

Prediktivni model za ranu identifikaciju ispada diskovnih jedinica

Vugrinec, Ivica

Undergraduate thesis / Završni rad

2019

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zagreb, Faculty of Organization and Informatics / Sveučilište u Zagrebu, Fakultet organizacije i informatike**

Permanent link / Trajna poveznica: <https://urn.nsk.hr/urn:nbn:hr:211:551403>

Rights / Prava: [Attribution 3.0 Unported/Imenovanje 3.0](#)

Download date / Datum preuzimanja: **2024-09-04**



Repository / Repozitorij:

[Faculty of Organization and Informatics - Digital Repository](#)



**SVEUČILIŠTE U ZAGREBU
FAKULTET ORGANIZACIJE I INFORMATIKE
VARAŽDIN**

Ivica Vugrinec

**PREDIKTIVNI MODEL
ZA RANU IDENTIFIKACIJU
ISPADA DISKOVNIH JEDINICA
ZAVRŠNI RAD**

Varaždin, 2019.

SVEUČILIŠTE U ZAGREBU
FAKULTET ORGANIZACIJE I INFORMATIKE
V A R A Ž D I N

Ivica Vugrinec

Matični broj: 31819/99-I

Studij: Informacijski sustavi

PREDIKTIVNI MODEL
ZA RANU IDENTIFIKACIJU
ISPADA DISKOVNIH JEDINICA
ZAVRŠNI RAD

Mentorica:

doc. dr. sc. Dijana Oreški

Varaždin, siječanj 2019.

Izjava o izvornosti

Izjavljujem da je moj završni rad izvorni rezultat mojeg rada te da se u izradi istoga nisam koristio drugim izvorima osim onima koji su u njemu navedeni. Za izradu rada su korištene etički prikladne i prihvatljive metode i tehnike rada.

Autor potvrdio prihvaćanjem odredbi u sustavu FOI-radovi

Sažetak

Ovaj završni rad, bavi se analizom velikog skupa podataka koji redovito objavljuje tehnološka kompanija BackBlaze na svojim javnim Internet stranicama. Analiziramo strukturu objavljenih podataka te vodeći se principima CRISP-DM metodologije radimo na razumijevanju, pripremi i čišćenju podataka.

U radu radimo različite upite nad bazom podataka, analizirajući izvorne podatke u kontekstu poslovanja tehnološke tvrtke BackBlaze, uspoređujući podatke u vremenu pokušavajući dokazati da je analizom tehničkih parametara diskovnih jedinica moguće dovoljno rano predvidjeti ispad pojedine diskovne jedinice. U ovom se radu pokušava pokazati da nadzorom vrijednosti nekih od velikog broja svakodnevno uzorkovanih podataka mogu s velikom vjerojatnošću predvidjeti ispad diskovne jedinice i samim time omogućiti donošenje poslovnih odluka o njihovoj zamjeni u cilju smanjenja poslovnih rizika, vremena oporavka, dodatnih troškova povezanih uz nekontrolirane ispade pojedinih diskovnih jedinica.

Dodatnu dimenziju ovog rada daje analiza potreba, priprema i isporuka tehničko/tehnološke platforme na kojoj su podaci analizirani. Istraživani su i primjenjeni tehnološki koncepti poput instalacije, konfiguracije i pripreme linux operativnog sustava, pripadajuće (MySQL) baze podataka, import velikog podatkovnog seta u bazu podataka, priprema složenih upita nad velikim podatkovnim setom, optimizacija upita i slično.

Korištenjem funkcija iz programskog jezika R radimo izravne upite nad relacijskom bazom podataka te upotrebom ugrađenih funkcija izrađujemo vizualizaciju pripremljenih podataka u kontekstu analize podataka i izvođenja zaključaka uz pomoć vizualizacije rezultata.

Zaključak rada potvrđuje pretpostavku da je analizom vremenskih serija podataka prikupljenih svakodnevnim uzorkovanjem moguće predvidjeti ispade. Iako su do samog zaključka moglo doći i jednostavnom logikom, cilj rada bio je i praktično provjeriti i pokazati isto.

Ključne riječi: Diskovne jedinice, S.M.A.R.T, CRISP-DM, priprema, čišćenje i analiza objavljenih podataka

1	Uvod	1
2	Metodologija i tehnike rada.....	2
2.1	CRISP-DM metodologija	2
2.2	Razumijevanje poslovanja	4
2.3	Razumijevanje podataka.....	5
2.4	Priprema podataka	5
2.5	Modeliranje podataka	6
2.6	Vrednovanje podataka.....	6
2.7	Korištenje podataka	6
3	Razumijevanje poslovanja na primjeru tvrtke Backblaze	8
3.1	Ciljevi i zahtjevi poslovnog okruženja.....	8
3.2	Backblaze - sukobljeni poslovni ciljevi.....	9
3.3	Primjena metoda rudarenja podataka u svrhu podrške odlučivanju	10
3.4	Definicija pretpostavke i hipoteza.....	11
3.5	Kriteriji uspjeha primjene rudarenja podataka u svrhu podrške odlučivanju ..	12
4	Razumijevanje podataka	14
4.1	Okolina u kojoj se podaci prikupljaju	14
4.2	Uzorkovani podaci i njihovo značenje	15
4.2.1	Struktura analiziranih podataka	15
5	Priprema podataka	17
5.1	Kreiranje podataka i priprema baze podataka.....	17
6	Čišćenje podataka	18
6.1	Nedostatak uzorkovanih podataka – prazna polja.....	18
6.2	Nekonzistentna polja i nejasne vrijednosti	18
6.3	Vrijednosti izvan očekivanih granica	18
7	Redukcija podataka	19
7.1	Redukcija diskovnih jedinica za operative sustave.....	19
7.2	Redukcija diskovnih jedinica koje imaju premali statistički uzorak.....	19
8	Istraživanje podataka.....	20
8.1	Učestalost kvara na temelju uzorkovanih podataka	20
9	Razumijevanje i vizualizacija podataka.....	23

9.1	Posljednji kvartal 2016. godine	23
9.1.1	Rezultat analize i zaključak obrađenih podataka.....	24
9.2	Dulji vremenski termin – cijela 2016. godina	24
9.2.1	Rezultat analize i zaključak obrađenih podataka.....	25
9.3	Dulji vremenski termin (1. 1. 2015. – 31. 12. 2016.).....	26
9.3.1	Rezultat analize i zaključak obrađenih podataka.....	27
10	Problem dva različita kraja životnog vijeka diskovne jedinice.....	28
10.1	Kaplan-Meierova krivulja doživljenja - primjena medicinskih statističkih metoda nad promatranim podacima.....	30
10.2	Alati za analizu podataka - „R“	30
10.3	Konstrukcija Kaplan-Meierovih krivulja u alatu „R“.....	31
10.3.1	Izvorni kod izrade Kaplan-meirovih krivulja doživljenja u programskom jeziku „R“	32
10.4	Rezultati analize podataka u programskom jeziku „R“	37
11	Zaključak istraživanja:	41
12	Popis literature:.....	43
13	Popis slika:	44
14	Popis tablica:.....	45

1 Uvod

U ovom radu razmatramo primjenu CRISP-DM metodologije rudarenja podataka u cilju povećanja potencijala za podršku odlučivanju.

