

Strojno učenje na velikih podatkih

Machine learning for big data

Jasna Urbančič

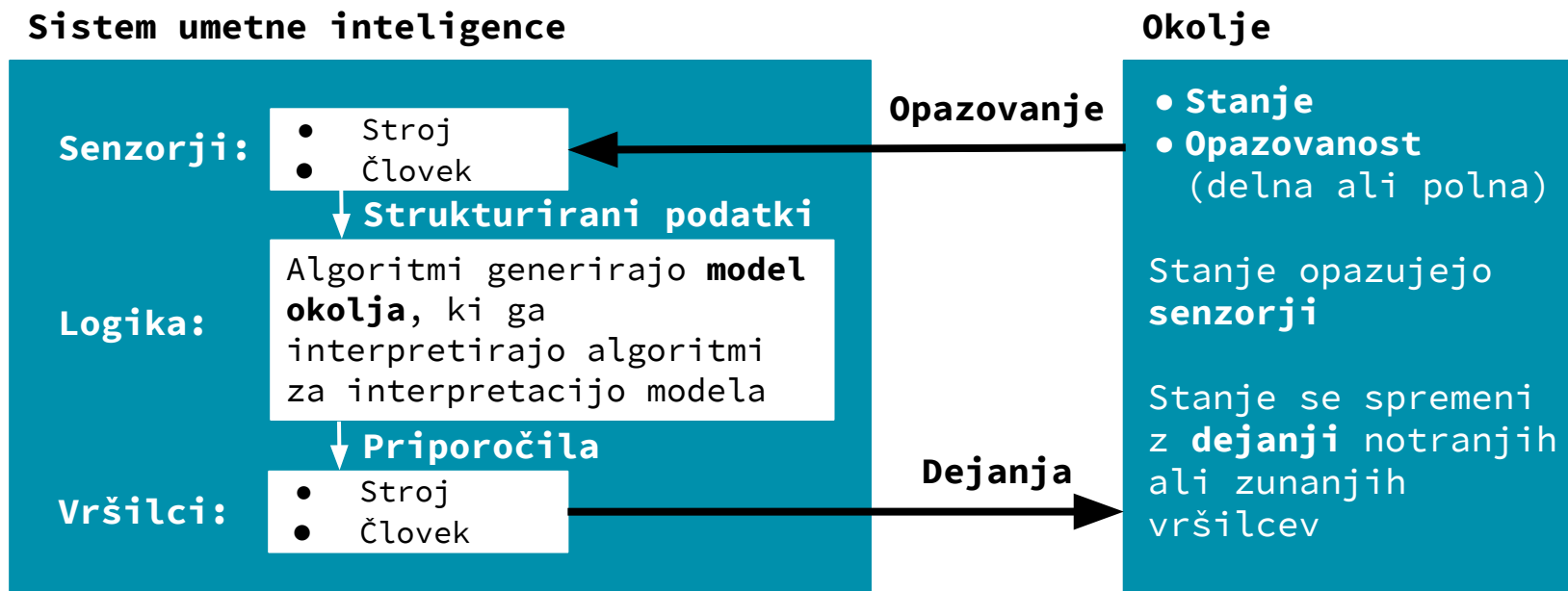
Laboratorij za umetno inteligenco

20. Junij 2019



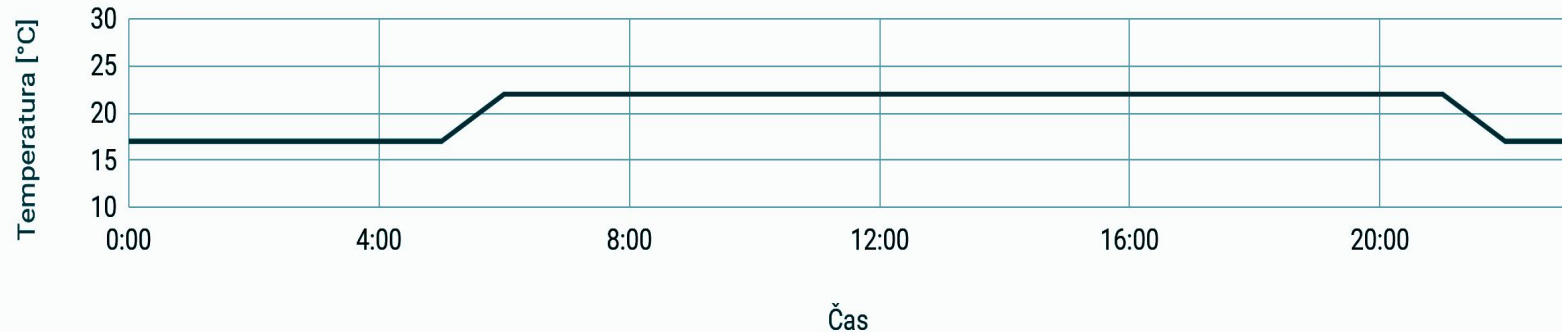
Strojno učenje

Umetna inteligenca in strojno učenje



Umetna inteligenca in strojno učenje - ilustrativni primeri

- Termostat: senzor in vršilec je stroj



- Odobritev kredita: vršilec je človek, ki odločitev sprejme (delno) na podlagi priporočila modela

Strojno učenje

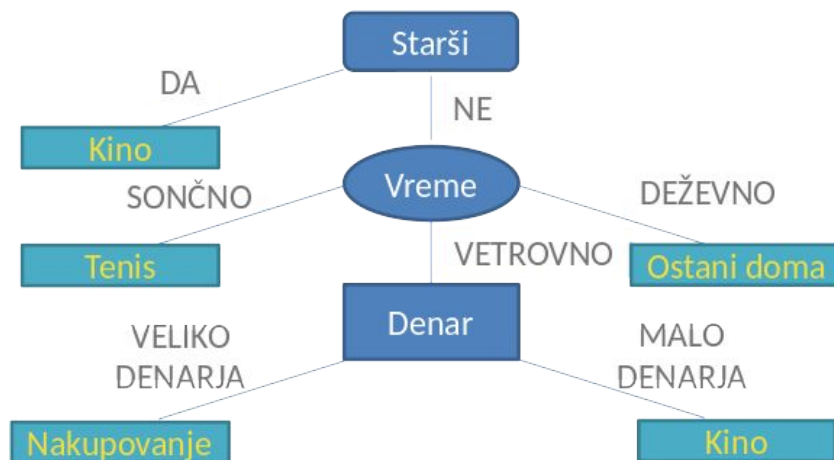
Podatki

Značilke

Tarčna spremenljivka

Teden	Vreme	Starši na obisku	Denar	Odločitev (kategorija)
T1	Sončno	Da	Veliko	Kino
T2	Sončno	Ne	Veliko	Tenis
T3	Vetrovno	Da	Veliko	Kino
