

Tri principa učenja

- Occam's Razor
- Sampling Bias
- Data Snooping

Jednostavna hipoteza

- An explanation of the data should be made *as simple as possible, but no simpler* – Albert Ajnštajn
- Okamova oštrica – simbol principa
 - Imamo objašnjenje podataka. Nastavljamo da „orezujemo“ objašnjenje sve dok ne stignemo do minimuma koji je i dalje konzistentan sa podacima
 - Kada dobijemo ovaj minimum, njega smatramo najboljim mogućim objašnjenjem

Okamova oštrica

- Najjednostavniji model koji odgovara podacima je najverodostojniji

1. Šta znači da je model jednostavan?
2. Kako znamo da je jednostavnije bolje?

Šta znači da je model jednostavan?

- Dva tipa merenja kompleksnosti:
 1. Kompleksnost pojedinačne hipoteze h
 2. Kompleksnost skupa hipoteza \mathcal{H}

Kompleksnost pojedinačne hipoteze h

- Ove mere se odnose na pojedinačan objekat – kompleksnost je svojstvo samog objekta
- MDL (*Minimum Description Length*)
 - Dati objekat pokušavamo da specificiramo sa najmanje moguće „bitova“. Što nam manje „bitova“ treba, objekat je jednostavniji
 - Npr. integer od 10^6 cifara. Kompleksnost pojedinačnih integera te dužine varira. Npr. $2^6 - 1$ je jednostavan jer smo ga mogli tako (jednostavno) opisati
- Stepen polinoma

Kompleksnost skupa hipoteza \mathcal{H}

- Kompleksnost se izračunava za *skup* objekata
- Entropija
 - Izvršite eksperiment
 - Razmotrite sve moguće ishode i verovatnoće koje idu uz te ishode
 - Formirate jedinstvenu kolektivnu funkciju koja obuhvata ovu verovatnoću
$$\sum p(x_i) \log \left(\frac{1}{p(x_i)} \right)$$
 - Izračunava se za *klasu* objekata – svaki ishod je jedan objekat
- VC dimenzija
 - Svojstvo skupa hipoteza
 - Posmatra skup hipoteza kao celinu i predstavlja jedinstven broj koji označava raznolikost tog skupa hipoteza
 - U ovom slučaju, raznolikost znači kompleksnost

Šta znači da je model jednostavan?

- Kada mislimo „jednostavno“, obično mislimo na jedan objekat h
 - Ne mislimo na alternative koje postoje u opisu podataka
- Dokazi Okamove oštrice su obično u terminima skupa hipoteza \mathcal{H}
 - Sa ovim smo se već susreli, npr. VC dimenzija
- Ovo je malo zabrinjavajuće – intuitivan koncept je jedno, a matematički dokaz drugo

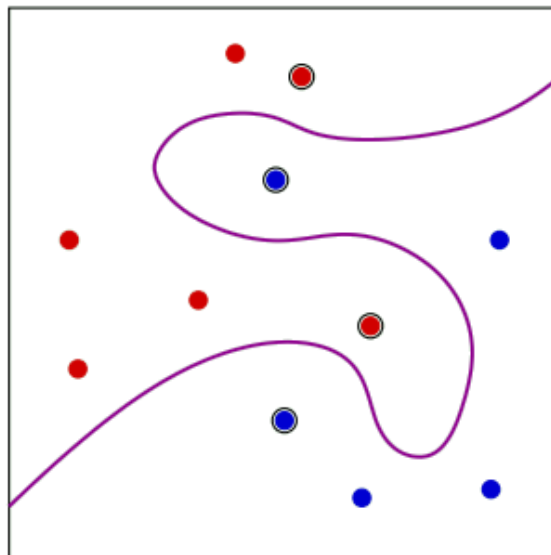
Koncept koja ih povezuje: prebrojavanje

- Dobra vest jeste da su koncepti kompleksnosti objekta i kompleksnosti skupa objekata veoma povezani (skoro identični)
 - l bitova specificiraju h
 - Ovo implicira da postoji 2^l elemenata sličnih h (koji se takođe mogu opisati sa l bitova)
 - Skup svih sličnih objekata možemo označiti sa \mathcal{H}
 - „Jedan od 2^l “ možemo koristiti kao opis kompleksnosti \mathcal{H}
- l bitova specificira $h \Rightarrow h$ je jedan od 2^l elemenata skupa \mathcal{H}
 - Koncept: objekat je kompleksan ako je jedan od mnogih. Objekat je jednostavan ako je jedan od nekolicine
- Šta je sa parametrima koji su realni brojevi (npr. polinom 17. stepena)?
 - I dalje odgovaraju našem opisu „jedan od mnogih“