Primjenom navedene metodologije u ovom radu pokušat ću:

- na praktičnom primjeru tehnološke kompanije Backblaze
- nad vrlo svježim, javno objavljenim, skupom otvorenih podataka vrlo velikog opsega na dan 31.12.2017.
- primijeniti faze procesa otkrivanja znanja u podacima – CRISP DM
- analizirajući ciljeve i zahtjeve poslovnog okruženja praktičnog primjera
- analizirajući skup podataka utvrditi pravilnosti
- stvoriti novo znanje i preporuke
- primjenjive u donošenju odluka u svakodnevnom poslovanju podatkovnih centara, tvrtki ali i kućnih korisnika.

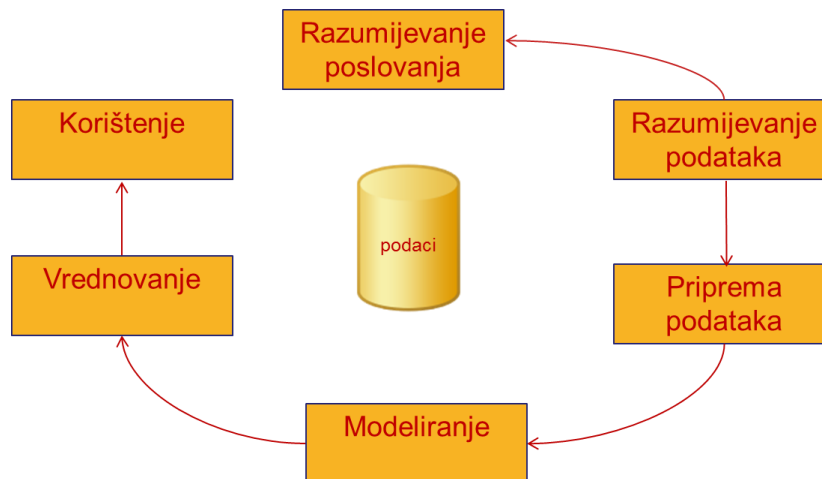
U nastavku, opisat ću osnovne elemente CRISP-DM metodologije te prateći faze procesa otkrivanja znanja u podacima na praktičnom primjeru tehnološke tvrtke Backblaze.

2 Metodologija i tehnike rada

U ovom poglavlju navodim osnovne faze metodologije koja je korištena (praćena) prilikom izrade ovog rada. Opisana metodologija u poglavlju dva sažetak je informacija preuzetih s predavanja doc. dr.sc. Dijana Oreški, te različitih izvora navedenih u popisu literature.

2.1 CRISP-DM metodologija ¹

CRISP-DM je akronim za *cross-industry process for data mining* odnosno predstavlja metodologiju rudarenja podataka. Metodologija je neovisna o industriji, alatima i daje potpuni plan u procesu rudarenja podataka.²



Slika 1: Faze procesa otkrivanja znanja u podacima, (Izvor: dr. Dijana Oreški, FOI, 2017.)

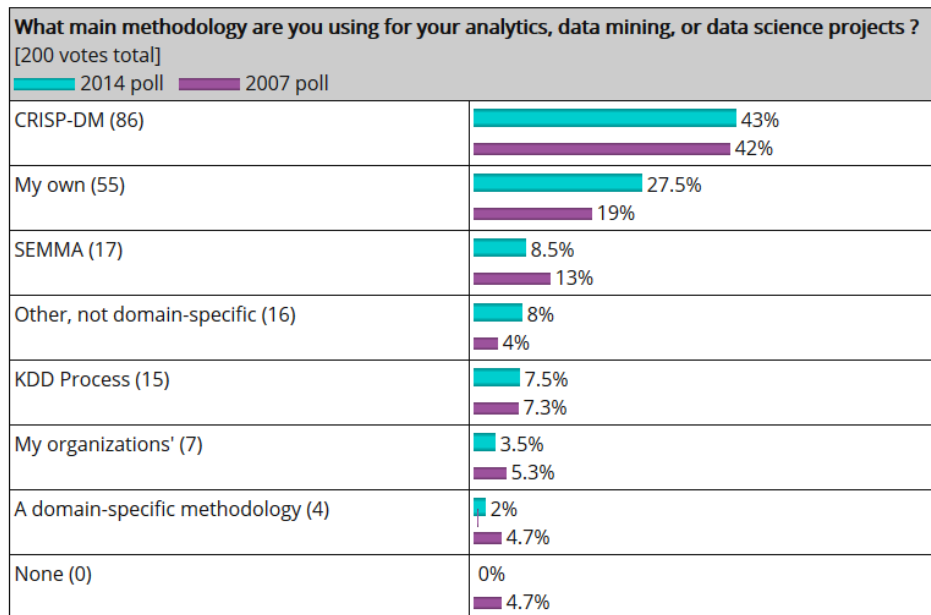
CRISP-DM metodologija opisuje strukturirani pristup planiranju i provedbi poduhvata rudarenje podataka. Idealizirani je to prikaz slijeda događaja (procesu). U praktičnoj primjeni često to nije slučaj već se gore navedeni procesi izvršavaju u drugačijem tijeku (redoslijedu). Često se iterativno vraćamo na prethodne korake ili ih ponavljamo. Robusna je to i dokazana metodologija. Odabrana je s obzirom na zastupljenost i preporuku mentorice.

¹ Wirth, R., & Hipp, J. (2000, April). CRISP-DM: Towards a standard process model for data mining. In *Proceedings of the 4th international conference on the practical applications of knowledge discovery and data mining* (pp. 29-39). Citeseer.,

<https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.198.5133&rep=rep1&type=pdf>, preuzeto 10. kolovoza 2018.

Prilikom odabira metodologije razmatrana je i SEMMA metodologija³.

Pokazuju to i nedavna anketiranja pa tako jedno od istraživanja na relevantnom portalu/trazilici www.kdnuggets.com iz 2014. gdje je postavljeno sljedeće pitanje: *What main methodology are you using for your analytics, data mining, or data science projects?*



Slika 2: Rezultati ankete i usporedba zastupljenosti CRISP-DM metodologije u 2007. i 2014. (N=2000)⁴

Rezultati ove ankete pokazuju da zajednica koja se bavi tematikom koja pokriva ovo područje najčešće koristi CRISP-DM metodologiju za analizu i rudarenja podataka.