T4	Deževno	Da	Malo	Kino
T5	Deževno	Ne	Veliko	Ostali doma
T6	Deževno	Da	Malo	Kino
T7	Vetrovno	Ne	Malo	Kino
T8	Vetrovno	Ne	Veliko	Nakupovanje
T9	Vetrovno	Da	Veliko	Kino
T10	Sončno	Ne	Veliko	Tenis

Model



Strojno učenje - model

Iščemo **funkcijo** $f(x)$ s **parametri** a, b, c, \dots , ki opiše **podatke** (x_i, y_i) .

Optimiziramo parametre, tako da je **napaka** čim manjša.

Definicija napake se razlikuje od aplikacije do aplikacije.

Model: Funkcija f in množica optimiziranih parametrov a, b, c, \dots

Strojno učenje - variacije

Nadzorovano učenje

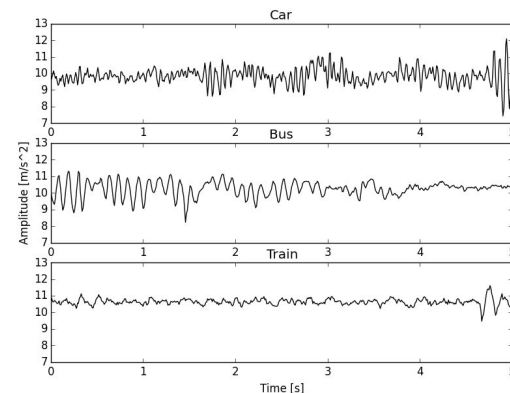
- **Znane** vrednosti tarčnih spremenljivk (y_i)
- Primeri: klasifikacija, regresija