Izuzetak od pravila

- Izuzetak od ovog pravila (koji izgleda kompleksan ali je samo „jedan od nekolicine“)
 - Namerni izuzetak – želeli smo kompleksan model koji može dobro da se prilagodi podacima. Ali ipak je jedan od nekolicine – nismo želeli da platimo punu cenu kompleksnosti

SVM



Pitanje: predikcija ishoda fudbalske utakmice

- Dobili ste pismo koje predviđa ishod utakmice koja se to veče održava
- Ispostavilo se da je predikcija tačna! Ali možda je samo srećan pogodak...
- Tokom sledećih 5 nedelja dobili ste još 5 ovakvih pisama. Sve predikcije su ispale tačne!
- U šestoj nedelji dobijate pismo: „Želite još? 50\$“
- Da li da platite?

0000 0000 0000 0000	1111 1111 1111 1111	0
0000 0000	1111 1111	1
	0000 1111	0
	0011	1
	01	1

Kako znamo da je jednostavnije bolje?

- Bolje ne znači elegantnije! Bolje znači bolje performanse van uzorka E_{out}
- Argument*:
 1. Postoji manje jednostavnih hipoteza (u poređenju sa brojem kompleksnih hipoteza) $m_{\mathcal{H}}(N)$
 2. Pošto ih ima manje, manje je verovatno da će savršeno odgovarati datom skupu podataka:

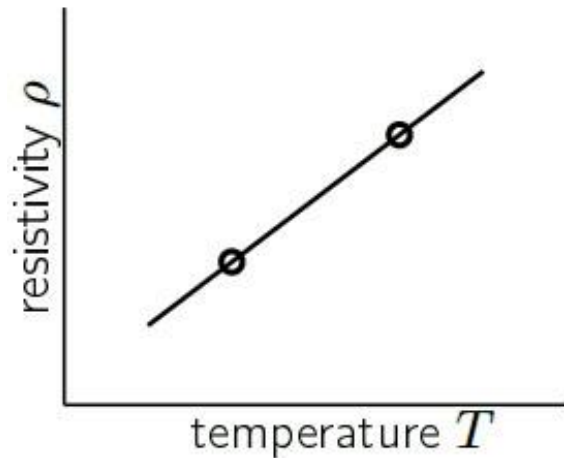
$$\frac{m_{\mathcal{H}}(N)}{2^N}$$

3. Ako je nešto manje verovatno, značajnije je kada se to desi
- Npr. u prevari sa pismima
 - Onaj ko je prevaren vidi samo jednu hipotezu i ona je perfektna – zato joj pridodaje veliki značaj jer je prilično neverovatno da će se to desiti
 - Onaj ko vidi širu sliku zna da je $m_{\mathcal{H}}(N) = 2^N$ - sigurno je da će se desiti, zbog toga je beznačajno

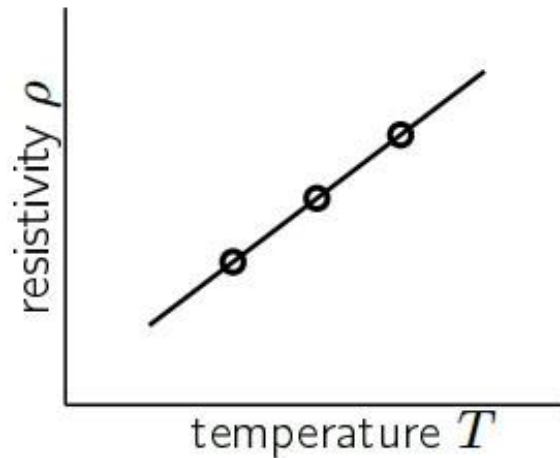
*Formalan dokaz postoji pod različitim idealizovanim uslovima

