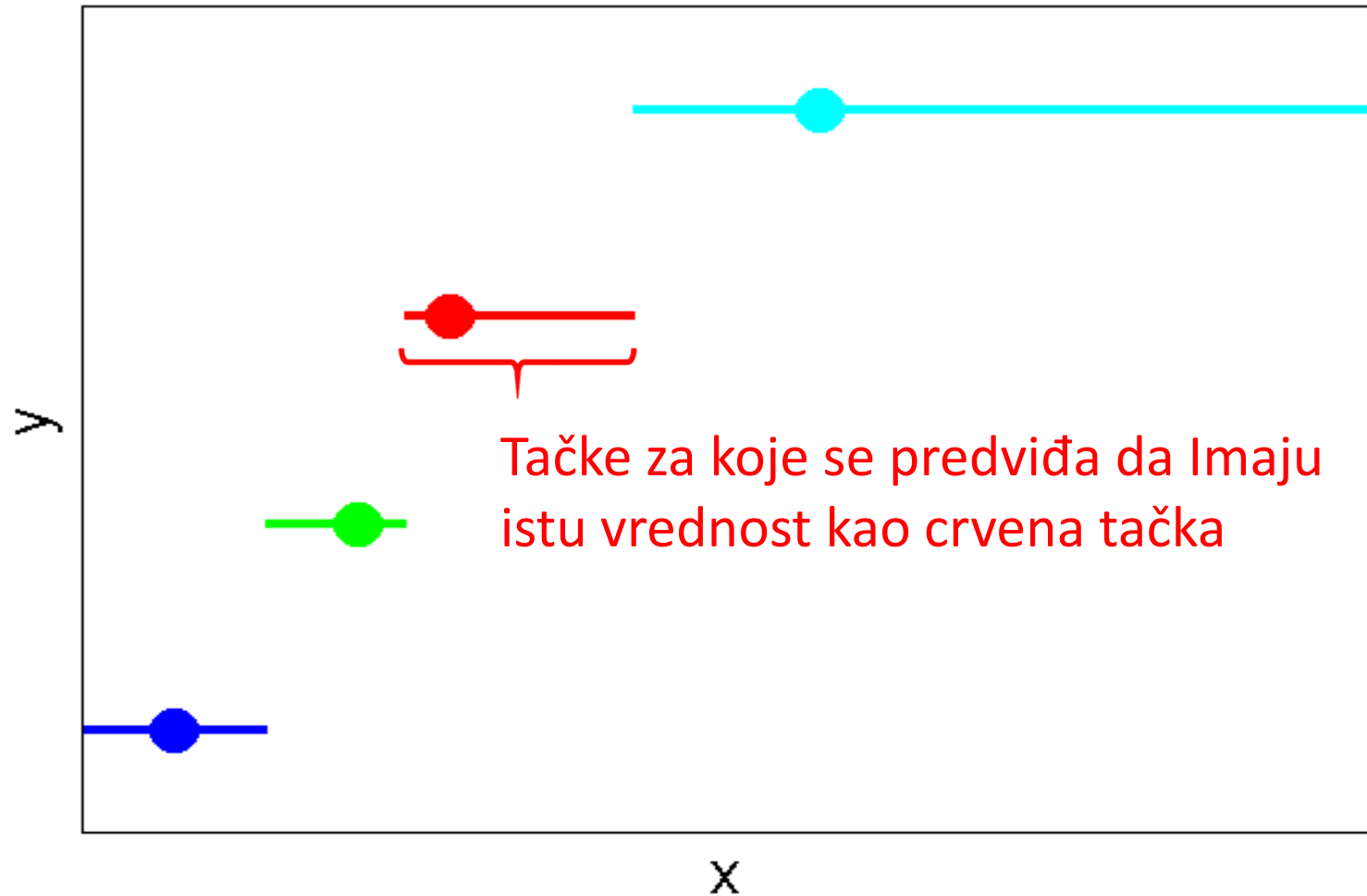


# Najbliži susedi (lokalno uprosečavanje)

Ovo se zove *nearest neighbors* ili *local averaging*



# 1-NN regresija



# 1-NN regresija

Ulaz	<ul style="list-style-type: none"><li>• <math>T = \{(x^{(i)}, y^{(i)}), i = 1, \dots, N\}</math></li><li>• <math>x^{(q)}</math>: tačka za koju treba da odredimo <math>y^{(q)}</math> vrednost</li></ul>
Postupak	<p>Inicijalizovati <math>Dist2NN = \infty, \hat{y}^{(q)} = NaN</math></p> <p>for <math>i = 1, 2, \dots, N</math></p> <p style="padding-left: 40px;"><math>\delta = distance(x^{(i)}, x^{(q)})</math></p> <p style="padding-left: 40px;">if <math>\delta &lt; Dist2NN</math></p> <p style="padding-left: 80px;"><math>Dist2NN = \delta</math></p> <p style="padding-left: 80px;"><math>\hat{y}^{(q)} = y^{(i)}</math></p>
Izlaz	$\hat{y}^{(q)}$