Nenadzorovano učenje

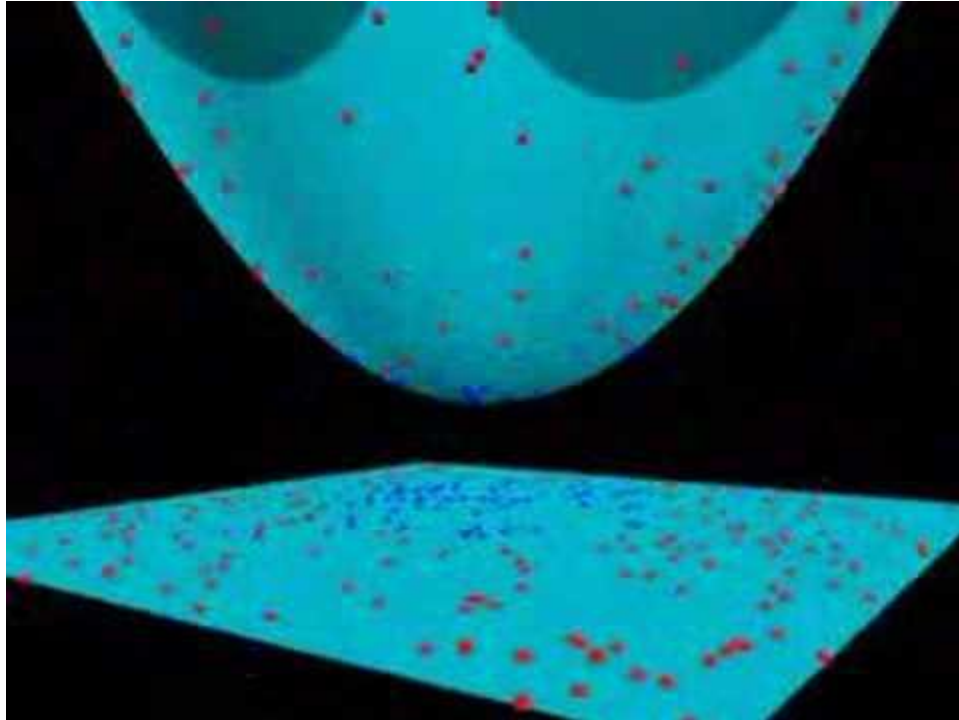
- Tarčne spremenljivke in njihove vrednosti **niso znane** - pogosto gre za iskanje podobnosti
- Primeri: razvrščanje v skupine, detekcija anomalij

Strojno učenje - klasifikacija

- **Razvrščanje v razrede**
- Primerjava na podlagi podobnosti ali razdalje
- V praksi:
 - Detekcija prevoznega sredstva
 - Razpoznava ročno zapisanih števk



Strojno učenje - klasifikacija



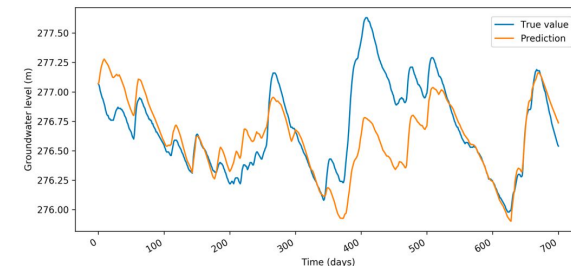
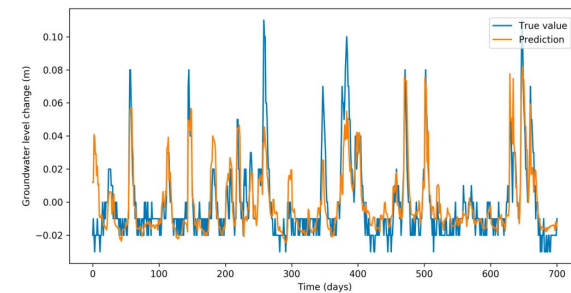
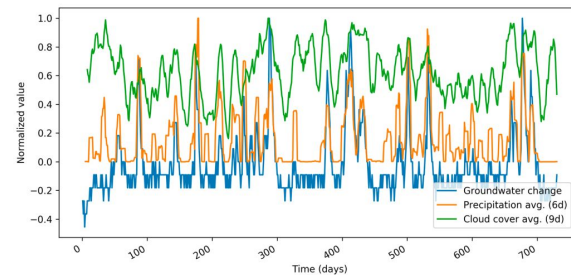
Strojno učenje - regresija