CRISP-DM metodologija opisuje ciklički proces u rudarenju podataka. 6 najčešćih faza (processa) opisuje osnovne korake u procesu rudarenja podataka. Standardizacijom procesa rudarenja podataka osiguravamo pouzdanost i ponovljivost procesa. CRISP-DM metodologija usmjerena je k primjeni u rješavanju poslovnih problema.

³ Izvor: Azevedo, A. I. R. L., & Santos, M. F. (2008). KDD, SEMMA and CRISP-DM: a parallel overview. *IADS-DM*. <http://recipp.ipp.pt/bitstream/10400.22/136/3/KDD-CRISP-SEMMA.pdf>

⁴ Izvor: <http://www.kdnuggets.com/polls/2014/analytics-data-mining-data-science-methodology.html>, prosinac 2018.

2.2 Razumijevanje poslovanja

U ovoj fazi želimo razumjeti što želimo postići sa stajališta poslovnih ciljeva. Cilj ove prve faze je definirati najvažnije faktore koji mogu utjecati na konačni cilj projekta. U ovoj fazi razmatramo sljedeće aktivnosti:

- Definiramo naša očekivanja od projekta;
- Pripremamo projektni plan koji bi trebao sadržavati korake koje je potrebno provesti kako bi bili usmjereni k cilju. Aktivnost uključuje inicijalni odabir tehnika i alata;
- Definiramo kriterije uspješnosti projekta s poslovne perspektive. Kriteriji, idealno, moraju biti jasni, jednostavni i mjerljivi;
- Procjena trenutnog stanja je aktivnost prve faze u kojoj radimo evaluaciju resursa, ograničenja i pretpostavki ali i ostale faktore koje moramo uzeti u obzir prilikom analize cilja i definiranja projektnog plana
- Izrađuje se popis resursa (osobe i stručnjaci, experti, tehnološka osnovica, software i alati, dostupni podaci)
- Definiraju se pretpostavke i ograničenja (dostupnost resursa ali i tehnološka ograničenja poput količine podataka koju je praktično modelirati)
- Izrađuje se popis rizika i planovi oporavka
- Izrađuje se popis terminologije koja je važna za projekt
- Izrađuje se cost-benefit analiza projekta u kojoj se uspoređuju troškovi projektnog poduhvata sa eventualnim benefitima pri završetku projektnog poduhvata
- Definiraju se ciljevi projekta rudarenja podataka - koji su različiti od poslovnih ciljeva! Primjer: Istražimo kako možemo predvidjeti kvar opreme prateći podatke o njegovom radu
- Definiraju se kriteriji uspješnosti projekta s poslovne perspektive
- Definiraju se kriteriji uspješnosti postupka rudarenja podataka (s perspektive rudarenja podataka)
- Izrađuje se projektni plan u obliku popisa aktivnosti, njihovog trajanja, potrebnih resursa, ulaza, izlaza.. Dinamički je to dokument koji će se nakon svake faze analizirati i ažurirati
- Izrađuje se inicijalna procjena tehnika i tehnologija. Vrlo je važno napraviti inicijalnu procjenu tehnika i tehnologija rano u fazi planiranja i projektiranja jer odabir loše ili krivih tehnika i tehnologija može drastično utjecati na cijeli projekt.

2.3 Razumijevanje podataka

Razumijevanje podataka podrazumijeva da prikupimo podatke koje smo popisali u popisu nužnih resursa iz projektnog plana. Aktivnosti u ovoj fazi su:

- Izrada popisa podatkovnih resursa zajedno s njihovim lokacijama, metodama koje su korištene pri sakupljanju, problemima koji su se javili prilikom prikupljanja podataka zajedno s rješenjima problema.
- Izrađuje se opis podataka – format, količina (primjerice broj slogova u tablicama), entiteti u poljima i svi ostali važni elementi nužni za razumijevanje prikupljenih podataka
- Istraživanje podataka koristeći razne tehnike poput pisanja upita, vizualizacije podataka, izrada izvještaja
- izrađuje se izvještaj o istraživanju podataka – uključujući inicijalne spoznaje, pretpostavke i hipoteze ali i procjenu njihovog učinka na cjelokupni poduhvat. Izvještaji idealno sadrže grafove i druge vizualizacije koje prikazuju karakteristike podataka
- Analizira se kvaliteta podataka te se traže odgovori na pitanja poput: Da li je podatkovni set kompletan? Da li je točan? Sadrži li greške i koliko su česte? Imali vrijednosti koje nedostaju? Ako da, kako se identificiraju, gdje se pojavljuju i koliko učestalo?
- Izrađuje se izvješće o kvaliteti podataka. Ukoliko postoje problemi sa kvalitetom podataka potrebno je pronaći i predložiti rješenja.

2.4 Priprema podataka

U ovoj se fazi odabiru podaci koji će se upotrijebiti za analize. Ova faza obuhvaća sljedeće aktivnosti:

- Definiraju se kriteriji za odabir podataka uključujući i relevantnost podataka u našem poduhvatu, kvalitetu podataka, volumen, tipove podataka
- Vršiti se odabir atributa i slogova
- Provodi se čišćenje podataka kojim se podiže kvaliteta podataka korištenih u analizi i procesima rudarenja. Ova aktivnost podrazumijeva i naprednije tehnike poput dopune nepoznatih i nedostajućih podataka modeliranjem.
- Izrada deriviranih atributa – novi atributi konstruirani iz izvornih atributa
- integracija podataka iz različitih baza podataka ili tablica unutar baza

- spajanje podataka – primjerice spajanje dvije ili više tablica koje sadrže različite podatke o istom objektu
- abrogacija podataka – aktivnost koja na temelju različitih tablica daje nove sumirane vrijednosti na temelju podataka iz više tablica, slogova ili čak baza podataka

2.5 Modeliranje podataka

Kao prvu aktivnost u fazi modeliranja odabire se tehnika modeliranja koja će se koristiti. Iako je možda u fazi razumijevanja poslovanja odabran alat, u ovoj se fazi odabire specifična tehnika modeliranja i to primjerice – stabla odlučivanja, neuronske mreže i druge tehnike. Pri tome ako se koristi više tehnika modeliranja sljedeće aktivnosti potrebno je provesti za svaku tehniku modeliranja:

- dokumentiranje stvarne tehnike modeliranja koja se koristi
- opisuju se pretpostavke modeliranja
- stvaraju se procedure i mehanizmi za testiranje kvalitete modela i njegove valjanosti
- izrađuje se plan treniranja, testiranja i evaluacije modela
- izgradnja modela pomoću odabrana alata
- procjenjuje se izrađeni model
- na temelju izrađene procjene modela radi se revizija pretpostavki i parametara, radi se fino podešavanje. Alatom za izradu modela iterativno ponavljamo zadnje aktivnosti tako dugo dok ne pronađemo najbolji model
- izrađuje se dokumentacija za svaku iteraciju i procjenu

2.6 Vrednovanje podataka

U ovoj se fazi vrednuje koliko izrađeni model uistinu pridonosi postavljenim poslovnim ciljevima. Tražimo elemente pomoću kojeg bi mogli donijeti zaključak da izrađeni model ne odgovara postavljenim poslovnim ciljevima. Uobičajeno uključuje sljedeće aktivnosti:

- evaluacija rezultata,
- ocjenjivanje procesa,
- određivanje slijedećih koraka.