Primer beznačajnog eksperimenta

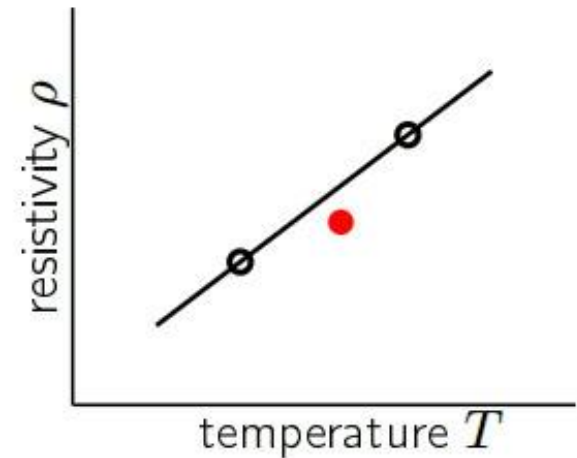
- Želimo da utvrdimo da li su provodljivost određenog metala i temperatura u linearnoj zavisnosti
- Kakav dokaz nam sledeći rezultati na slici pružaju?
- **Aksiom neopovrgljivosti (*the Axiom of Non-Falsifiability*)** – ako sa datim podacima nemamo mogućnost da opovrgnemo tvrdnju, oni zbog toga ni ne mogu pružiti nikakav dokaz u korist te tvrdnje



Scientist A



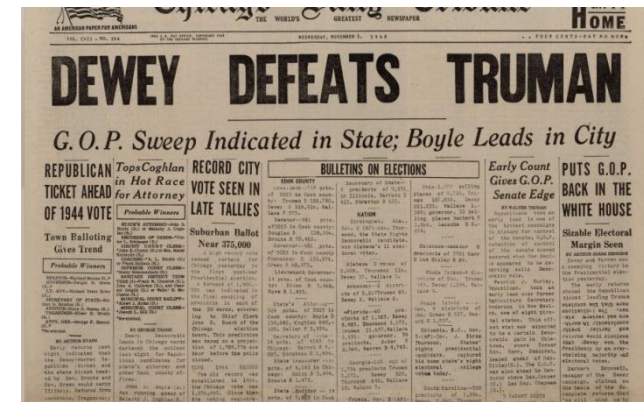
Scientist B



"falsifiable"

Sampling bias – izbor predsednika

- 1948, prvi predsednički izbori nakon Drugog svetskog rata – Truman i Dewey
- Prema anketama, kandidati su veoma blizu i nije jasno ko će pobediti
- Nakon što su izbori završeni, ali glasovi još nisu prebrojani, jedna redakcija je sprovedla telefonsku anketu i pitala ljude kako su glasali
- Dobili su rezultat da je Dewey neosporno pobedio
- Rezultat je izgledao toliko očigledan da su odlučili da budu prvi koji će izvestiti o tome i ištampali novine sa naslovom da je Dewey pobedio
- Pobednik (koji na slici drži novine) je Truman



Sampling bias – izbor predsednika

- Šta je pošlo po zlu? Da li je ovo posledica probabilističkih garancija?

$$P[|E_{in} - E_{out}| > \epsilon] \leq \delta$$

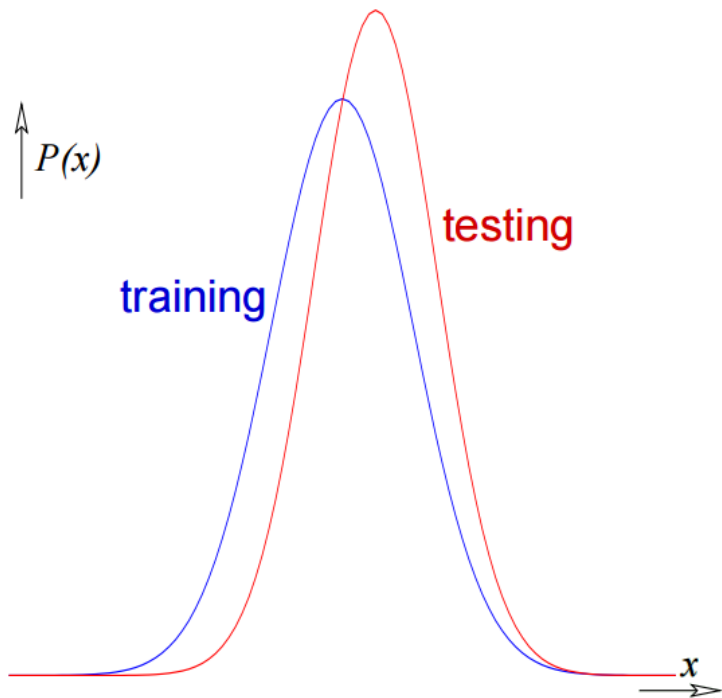
- U ovom slučaju, nije δ krivo. Dobijeni rezultat je bio posledica sistematske greške (*bias*) u uzorkovanju:
 - 1948 telefoni su bili skupi – kuće sa telefonom su obično bile bogatije
 - Bogati ljudi su preferirali Dewey-a