- Uporabno pri **napovedovanju številskih vrednosti**

- Občasno lahko uporabimo regresijo, nato pa izhod interpretiramo kot kategorije - klasifikacija

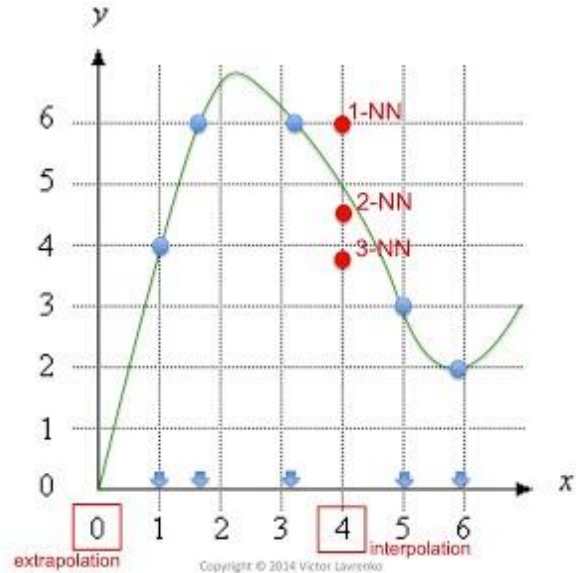
- V praksi

- Napovedovanje porabe energije
- Ocenjevanje nivoja podtalnice



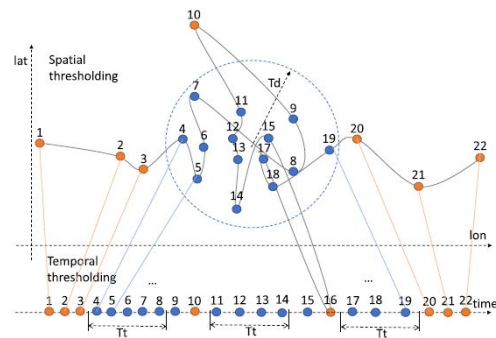
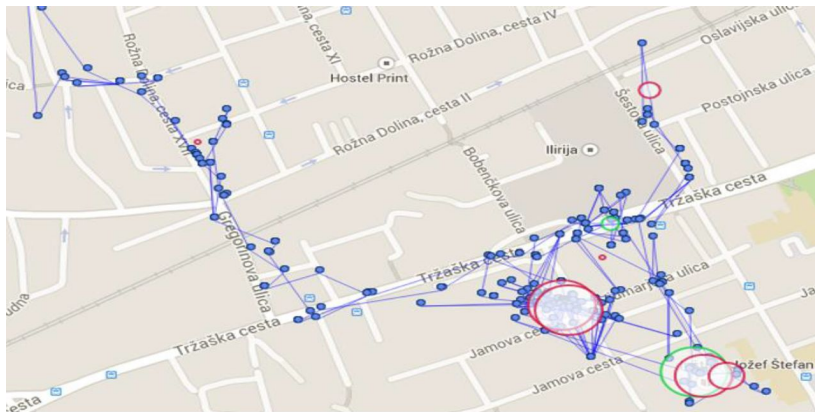
Strojno učenje - regresija