2.7 Korištenje podataka

U fazi korištenja izrađuje se plan implementacije modela. Uključuje sljedeće aktivnosti:

- izrada plana implementacije

- izrada plana nadzora i održavanja
- izrada finalnog izvještaja poduhvata rudarenja podataka
- finalna prezentacija rada i projekta
- analizira se provedba projekta te se izrađuje dokumentacija iskustva u kojoj se nastoje pobrojati loše odluke, krivi pristupi, problemi u provedbi projekta i druge spoznaje koje mogu pomoći u nekom budućem poduhvatu rudarenja podataka

3 Razumijevanje poslovanja na primjeru tvrtke Backblaze

3.1 Ciljevi i zahtjevi poslovnog okruženja

Tvrtka Backblaze⁵ pruža usluge po paradigmi PaaS (Platform as a Service). Ta paradigma opisuje usluge koje potencijalnim klijentima omogućavaju da unajme tehnološku osnovicu nužnu za sigurnosno kopiranje podataka na vanjske podatkovne podsustave. Tvrtka Backblaze posvećena je pružanju cjenovno prihvatljive usluge neograničene količine backup-a za cijenu od **\$5 po računalu za neograničenu količinu podataka** koju prenesemo na njihovu platformu.

Poslovni ciljevi tvrtke Backblaze definirani su prema održivosti takvog modela (jeftine i pouzdane platforme za izradu sigurnosnih kopija podataka u oblaku. Pri tome jamče izuzetnu **visoku dostupnost sustava od 99.9%**).

Misija⁶ tvrtke Backblaze je u bilo kojem trenutku jamčiti dostupnost svoje usluge po gore navedenim pogodnostima (niska cijena za neograničeni kapacitet).

Tvrtka Backblaze pruža usluge tehnološke osnovice kao usluge svim zainteresiranim kupcima putem interneta. Kupac putem online sustava ugovori pretplatu na njihovu platformu, preuzme odgovarajući software te svoje računalo putem interneta uključi u Backblaze platformu. Korisnik izrađuje sigurnosnu kopiju podataka sa svog računala na njihov sustav, bez obzira na količinu podataka. Gore navedeni poslovni model inovativan je načne pružanja takve usluge na tržištu no za samu tvrtku takav je poslovni model vrlo zahtjevan. Potrebno je kontinuirano pratiti opterećenje njihova sustava i usprkos utjecajima i smetnjama koji djeluju na poslovanje tvrtke. Jezgra poslovanja tvrtke Backblaze je podatkovni sustav. Podatkovni sustav je klaster podatkovnih centara. U podatkovnim centrima na diskovnim jedinicama, spremljeni su podaci korisnika (kupaca).

Osnovna gradivna jedinica podatkovnih podsustava **su diskovne jedinice (tvrđi diskovi)**. Diskovna jedinica nedjeljiva je komponenta od kojeg je sagrađen sav diskovni

⁵ Backblaze, Inc. je pružatelj usluga pohrane sigurnosnih kopija podataka putem interenta. Osnovana je 2007. godine i posluje iz grada San Mateo, Kalifornija, Sjedinjene Američke Države. (500 Ben Franklin Court, San Mateo, CA 94401, <https://www.backblaze.com/company/about.html>).

⁶ Misija tvrtke backblaze, <https://www.backblaze.com/blog/an-intro-to-backblaze/>, prosinac 2018

kapacitet tvrtke. Krajem 2018. godine tvrtka je u svojim podatkovnim centrima koristila preko 100 000 diskovnih jedinica različitih kapaciteta i različitih proizvođača.⁷

3.2 Backblaze - sukobljeni poslovni ciljevi

Kako bi tvrtka Backblaze mogla ostvariti svoju misiju (visoka dostupnost usluge, jamčene kroz SLA po niskoj cijeni za neograničeni kapacitet) nužno je da pažljivo odabire osnovne gradivne blokove svojeg poslovanja - diskovne jedinice!. Kako bi ostvarili svoju misiju vodstvo mora uspješno balansirati između dva suprotstavljena poslovnih cilja: Nisku cijenu usluge za korisnika i visoku dostupnost. Kako su ta dva cilja (i zašto) suprotstavljena?

Kako su diskovne jedinice osnovni gradivni element poslovanja tvrtke Backblaze cijena same diskovne jedinice mora biti niska kako bi operativni troškovi poslovanja mogli podnijeti za korisnika uslugu po zajamčenoj cijeni. Visoka dostupnost usluge ključan je faktor u poslovanju tvrtke Backblaze. Korisnici (kupci) pouzdaju se u kvalitetu usluge. Korisnici izrađuju sigurnosne kopije svojih podataka na Backblaze platformu te očekuju da će ih u bilo kojem trenutku moći povratiti iz sigurnosnih kopija. Korisniku je gubitak podataka havarija koju si ne može dopustiti. Kada se to pak dogodi, korisnik će se osloniti na potpisani SLA ugovor i pokušati vratiti svoje podatke iz sigurnosnih kopija. Samo kratkotrajna nedostupnost Backblaze usluge može trajno uništiti povjerenje u Backblaze usluge. Upravo je to povjerenje u dostupnost još jedan od ključnih gradivnih elemenata njihova poslovnog modela.

Pouzdaniji tvrdi diskovi skuplji su od manje pouzdanih. Proizvođači koriste različite proizvodne procese i različite tehnologije za izradu diskovnih jedinica. Kvalitetniji materijali, noviji i pouzdaniji tehnološki procesi, nova znanja, patenti, procesi osiguranja i praćenja kvalitete – sve to utječe na kvalitetu finalnog proizvoda a samim time i njegovu cijenu.

Idealno, Backblaze bi svoje podatkovne centre gradio od najpouzdanijih diskovnih jedinica koje su u tom trenutku dostupne na tržištu. Posredno, troškovi nabave nužno bi utjecali na cijenu usluge za krajnjeg kupca.

Idealno, Backblaze bi svoje podatkovne centre gradio od najjeftinijih diskovnih jedinica koje su u tom trenutku dostupne na tržištu. Posredno, niska cijena samih jedinica negativno bi utjecala na dostupnost usluge no još više na same operativne troškove upravljanja podatkovnim centrom. Neispravne diskovne jedinice izuzetno jako utječu na sam poslovni model.

⁷ Izvor: <https://www.backblaze.com/blog/hard-drive-stats-for-2018/>, prosinac 2018.