# 1-NN regresija

Euclidean distance



*Voronoi*-ev dijagram:  
vizualizacija 1-NN u više  
dimenzija

Ne moramo eksplicitno  
formirati regije, dovoljna  
nam je definicija  
udaljenosti

Svaka regija sadrži tačno  
jednu tačku  $x^{(i)}$

svakoj tački regije „najbliža“  
tačka  $x^{(i)}$

# Metrika udaljenosti

- Euklidska udaljenost  $distance(x^{(j)}, x^{(q)}) = |x^{(j)} - x^{(q)}|$
- U višedimenzionom prostoru možemo pripisati težine dimenzijama (neke varijable su bitnije)
  - Npr. u predikciji životnog veka , vakcinacija i GDP su važnije od stepena kriminala

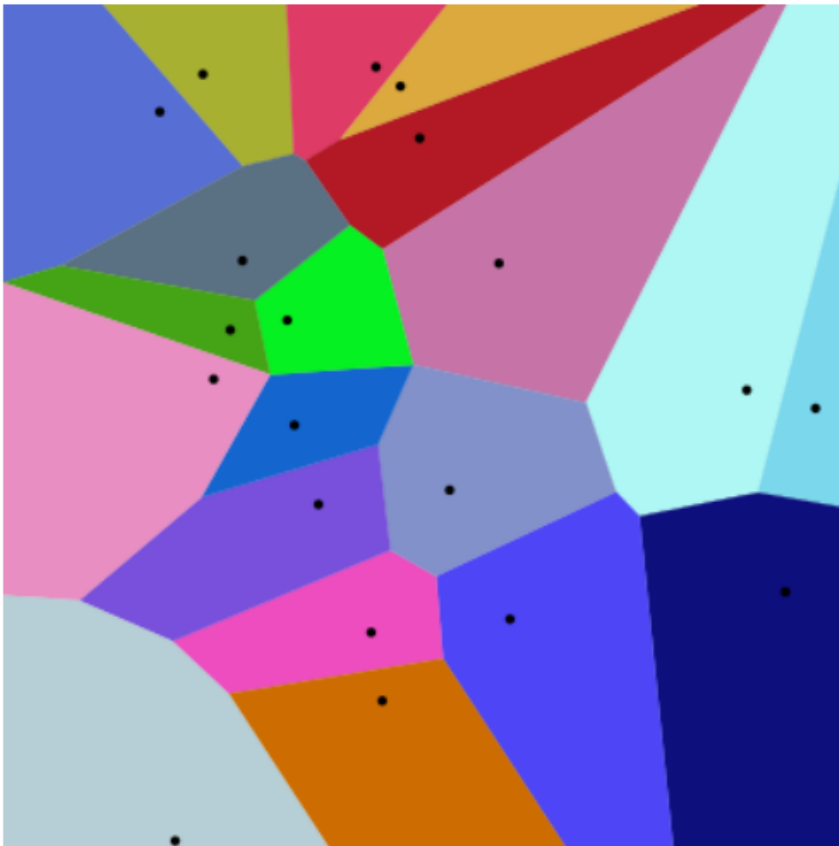
$$distance(x^{(j)}, x^{(q)}) = \sqrt{w_1(x_1^{(j)} - x_1^{(q)})^2 + \dots + w_d(x_d^{(j)} - x_d^{(q)})^2}$$

- Druge metrike udaljenosti: Mahalanobis, rank-based, correlation-based, cosine similarity, Manhattan, Hamming,...