Sampling bias principle

- If the data is sampled in a biased way, learning will produce a similarly biased outcome
- Moramo se obezbediti da podaci budu reprezentativni u odnosu na ono šta želimo da pokažemo
- Praktičan primer: financial forecasting
 - Situacija je prilično nepredvidiva, pod uticajem neke glasine tržište može drastično da se promeni
 - Ako želimo da pronađemo šablon koji postoji u podacima, razmatramo „normalne“ periode u kojima se taj šablon vidi
 - Kada testiramo model, testiramo ga na pravom tržištu. Može se desiti da uvidimo da postoji sistematsko odstupanje

Poklapanje distribucija

- Jedan način da se borimo sa sistematskim greškama u uzorkovanju jeste poklapanje distribucija
- Pretpostavka uvedena kako bi *Hoeffding*-ova nejednakost važila:
 - Instance korišćene za obučavanje su odabrane iz iste distribucije kao i instance koje ćemo koristiti za testiranje
 - Ako imamo sistematsku grešku u uzorkovanju, ovo je narušeno



- Ako imamo pristup trening i test distribuciji, možemo proveriti da li se poklapaju, npr., u datom primeru su trening i test distribucije donekle različite
- Ako poznajemo ove distribucije, možemo:
 - dodeliti različite težine primerima trening skupa
 - Iz datih podataka ponovo uzorkovati trening skup tako da ispadne kao da je izvučen iz druge distribucije

Poklapanje distribucija

- Ovaj metod radi i u praksi, čak i ako ne znamo konkretne distribucije, možemo da ih procenimo
- Međutim, metod ne radi ako postoji regija u ulaznom prostoru gde je $P = 0$ za trening, ali je $P > 0$ za testiranje
 - Npr. ljudi koji nisu imali telefon u 1948
 - Ne možemo ništa da uradimo pomoću poklapanja distribucija jer nemamo predstavu šta se u tom delu dešava
- Dakle, u određenim situacijama postoji rešenje za sampling bias
- Ali, u nekim situacijama, sve što možemo da uradimo jeste da priznamo da ne možemo garantovati performanse našeg rešenja u delovima koji nisu pokriveni uzorkom

Pitanje – pokušajte da detektujete sampling bias

- Odobravanje kredita mušteriji
- Istorijski podaci mušterija iz prethodne 2-3 godine
 - Dostupne su nam informacije koje svaka mušterija daje prilikom aplikacije za kredit (jer su to jedine informacije koje ćemo imati za nove mušterije)
 - I dostupno nam je ciljno obeležje – u retrospektivi, da li je banka zaradila na ovim mušterijama
- Gde je ovde sampling bias?
- Koristimo podatke o mušterijama koje smo ranije odobrili (zato što su to jedine mušterije za koje imamo podatke o vraćanju kredita)
- Mušterije koje smo odbili nisu deo ovog trening skupa
- Za novu mušteriju ne znamo da li bi ova mušterija bila odbijena ili ne prema starim kriterijumima banke, dakle, ona bi mogla biti deo skupa koji nikada nije pokriven trening skupom
- Međutim, u ovoj primeni, sampling bias nema katastrofalne posledice
 - Banke prilično agresivno dodeljuju kredit pa imamo i dovoljno primera mušterija na kojima je banka pogrešila