Ove dvije krajnosti nameću potrebu pronalaženja „zlatne sredine“. Naravno – to nije nimalo trivijalan zadatak. Primjerice, pouzdanost tvrdih diskova nije u uvijek u izravnoj vezi sa cijenom. Tu su mnogi drugi elementi koji utječu. Primjerice, ponuda i potražnja u određenom trenutku, kapacitet diskovnih jedinica, promjena tehnologije (primjerice uvođenje helija umjesto dušika), vremenske nepogode koje uzrokuju poremećaje na tržištu, inovacije u nove proizvodne linije (povećanje pouzdanosti zbog nove proizvodne linije)...

3.3 Primjena metoda rudarenja podataka u svrhu podrške odlučivanju

Tehnološka kompanija Backblaze čiji primjer razmatramo, kao jedan od načina ostvarenja svojih poslovnih ciljeva koristi metode rudarenja podataka. Tvrtka cjelokupno svoje poslovanje bazira na osnovnim elementima – diskovnim jedinicama (mehanički tvrdi disk) uobičajeni u IT industriji i računarstvu. Backblaze u svojim podatkovnim centrima s 31. prosinca 2016. koristi preko 70000 diskovnih jedinica. Svaka diskovna jedinica nadzire se. Za nadzor svake pojedine diskovne jedinice koristi se S.M.A.R.T a često još radi jednostavnije piše i kao SMART a predstavlja akronim za *Self-Monitoring, Analysis and Reporting Technology*. Tehnologija je to koja je ugrađena u same diskovne jedinice. Tehnologija detektira i u svoje registre pamti različite parametre kao rezultata analize i monitoringa a sve u cilju omogućavanja i predviđanja kvara diskovne jedinice. Diskovne jedinice su mehanički uređaji. S vremenom i upotrebom povećava se vjerojatnost kvara bilo koje mehaničkog uređaja u upotrebi pa tako i same diskovne.

Backblaze, svakodnevno, jednom dnevno prikuplja S.M.A.R.T podatke sa svake pojedine diskovne jedinice (za svih njih 70000) te primjenjuje metode rudarenja podataka, analizirajući vremenske serije a u cilju predviđanja ispada (kvara) pojedine diskovne jedinice.

3.4 Definicija pretpostavke i hipoteza

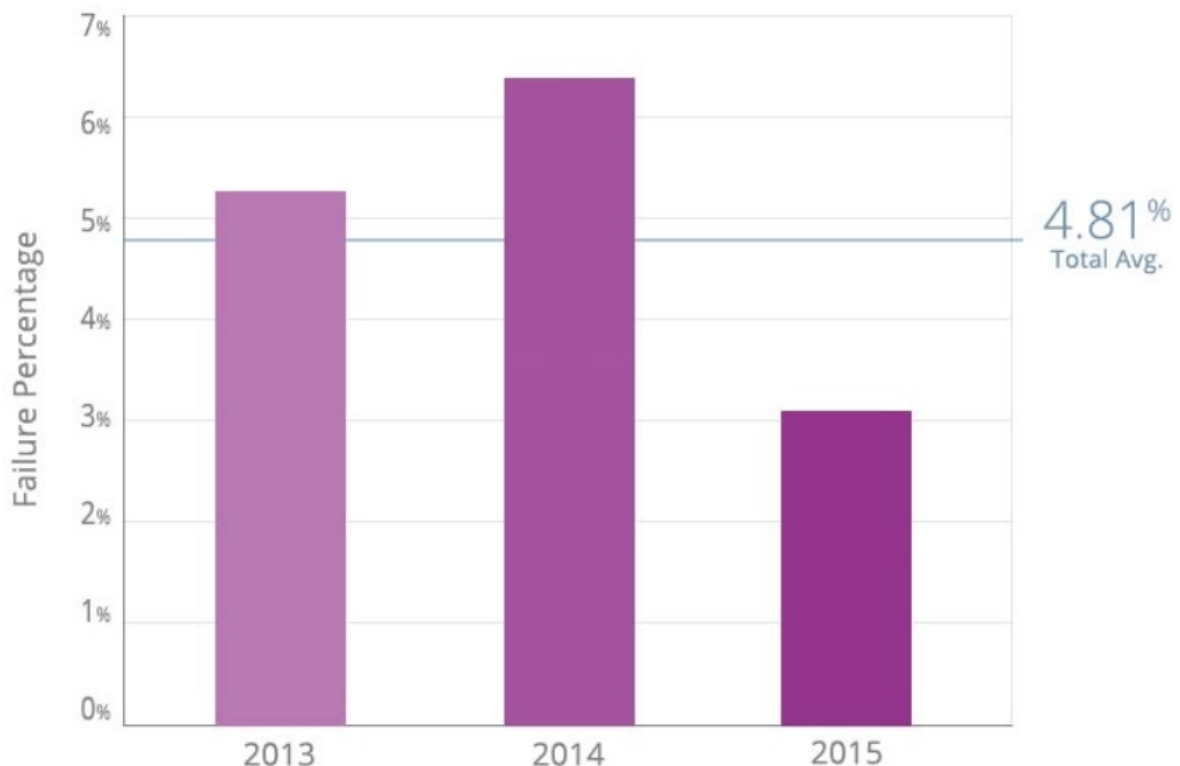
Tvrtka Backblaze počela je prikupljati stanja S.M.A.R.T registara s pojedinih diskovnih jedinice u ožujku 2013. godine. na kraju 2016. godine posjeduje respektabilan set izvornih (raw) podataka nad kojima može provoditi metode rudarenja podataka, analize vremenskih serija te druge metode prikupljanja znanja. Ovo su pretpostavke i hipoteze:

- Prikupljanje SMART podataka omogućit će bolje razumijevanje ponašanja rada diskovnih jedinica.
- Analiza prikupljenih SMART podataka omogućit će predviđanje ispada diskovne jedinice.
- Korištenje metoda rudarenja podataka i primjena stečenog znanja izravno (pozitivno) utječe na uspješnost poslovanja tvrtke (smanjene troškova, povećanje pouzdanosti i dostupnosti, smanjenje rizika i druge pozitivne učinke).

3.5 Kriteriji uspjeha primjene rudarenja podataka u svrhu podrške odlučivanju

Kako bi metode rudarenja podataka smatrane uspješnima nužno je definirati ciljeve takve analize. Backblaze kontinuirano provodi metode rudarenja prikupljenih podataka. Spoznaje i nova znanja prikupljena na taj način primjenjuje u svom poslovanju.

Osnovni kriterij uspjeha primjene rudarenja podataka na našem primjeru je ukupna učestalost kvara pojedine diskovne jedinice (Hard Drive Failure Rate). Znanja dobivena rudarenjem podataka koriste pri donošenju strateških poslovnih odluka a sve u cilju smanjenja učestalost kvara diskovnih jedinica.