# Metrika udaljenosti

Različite metrike udaljenosti rezultuju različitim prediktivnim površinama

Euclidean distance



Manhattan distance

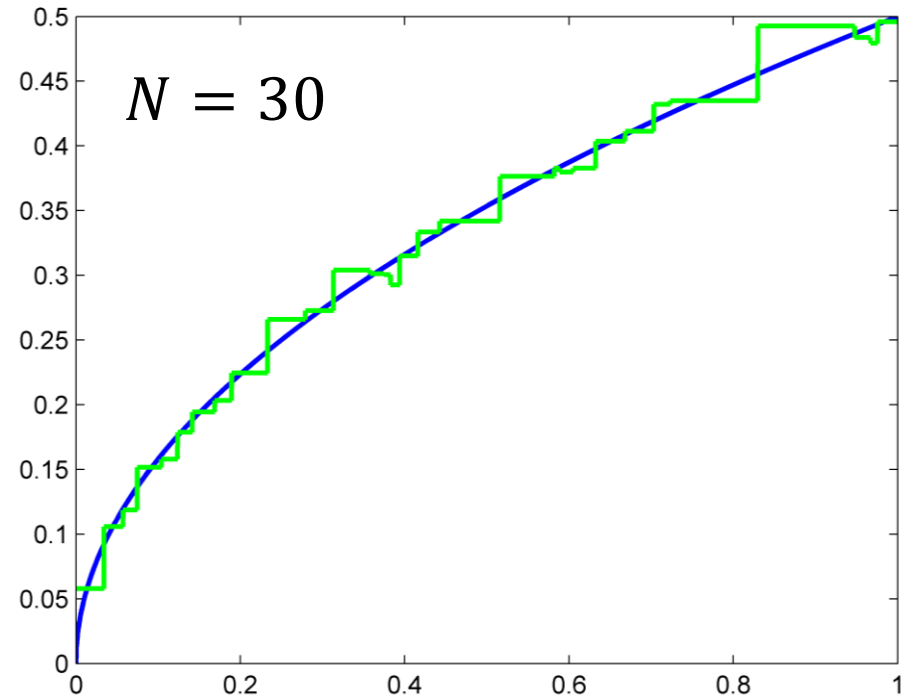
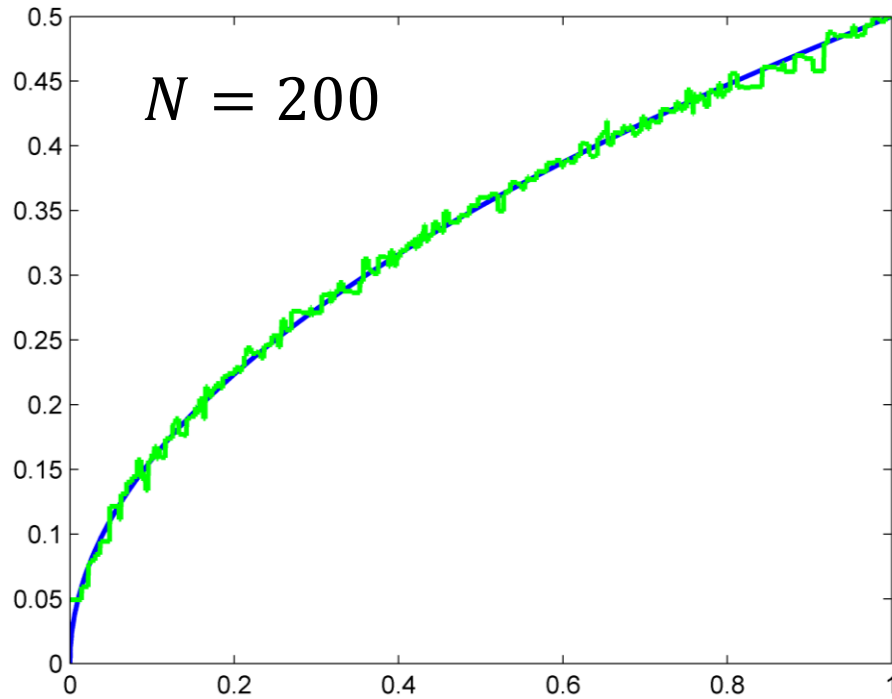