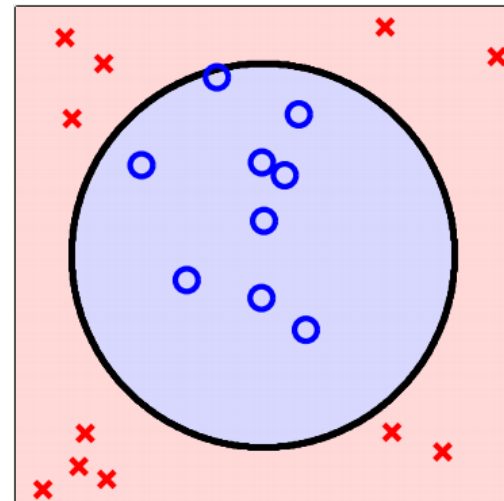
Data snooping

- Princip: ako je skup podataka imao uticaja na bilo koji korak učenja, onda je mogućnost ovog skupa podataka da proceni ishod učenja kompromitovana
- Ovo je zamka u koju upadaju mnogi
 - Ima mnogo načina da se uhvatimo u nju
 - I veoma je privlačno da upadnemo u nju jer dobijamo bolje performanse
 - Manifestuje se na mnogo različitih načina

Data snooping primeri

1. Pogledali smo podatke

- Recimo da imamo primere sa slike
- Bez gledanja smo rešili smo da primenimo transformacijo drugog stepena $z = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$. Greška na uzorku je mala, ali plaćamo cenu generalizacije jer imamo više parametara
- Ali, nakon inspekcije, shvatimo da nam ne trebaju sva obeležja – samo $z = (1, x_1^2, x_2^2)$ ili čak $z = (1, x_1^2 + x_2^2)$



Data snooping primeri

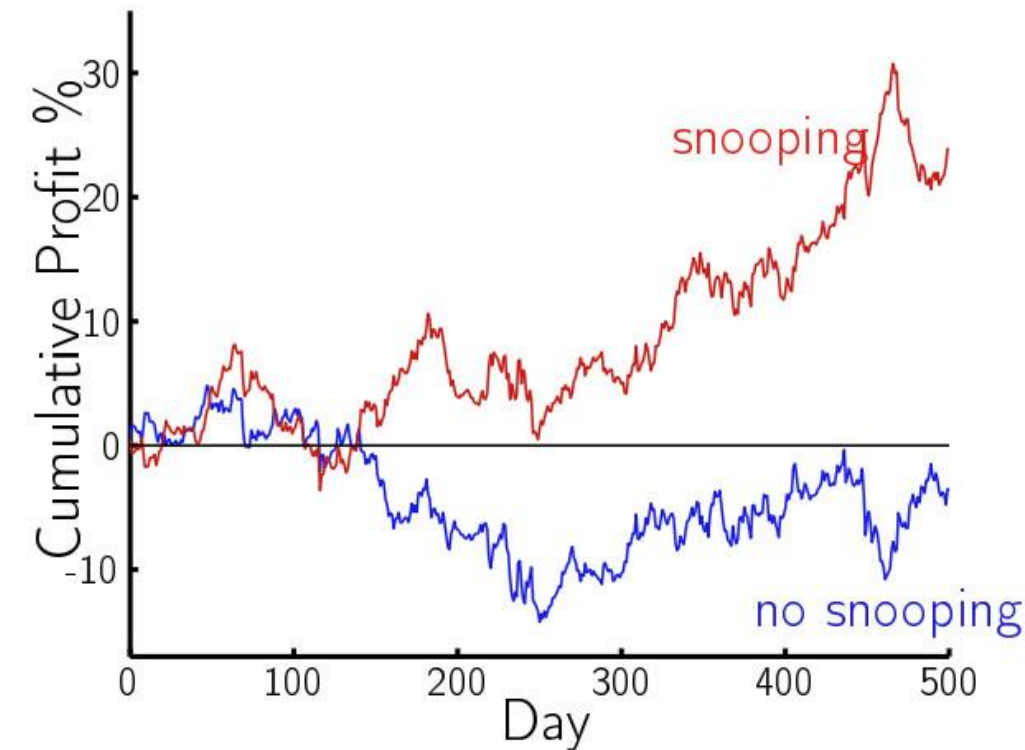
1. Pogledali smo podatke

- Ako pogledamo skup podataka, lako se može desiti da dizajniramo model da bude prilagođen tom konkretnom skupu podataka
- Iako radimo dobro na tom skupu podataka ne znamo da li ćemo jednako dobro raditi na drugom skupu podataka (generisanom iz iste distribucije)
- Data snooping uključuje skup podataka T , ali ne i druge informacije
 - U proces obučavanja možemo (i dobro je) uključiti domensko znanje (npr. koliko imamo ulaza, koliki su im opsezi, kako su mereni, da li su fizički korelirani, ...)
 - Samo ne treba razmatrati konketan skup podataka