Slika 3: Udio ispada diskovnih jedinica za sve diskovne jedinice, svih (proizvođača). 2013. - 2015.⁸

Slika 3. preuzeta je s Internet stranica tehnološke kompanije backblaze koja prikazuje vrlo generalne podatke – udio ispada diskovnih jedinica u odnosu na cijelokupan broj aktivnih diskovnih jedinica u njihovim podatkovnim centrima. Ako se prisjetimo da tvrtka backblaze mna kraju 2018. godine ima preko 100.000 diskovnih jedinica u upotrebi (i aktivnom nadzoru) -

⁸ Izvor: <https://www.backblaze.com/blog/hard-drive-reliability-q3-2015/>, prosinac 2018.

prosječan broj ispada od 4,81% pokazuje da iako su diskovne jedinice (tvrđi diskovi) izuzetno zrela tehnologija razvijana i poboljšavana od 1950.-tih godina – još uvijek gotovo 5000 diskovnih jedinica doživi kvar tijekom predviđenog životnog vijeka trajanja. Analiza kvarova, uvjeta u kojima se oni događaju i njihovih posljedica stalna su tema istraživanja zadnjih 50tak godina. Tvrtka backblaze vrlo transparentno objavljuje korisne podatkovne skupove (S.M.A.R.T) podataka za vrlo velik skup diskovnih jedinica koji rade u vrlo sličnim uvjetima. Ti su podaci vrlo korisni te se rudarenjem mogu proizvesti nova znanja koja olakšavaju donošenje poslovnih odluka.

U nastavku opisana je poslovna okolina, motivacija, načini prikupljanja, struktura objavljenih podataka i mehanizmi koji su korišteni u ovom radu kako bi se preuzeti podaci pročistili preoblikovali, analizirali i vizualizirali u cilju generiranja novog znanja i zaključaka u svrhu donošenja boljih poslovnih odluka.

4 Razumijevanje podataka

U podatkovnom centru tvrtke Backblaze koriste se diskovne jedinice nekoliko proizvođača. Nažalost diskovne jedinice različitih proizvođača ne bilježe sve S.M.A.R.T parametre pa ih u skladu s time nije moguće ni uzorkovati. Ti se podaci razlikuju i među modelima što je razumljivo s obzirom da nove arhitekture diskovnih jedinica uključuju nove i senzore pa je u skladu s time, promatrano kroz vremensko razdoblje, diskovne jedinice novije arhitekture imaju više senzorski prikupljenih podataka u vlastitim S.M.A.R.T registrima.

S druge strane poslovna politika tvrtke Backblaze je takva da na tržištu žele ponuditi neograničenu količinu pohrane podataka za relativno malen novac. Samim time tvrtka Backblaze uobičajeno u svoje podatkovne centre ugrađuje „*consumer grade*“ diskovne jedinice što je termin za prodajnu liniju tvrdih diskova namijenjenih za upotrebu u kućnim i poslovnim računalima široke primjene. Za razliku od *consumer grade* na tržištu postoje i *enterprise grade* linije tvrdih diskova koji su namijenjeni visoko poslovnim primjenama i podatkovnim centrima a ugrađuju se podatkovne podsustave (*storage systems*). *Enterprise grade* diskovne jedinice uobičajeno imaju višu cijenu od *consumer grade* diskovnih jedinica no isto tako trebali bi nuditi veću pouzdanost a na temelju drugačijih proizvodnih procesa, kvalitetnijih materijala i naprednijih arhitektura. Analizom prikupljenih podataka želimo provjeriti da li je to uistinu i stvaran slučaj ili samo mehanizam za prodaju diskovnih jedinica po višim cijenama a bez konkretnog utjecaja na stvarnu pouzdanost samim *enterprise* diskovnih jedinica.

4.1 Okolina u kojoj se podaci prikupljaju

Sve promatrane diskovne jedinice dio su podatkovnog centra u kontroliranim klimatskim uvjetima. Sve diskovne jedinice promatraju se zasebno. Svakodnevno, jednom dnevno tvrtka Backblaze uzrokuje S.M.A.R.T podatke za svaku pojedinu promatranu diskovnu jedinicu (hard disk drive).

4.2 Uzorkovani podaci i njihovo značenje

U našoj analizi promatramo uzorkovani S.M.A.R.T parametre za promatrane diskovne jedinice u vremenskom razdoblju od 1. siječnja 2015. – 30. prosinca 2016., što je ukupno 731 dan.

4.2.1 Struktura analiziranih podataka⁹

- **Date** – datum uzorkovanja SMART vrijednosti
- **Serial Number** – serijski broj uređaja koji mu je dodijelio proizvođač
- **Model** – model uređaja.
- **Capacity** – podatkovni kapacitet diskovne jedinice u bajtovima
- **Failure** – Sadrži vrijednost “0” ako je diskovna jedinica ispravna. Sadrži vrijednost “1” ako je to posljednji dan kako je diskovna jedinica u upotrebi
- **2015 SMART Stats** – 90 stupaca uzorkovanih podataka, i to izvorne (raw) i normalizirane vrijednosti 45 različitih SMART vrijednosti koje prijavljuje diskovna jedinica. Sljedeća tablica¹⁰ pobliže opisuje značenja nekih važniji SMART parametara.