Example: kNN regression in 1-d

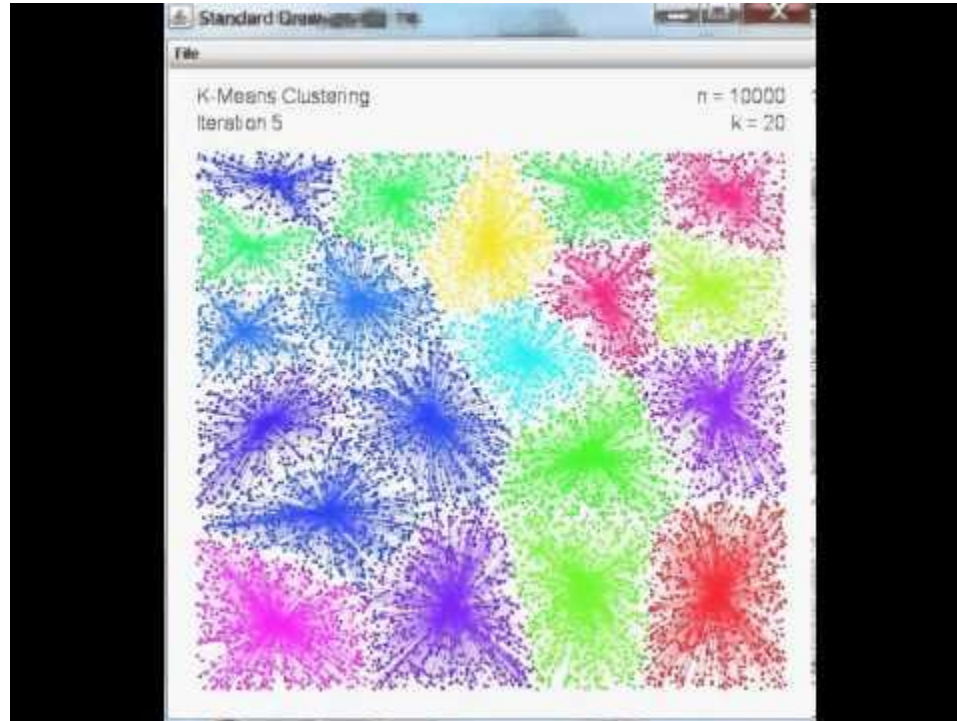


Strojno učenje - razvrščanje v skupine

- Vnose želimo razporediti v skupine **glede na podobnost** njihovih značilk
- V praksi:
 - Odkrivanje interesnih točk iz sledi GPS