Pitanje

- Da li možete identifikovati data snooping u ovom slučaju
- Financial forecasting: predviđanje odnosa između \$ i £
 - Imamo osam godina podataka o odnosu na dnevnom nivou (oko 2000 tačaka)
 - Izlaz je promena odnosa \$ i £ u datom danu u odnosu na prethodni Δr_0 , a ulazi predstavljaju promene tog odnosa u prethodnih 20 dana: $\Delta r_{-20}, \Delta r_{-19}, \dots, \Delta r_{-1} \rightarrow \Delta r_0$
 - Prvo, normalizujemo podatke ($\mu = 0, \sigma^2 = 1$)
 - Zatim ih podelimo na trening T_{train} ($N_{train} = 1500$) i test skup T_{test} ($N_{test} = 500$)
 - Za test skup smo odabrali primere na slučajan način, ne samo poslednje zabeležene dane
 - Ni u jednom trenutku nismo gledali podatke, sve analize smo izvršili automatski

Pitanje



- Data snooping se desio kada smo normalizovali podatke pre podele na trening i test skup
 - Koristili smo srednju vrednost i standardnu devijaciju test skupa!

- Korektno:

1. Podeliti na T_{train} i T_{test}
2. Normalizovati T_{train} . Sačuvati μ_{train} i σ_{train}^2 kako bismo identičnu transformaciju primenili na test podatke

Data snooping

3. Recikliranje skupa podataka

- *If you torture the data long enough, it will confess*
- Isprobavamo mnogo obučavajućih algoritama na istom skupu podataka
 - Podelili smo podatke na T_{train} i T_{test}
 - Obučavamo različite modele na T_{train} i evaluiramo ih na T_{test} (E_{test})
 - Kao rezultat vratimo model koji je imao najbolje performanse na T_{test} i kažemo da su njegove performanse E_{test}
- Problem jeste što uvećavamo VC dimenziju, bez da to shvatimo - prava VC dimenzija je **unija** svih modela koje smo isprobali
- Ovo može da obuhvati i ono šta su drugi probali!
 - Ako koristimo javno dostupan skup podataka i drugi ljudi su već isprobali stvari na njemu
 - Mi pročitamo te radove. Npr. saznamo da se najbolje pokazao SVM sa polinomijalnim kernelom

Data snooping

- Ključni problem u svim primerima je što se prilagođavamo konkretnom skupu podataka – počinjemo da se prilagođavamo šumu koji postoji u njemu

Dva rešenja za data snooping

1. Izbegavanje

- Disciplina – stavite test podatke u sef i nemojte ga otvarati sve dok nemate **finalnu** hipotezu

2. Uračunajte i efekat data snooping-a u performanse

- Kolika je kontaminacija podataka (VC dimenzija, ...)
- Najteže je uraditi ako ručno pogledamo podatke – teško je modelovati sebe (koliki skup hipoteza je razmatran)

Pitanje: bias via snooping

- Testiranje dugoročnih performansi „buy and hold“ kupovine akcija. Hoćemo da predvidimo kako ćemo proći
 - „buy and hold“ - ne možemo prodavati/menjati u nekom povoljnom trenutku pa kasnije ponovo kupovati, moraju ostati u našem posjedstvu od početka do kraja
- Koristićemo podatke od prethodnih 50 godina
 - Hoćemo da test bude što širi – uzmemo sve *Standard & Poor's 500* kompanije (uobičajen reper za U.S. akcije)
 - Pretpostavimo striktno „buy and hold“ model za sve akcije
- Recimo da smo predvideli fantastičan profit. Da li problem postoji?

Pitanje: bias via snooping

- Postoji. Sampling bias – gledamo akcije koje se trenutno preprodaju. Ne razmatramo sve one koje su propale
- Ljudi ovo često ne tretiraju kao sampling bias, već kao data snooping (iako se ne uklapa sasvim u našu raniju definiciju)
 - Jeste snooping – kao da gledamo 50 godina u budućnost i neko nam kaže kojim akcijama se još trguje u toj tački
 - Ali je više sampling bias prouzrokovan pomoću data snooping

Zaključak

- Sva tri koncepta koja smo prešli predstavljaju „zamke“ sa kojima se možemo sresti u primeni mašinskog učenja
 - Npr. Okamova oštrica – vodite računa o kompleksnosti modela koji primenjujete, prilagodite je resursima