# 1-NN regresija u praksi



1-NN regresija ima dobre performanse ako je  $N$  veliko, a šum mali

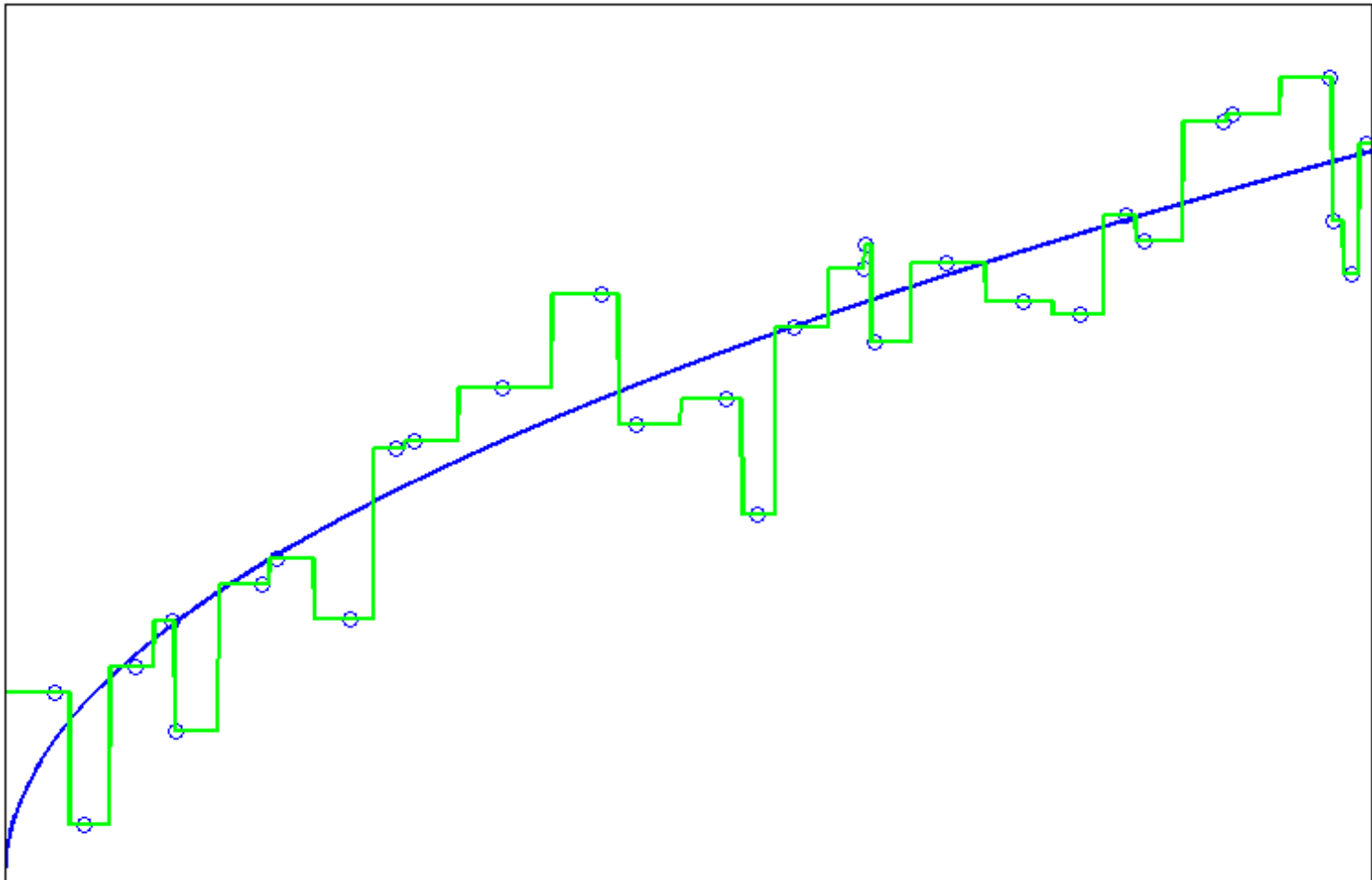
# 1-NN regresija u praksi



Šum: 0.05%

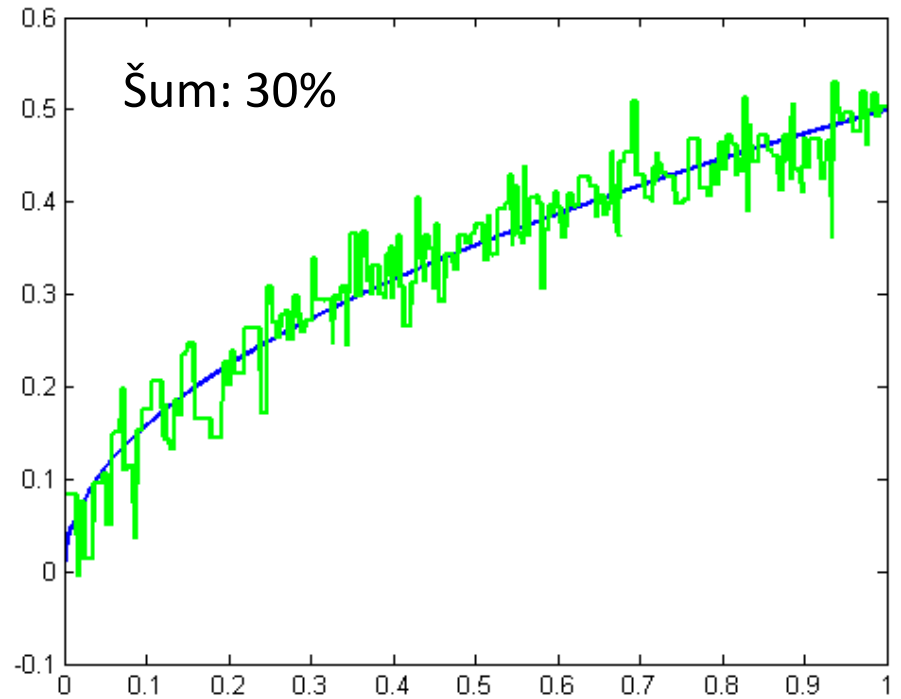
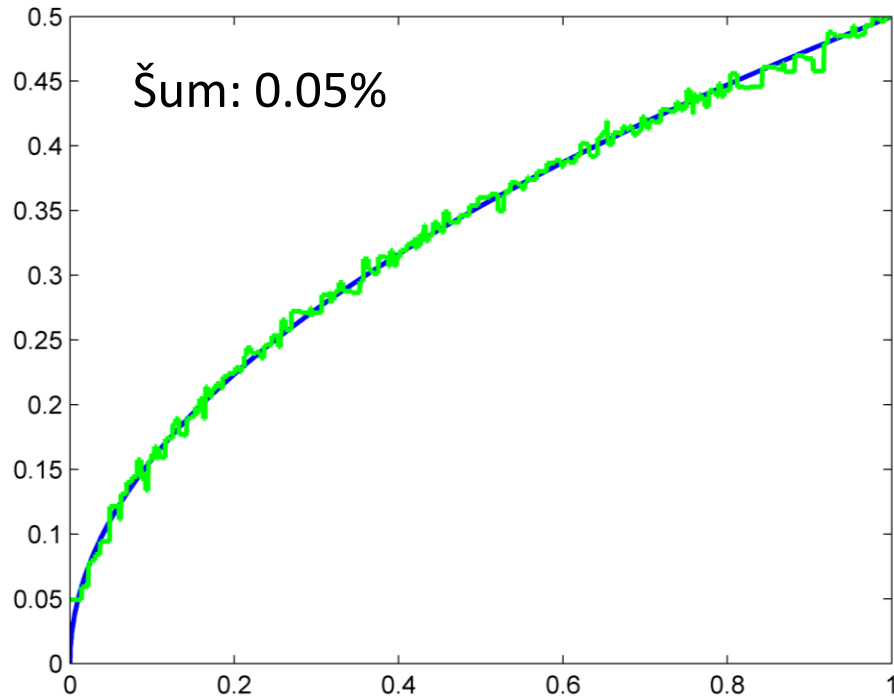