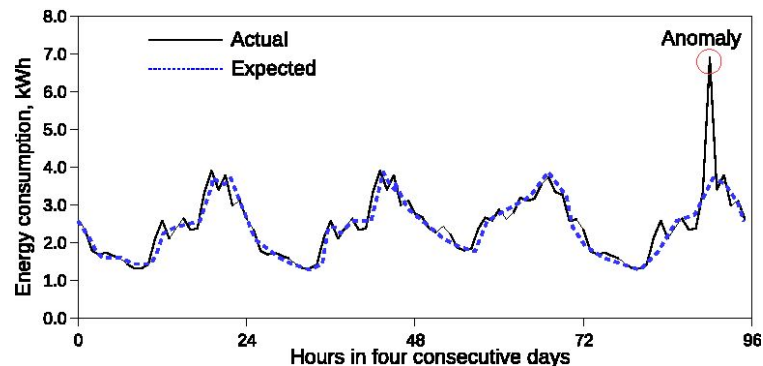


Strojno učenje - razvrščanje v skupine

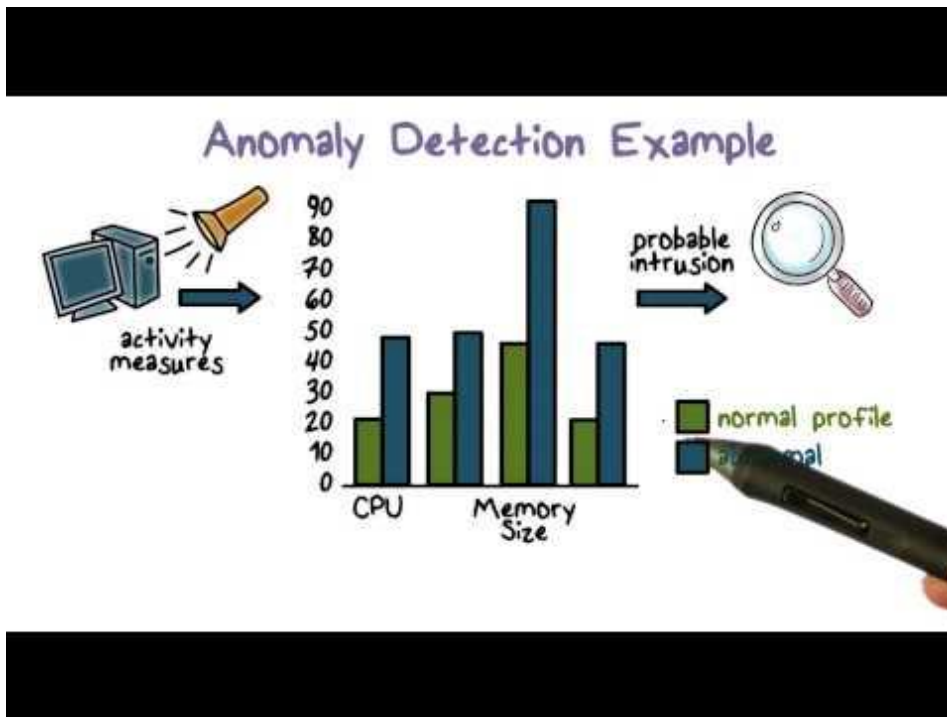


Strojno učenje - detekcija anomalij

- V podatkih želimo **identificirati redke dogodke**, elemente, pojave, ki se **močno razlikujejo** od ostalih vrednosti
- Anomalije so lahko tudi osamelci, novitete, šum ...
- V praksi:
 - Bančne goljufije
 - Spletni incidenti in vdori
 - Strukturni defekti



Strojno učenje - detekcija anomalij



Primeri strojnega učenja na velikih podatkih

Primer 1: Modeliranje verjetnosti neplačila

Opis problema:

- Podjetje nosi tveganje, da stranke ne poplačajo svojih obveznosti.
- Podjetje lahko omeji svoje tveganje:
 - Zavarovanje obveznosti
 - **Preventivni ukrepi**

Cilj:

- Modeliranje verjetnosti neplačila posamezne stranke
- Priporočilo višine limita kredita glede na izračunano verjetnost neplačila

Primer 1: Modeliranje verjetnosti neplačila - podatki

- **Finančni podatki** - javnodostopni podatki o prihodkih in odhodkih podjetja, bilance, letna poročila podjetij
- **Podatki o trgovanju** med podjetjem in njegovimi strankami - mesečni podatki o obsegu trgovanja, obstoječih dolgovih, izpodbijanih terjatvah in zapoznelih plačilih

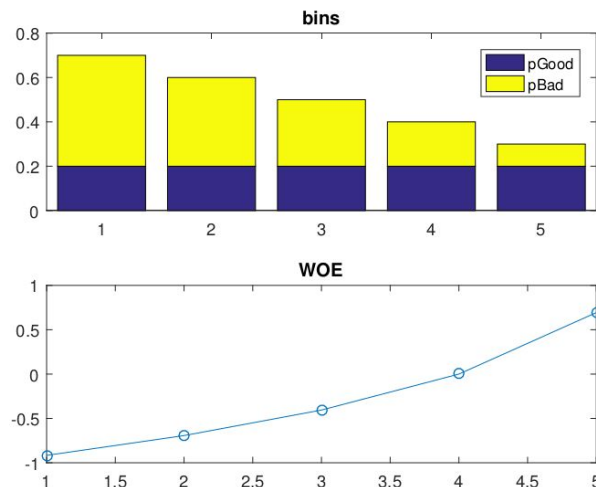
Primer 1: Modeliranje verjetnosti neplačila - značilke

- Finančni podatki \mapsto finančni indikatorji

- Finančni indikatorji \mapsto značilke

- Razvrščanje v razrede
- Transformacija v *weight of evidence* (WOE)

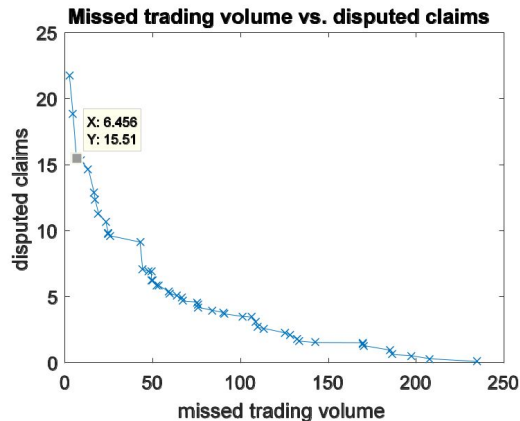
$$WOE = \log \frac{P(\text{company}=\text{good})}{P(\text{company}=\text{bad})}$$