⁹ <https://www.backblaze.com/b2/hard-drive-test-data.html>

Tablica 1: Pojašnjenja pojedinih S.M.A.R.T atributa i parametara¹¹

Critical device status attributes	
Reallocated sectors count	Indicates how many defective sectors were discovered on the drive and remapped using a spare sectors pool. Low values in absence of other fault indications point to the disk surface problem. Raw value indicates the exact number of such sectors.
Current pending sectors count	Indicates how many suspected defective sectors are pending "investigation". These will not necessarily be remapped. In fact, such sectors may be not defective at all (e.g. if some transient condition prevented reading of the sector, it will be marked "pending") - they will be then re-tested by the device off-line scan ¹ procedure and returned to the pool of serviceable sectors. Raw value indicates the exact number of such sectors.
Off-line uncorrectable sectors count	Similar to "Reallocated sectors count". Indicates how many defective sectors were found during the off-line scan ¹ .
Read error rate	Logs rate at which specified events (errors) occur. Lower value indicates more events (errors). Retries are not necessarily indicate a persistent problem, but one should proceed with caution if any of these attributes is degraded.
Read error retry rate	
Write error rate	
Seek error rate	
Recalibration retries	Indicates how often the drive is unable to recalibrate at the first attempt. Raw value may show the exact number of recalibration events (at least with some vendors) but this should be taken with a grain of salt.
Spin up time	Low value indicates that a drive takes longer than expected to spin up to its rated speed. Might indicate either a controller or a spindle bearing problem.
Spin retry count	Spin retry event is logged each time the drive was unable to spin its platters up to the rated rotation speed in the due time. Spin-up attempt was aborted and retried. This typically indicates severe controller or bearing problem, but may be sometimes caused by power supply problems.
Drive lifetime information	
Drive start/stop count	These two provide the estimation of the drive wear. Vendor estimates the supposed device lifetime and the number of cycles. The value for these attributes is then computed based on this estimation. The T.E.C. condition with one of these attributes does not necessarily indicate a drive failure, but rather suggests that a drive should be considered unreliable due to the wear and tear. Raw values are typically just the count of events.
Power off/retract cycle count	
Power on hours count	
Head flying hours count	
Normalized values are computed similar to the above. Despite what the name suggests, the raw value of the attribute is stored using all sorts of measurement units (hours, half-hours, or ten-minute intervals to name a few) depending on the manufacturer of the device.	
Operating conditions information	
Temperature	Indicates the device temperature, if the appropriate sensor is fitted. Lowest byte of the raw value contains the exact temperature value (Celsius degrees).
Ultra DMA CRC error rate	Low value of this attribute typically indicates that something is wrong with the connectors and/or cables. Disk-to-host transfers are protected by CRC error detection code when Ultra-DMA 66 or 100 is used. So if the data gets garbled between the disk and the host machine, the receiving controller senses this and the retransmission is initiated. Such a situation is called "UDMA CRC error". Once the problem is rectified (typically by replacing a cable), the attribute value returns to the normal levels pretty quick.
G-sense error rate	Indicates if the errors are occurring attributed to the drive shocking (either due to the environmental factors or due to improper installation). The hard drive must be fitted with the appropriate sensor to get information about the G-loads. This attribute is mainly limited to the notebook (2.5") drives. Once the operation conditions are corrected, the attribute value will return to normal.

¹¹ <https://www.z-a-recovery.com/manual/smart.aspx>

5 Priprema podataka

5.1 Kreiranje podataka i priprema baze podataka

Kako bi mogli analizirati uzorkovane podatke najprije ih treba prikupiti. U našoj analizi i traženju znanja u podacima koristiti će se podaci koje tvrtka Backblaze objavljuje na svojim internet stranicama¹². Backblaze je dobrim dijelom već pripremio strukturirane podatke no za daljinu analizu potrebno je podatke dodatno pripremiti. Podatke dostupne na njihovim internet stranicama potrebno je proveli smo kroz nekoliko koraka kako bi ih pripremili na daljinu obradu:

- Podaci su u izvornom obliku preuzeti s poslužitelja tvrtke Backblaze i spremljeni na korisničko računalo
- Izvorne računalne datoteke tvrtka Backblaze radi bržeg preuzimanja i komprimirane su te ih je potrebno dekomprimirati nakon preuzimanja
- Podaci su objavljeni po kvartalima, no za potrebe analize u ovom seminaru podaci su objedinjeni u godišnje
- Analizu podataka nije moguće raditi na računalnim datotekama u prikupljenom obliku (datoteke su preuzete u obliku gdje su diskretne vrijednosti odvojene zarezom – CSV format zapisa)
- **Za potrebe analize podataka pripremljena je zasebna poslužiteljska tehnološka osnovica sa snažnim procesnim i podatkovnim podsustavima na virtualizacijskoj platformi.**
- Preuzeti podaci translaterani si u sustav za upravljanje relacijskom bazom podataka - MySQL
- Podaci su spremni za analize i rudarenje.

¹² <https://www.backblaze.com/b2/hard-drive-test-data.html>

6 Čišćenje podataka

6.1 Nedostatak uzorkovanih podataka – prazna polja

Consumer grade diskovne jedinice često ne prikupljaju sve podatke u svojim S.M.A.R.T registrima. Posljedično, tih podataka u uzorkovanim podacima nema. Uzorkuju se podaci koje tvrdi diskovi sadrže u svojim registrima. Iz tih razloga svaki uzorak ima i praznih polja.

6.2 Nekonzistentna polja i nejasne vrijednosti

Proizvođači diskovnih jedinica ne objavljuju prava značenja vrijednosti S.M.A.R.T registara. Prikupljeni podaci tako mogu imati različita značenja i ta značenja mogu varirati među pojedinim proizvođačima i modelima diskovnih jedinica.

6.3 Vrijednosti izvan očekivanih granica

Uzorkovane vrijednosti su vrijednosti koje prijavljuju sami diskovi. Neki od uzoraka nose vrijednosti izvan očekivanih granica. Primjerice S.M.A.R.T vrijednost (SMART9) uobičajeno predstavlja vrijeme rada tvrdog diska u satima. Primijećeni su uzorci sa vrijednostima koji bi značili da su pojedine diskovne jedinice u upotrebi više od desetak godina, što naravno nije moguće.

7 Redukcija podataka

7.1 Redukcija diskovnih jedinica za operativne sustave

Na dan 31. prosinca 2016. godine Backblaze podatkovni centar koristio je 73653 diskovne jedinice¹³. Diskovne jedinice se Backblaze-u kao i u drugim podatkovnim centrima koriste u dvije uloge:

- diskovi čija je dedikirana namjena smještaj operativnih sustava te pokretanje operativnih sustava (eng: *boot drives*)
- diskovi čija je dedikirana namjena smještaj podataka (eng: *data drives*)

Krajem 2016. godine Backblaze podatkovni centar koristio je 1553 diskovne jedinice za pokretanje operativnih sustava (*boot drives*) i 72100 diskovne jedinice za smještaj podataka (*data drives*). Ove dvije kategorije diskova posebno su izdvojene te se zbog vrlo različitog načina korištenja a posljedično tome i utjecaja načina korištenja na njihovu pouzdanost, broj kvarova, vrijeme korištenja – promatraju zasebno. U ovom seminaru bit će izolirane sve diskovne jedinice iz kategorije *boot drives*. Kriterij za selekciju je kapacitet diskovne jedinice. Sve diskovne jedinice manje od 1.5TB smatraju se diskovnim jedinicama namijenjenim podizanju operativnih sustava i kategorizirani su kao *boot drives* te ih nećemo analizirati. Ukupna smo reducirali 1553 promatrane diskovne jedinice.

7.2 Redukcija diskovnih jedinica koje imaju premali statistički uzorak

Iz promatranog uzorka isključit ćemo sve modele diskovnih jedinica čiji je ukupan broj manji od 45 po modelu. Ovako malen broj diskovnih jedinica nema značajnijeg utjecaja na procese analize, ne dodaje dodatnu vrijednost točnosti podataka dok s druge strane može navesti na krive zaključke na temelju premalog statističkog uzroka.

Ukupno smo reducirali 161 tvrdi disk. Ukupni uzorak promatranih diskovnih jedinica sada je 71939.