# 1-NN regresija u praksi



1-NN metod je osjetljiv na šum u podacima – *overfitting*

# 1-NN regresija u praksi



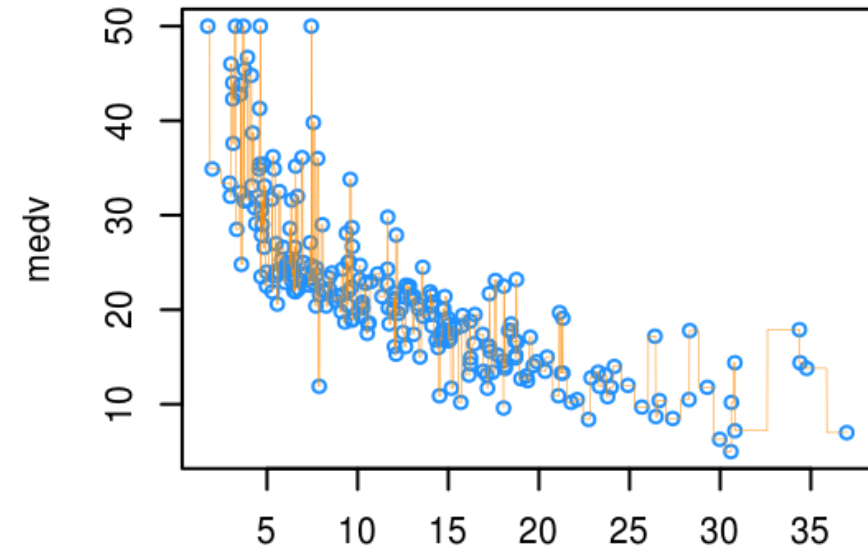
$$N = 200$$

# Kako da rešimo preprilagođavanje?

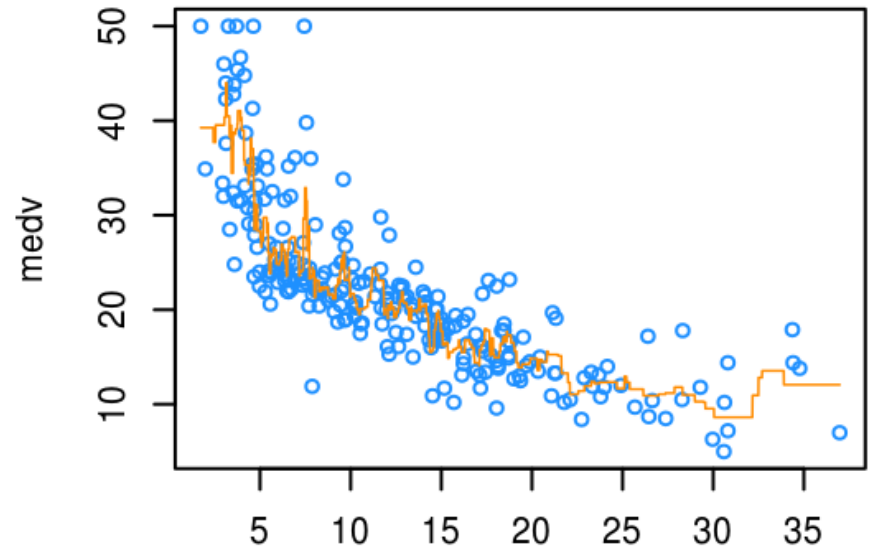
- Kod 1-NN metoda, najbliža tačka koja postoji u trening skupu može imati abnormalno veliku/malu vrednost u odnosu na stvarni trend u podacima
- Umesto da uzmemo samo jednog najbližeg suseda, možemo uzeti više ( $k$ )
- Ovo može rezultovati boljom aproksimacijom