Primer 1: Modeliranje verjetnosti neplačila - model

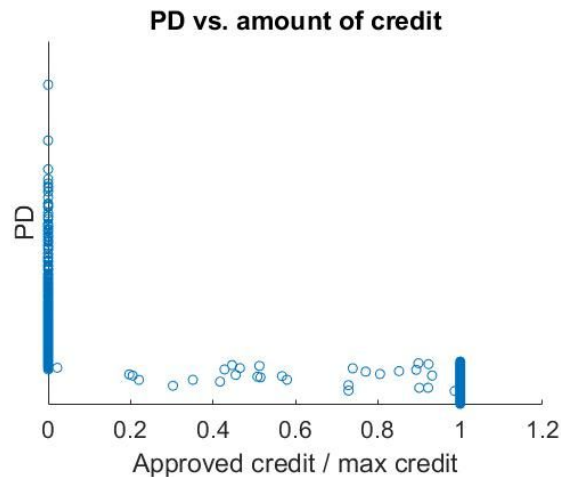
Verjetnost neplačila za podjetje

$$PD(x) = \frac{1}{1 + \exp[-(\beta_0 + \beta_1 WOE(x_1) + \beta_2 WOE(x_2) + \dots + \beta_n WOE(x_n))]}$$



Primer 1: Modeliranje verjetnosti neplačila - kreditni limit

Optimizacija izgubljenega prometa proti izgubi zaradi neplačil



Primer 2: Napovedovanje porabe električne energije

Opis problema:

- Modelirati želimo energetske fenomene (npr. poraba električne energije, cena energije na trgu, proizvodnja električne energije)
- Senzorski podatki prihajajo v sistem z visoko frekvenco (npr. električni tok, električna moč ...)
- Podatki iz napovedi (npr. vremenska napoved) se običajno spreminjajo
- Statične podatke lahko izračunamo vnaprej (npr. dan v tednu, čas v dnevu, dan v letu, lunine mene, prazniki ...)

Primer 2: Napovedovanje porabe električne energije

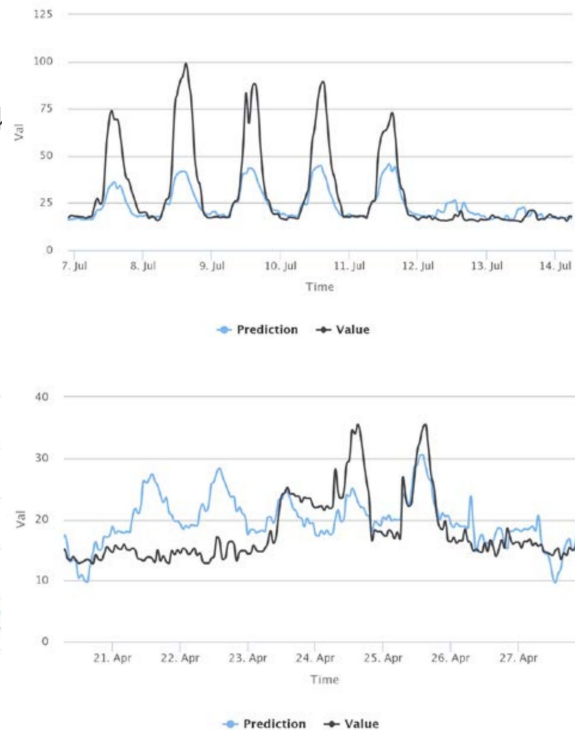
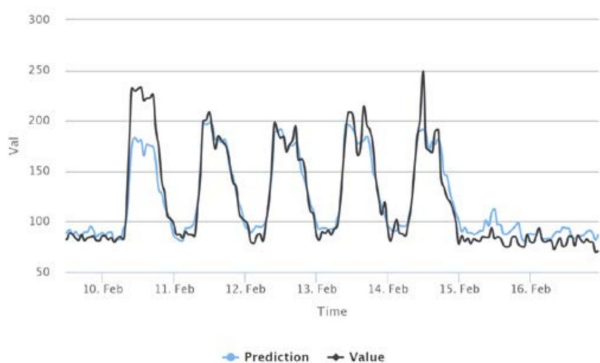
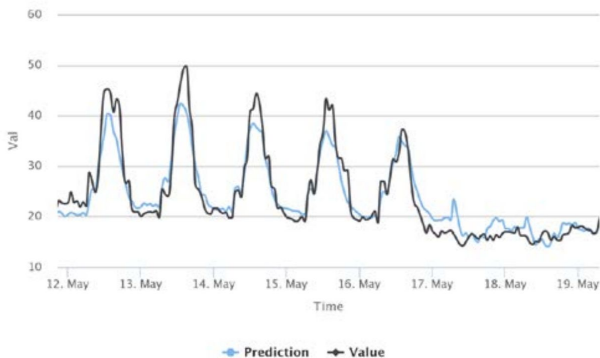
Cilj:

- Čim bolj natančno in zanesljivo želimo modelirati (**napovedovati**) energetske fenomene ob upoštevanju senzorskih in statičnih podatkov ter podatkov iz napovedi

Primer 2: Napovedovanje porabe električne energije

Poraba javne ustanove

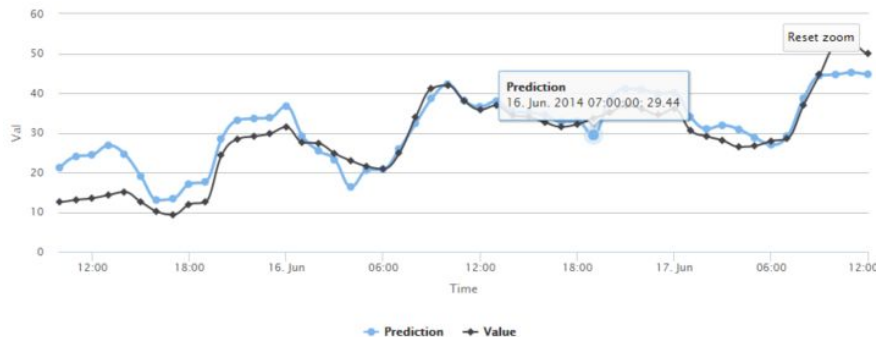
- **Dobre značilke:** tedenski agregati senzorskih značilni količina sončnega obsevanja, delovni čas, kurilna sezona, etc.
- **Kakovost podatkov je včasih vprašljiva:** manjkajoče in napačne vrednosti



Primer 2: Napovedovanje porabe električne energije

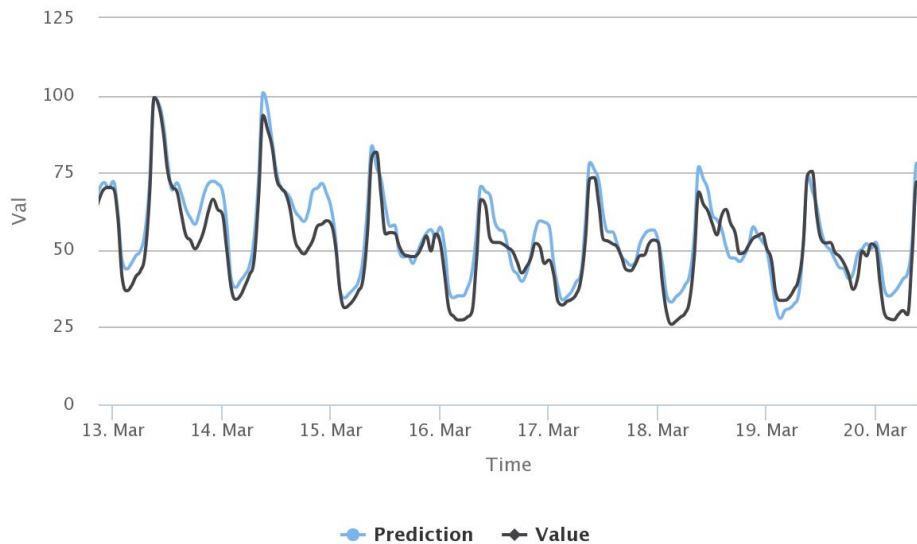
Cena električne energije (Energy Spot Market – EPEX)

- Močno odvisna od proizvodnje elektrike iz alternativnih virov
 - Cena proizvodnje je nizka, a omrežje ni optimizirano za neredne vire, zato se cena elektrike zniža
- **Dobre značilke:** cena v prejšnjih dneh, povprečja v tednu/mesecu, smer vetra



Primer 2: Napovedovanje porabe električne energije

Modeliranje proizvodnje v termoelektrarni



Vprašanja in diskusija