# Uticaj $k$

$k = 1$

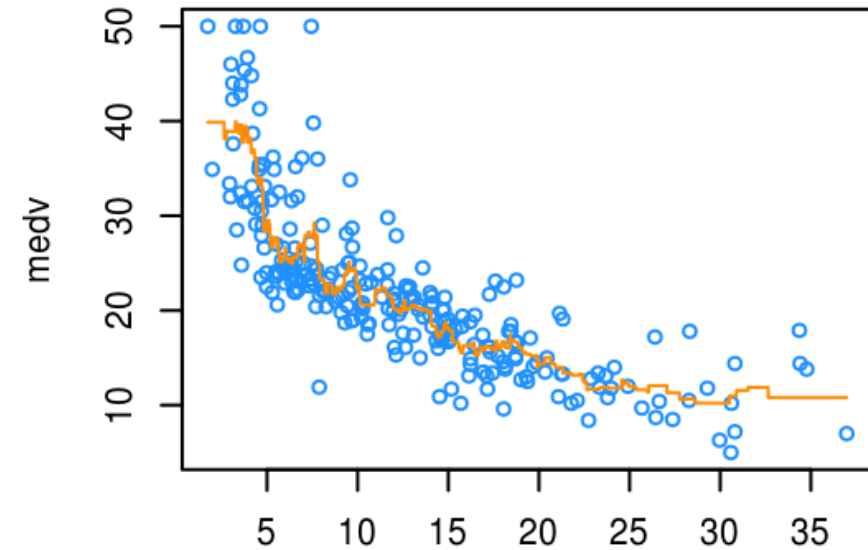


$k = 5$

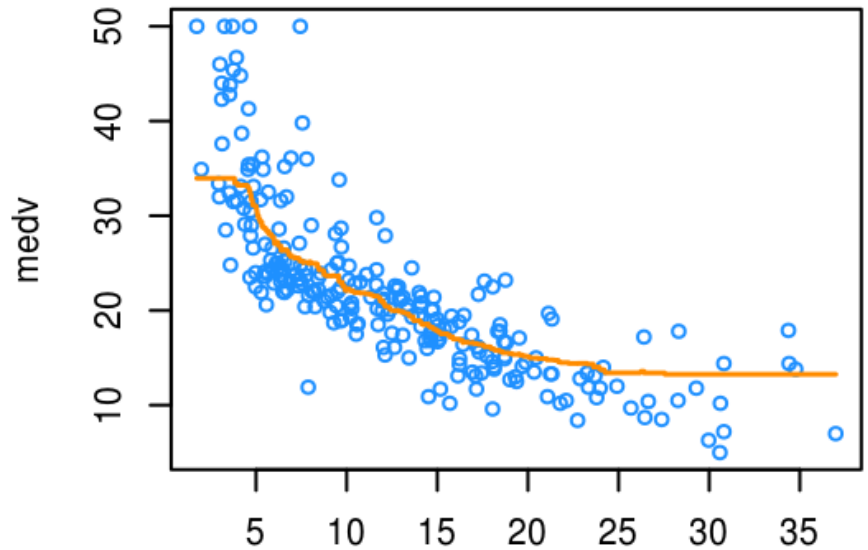


# Uticaj $k$

**$k = 10$**

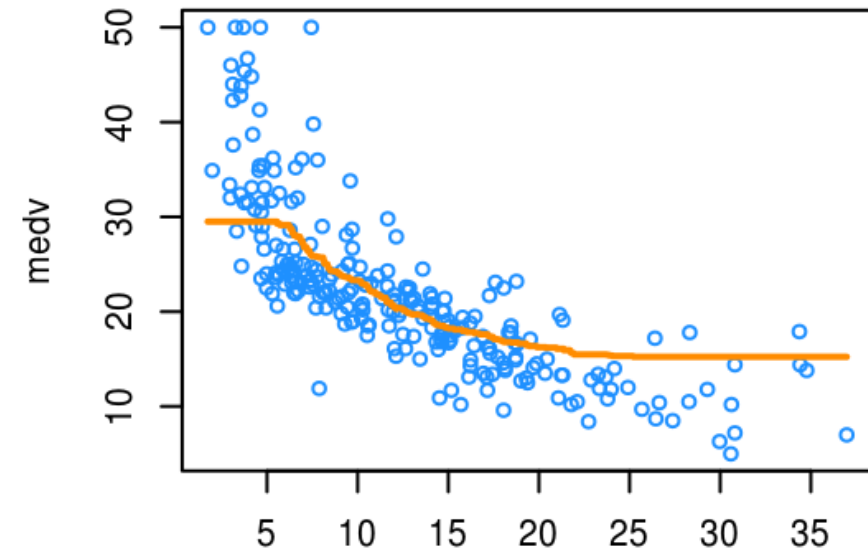


**$k = 25$**

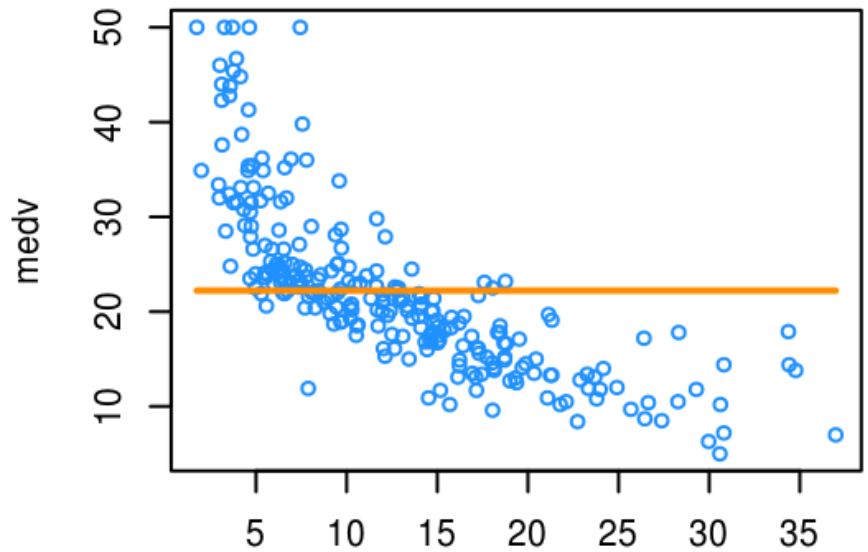


# Uticaj $k$

$k = 50$



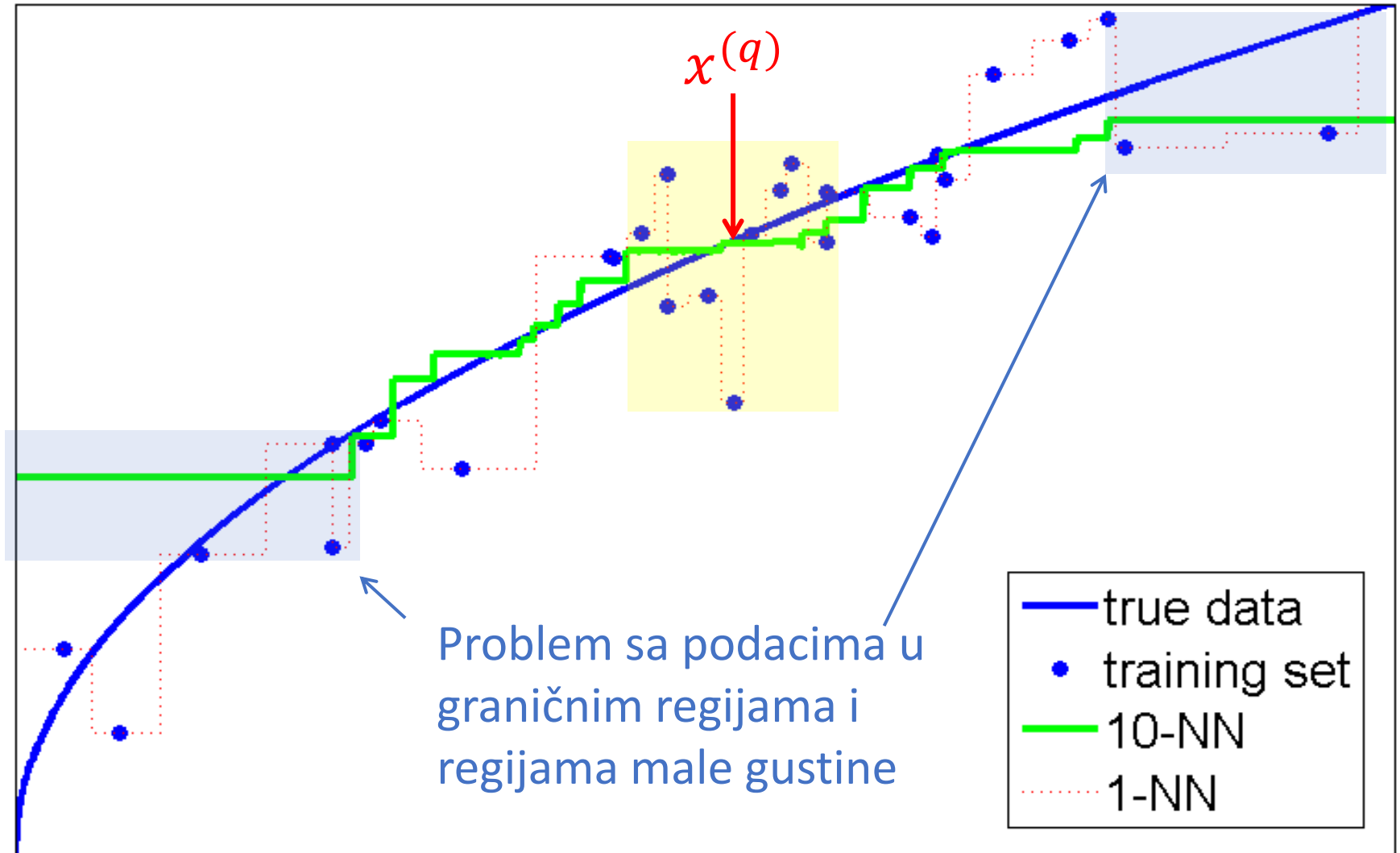
$k = 250$



# $k$ -NN regresija

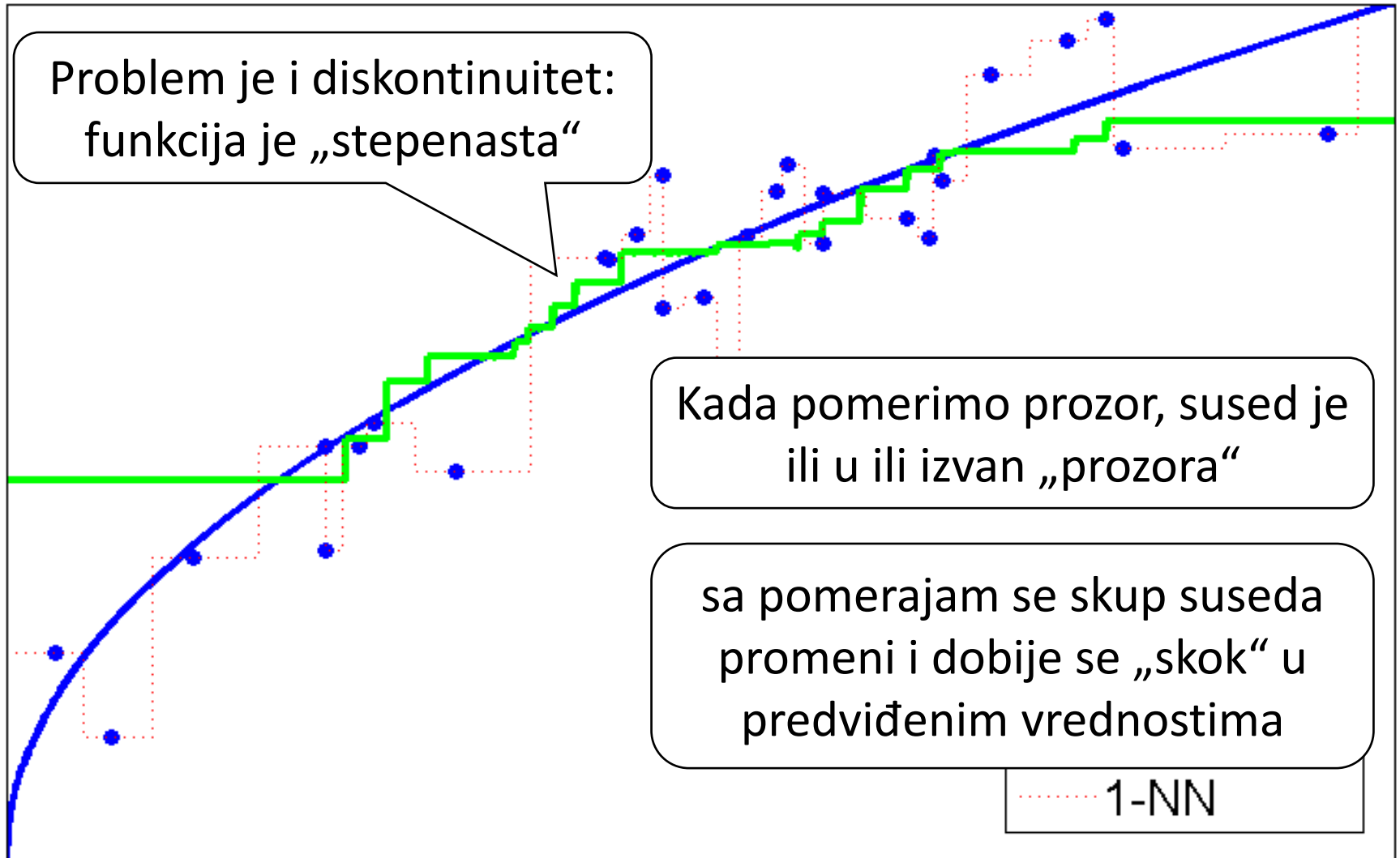
<b>Ulaz</b>	<ul style="list-style-type: none"><li>• <math>T = \{(x^{(i)}, y^{(i)}), i = 1, \dots, N\}</math></li><li>• <math>x^{(q)}</math> – tačka za koju treba da odredimo <math>y^{(q)}</math> vrednost</li></ul>
<b>Postupak</b>	<ol style="list-style-type: none"><li>1. U trening skupu pronaći <math>k</math> tačaka najbližih tački <math>x^{(q)}</math>: <math display="block">\{(x^{(NN1)}, y^{(NN1)}), \dots, (x^{(NNk)}, y^{(NNk)})\}</math></li><li>2. Predvideti <math display="block">\hat{y}^{(q)} = \frac{1}{k} (y^{(NN1)} + \dots + y^{(NNk)})</math></li></ol>
<b>Izlaz</b>	$\hat{y}^{(q)}$

# $k$ -NN regresija





# $k$ -NN regresija



# Problem sa diskontinuitetom

Globalna tačnost nije loša, ali...