¹³ <https://www.backblaze.com/blog/hard-drive-benchmark-stats-2016/>

Dnevna učestalost kvara je na ovom primjeru 2%. Da bi dobili godišnju učestalost kvara, pomnožimo tu vrijednost sa 365 (broj dana u godini). Kada bi se naš tjedan iz gornjeg primjera nastavio, godišnja učestalost kvara bila bi 730%.

Godišnja učestalost kvara od preko 100% je naravno moguća. Naime, ako pretpostavimo da u našem podatkovnom centru koristimo 100 diskovnih jedinica. Ako je dnevna učestalost kvara 2% to bi značilo da dnevno ispadnu 2 diska. Ako diskove redovito mijenjamo odmah po ispadu, tijekom godine zamijenili bi 730 diskova.

Analiza ispada diskova u podacima svodi se na brojenje disk-dana i brojenja „ispada“. U podacima je ispad označen vrijednošću „1“ u stupcu *failure*. Jednom kad je disk označen oznakom 1, više podaci iz njegovih SMART registara više se ne uzorkuju. Dan kada je vrijednost polja *failure* = 1, posljednji je dan te diskovne jedinice u upotrebi.

Kako su podaci pripremljeni u relacijskoj bazi podataka možemo jednostavnim SQL upitima dobiti nove informacije:

```
--  
-- Stvorimo tablicu koja sadrži broj disk-dana za svaki model  
-- to je zapravo broj slogova u tablici drive_stats_2016_q4  
-- za taj model  
--  
CREATE TABLE drive_days_2016_q4 AS  
  SELECT model, count(*) AS drive_days_2016_q4  
  FROM drive_stats_2016_q4  
  GROUP BY model;
```

Ako pogledamo dobivene podatke možemo vidjeti koliko disk-dana je određeni model u upotrebi. Ovi nam podaci zapravo još ništa ne govore jer je različit broj diskovnih jedinica određenog modela u upotrebi. Kreirajmo tablicu koja sadrži broj ispada diskovnih jedinica u promatranom razdoblju (Q4 2016).

```
CREATE TABLE failures_2016_q4 AS  
  SELECT model, count(*) AS failures_2016_q4  
  FROM drive_stats_2016_q4  
  WHERE failure = 1  
  GROUP BY model;
```

Spojimo li tablice zajedno i:

- izračunamo godišnju učestalost kvara
- "disk-godina" izračunamo kao broj disk-dana podijeljeno sa 365
- godišnja učestalost kvara je jednostavno broj kvarova podijeljena sa disk-godina
- rezultat pomnožimo sa 100 kako bi dobili postotak
- promatramo samo modele čiji je broj diskovnih jedinica veći od 45

```
CREATE TABLE failure_rates_2016_q4 AS
SELECT drive_days_2016_q4.model AS model,
       drive_days_2016_q4.drive_days_2016_q4 AS drive_days_2016_q4,
       failures_2016_q4.failures_2016_q4 AS failures_2016_q4,
       100.0 * (1.0 * failures_2016_q4) / (drive_days_2016_q4 / 365.0) AS annual_failure_rate
FROM drive_days_2016_q4, failures_2016_q4, model_count_2016_q4
WHERE drive_days_2016_q4.model = failures_2016_q4.model
      AND model_count_2016_q4.model = failures_2016_q4.model
      AND 45 <= model_count_2016_q4.count
ORDER BY model;
```

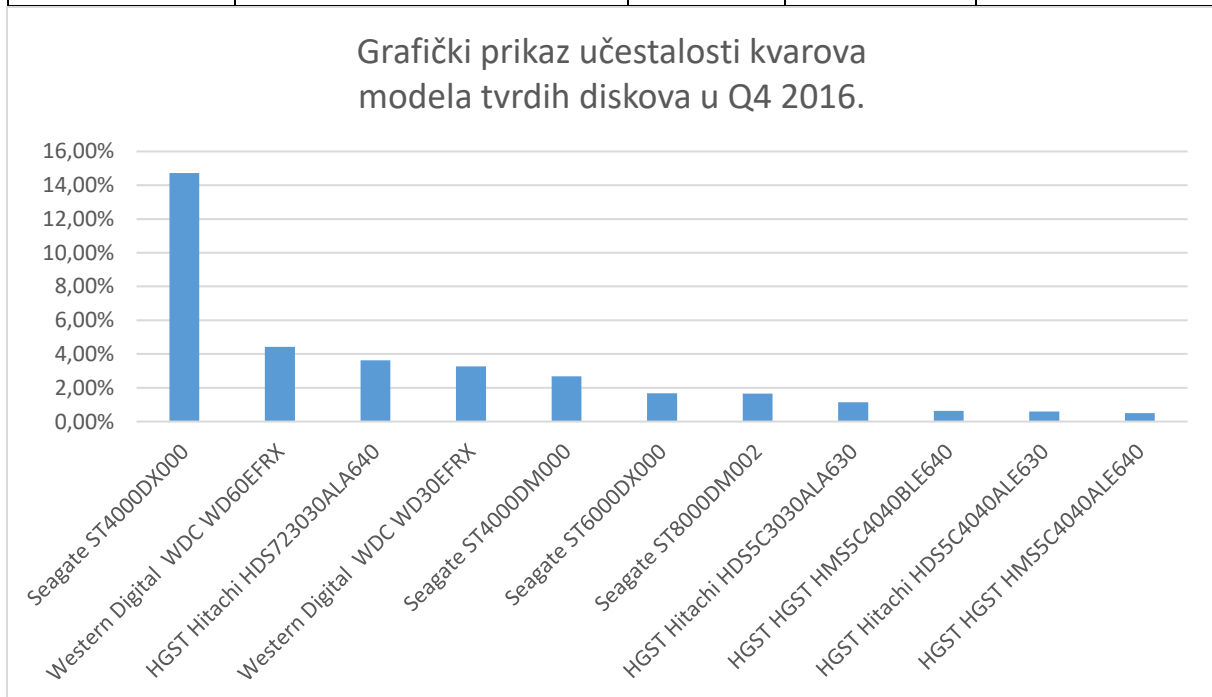
9 Razumijevanje i vizualizacija podataka

9.1 Posljednji kvartal 2016. godine

Izvršimo li pripremljene upite nad našim podacima koje smo pripremili i smjestili u relacijsku bazu podataka kreirat će se tablica sa pregledom učestalosti kvara pojedinih modela diskovnih jedinica. Pogledajmo kava je pouzdanost diskovnih jedinica prema modelima u četvrtom kvartalu 2016. godine.

Tablica 2: Učestalost kvarova diskovnih jedinica prema modelu u Q4, 2016.

Proizvođač	Model	disk-dana	broj-kvarova	učestalost kvara
Seagate	ST4000DX000	17354	7	14,72%
Western Digital	WDC WD60EFRX	41304	5	4,42%
HGST	Hitachi HDS723030ALA640	90415	9	3,63%
Western Digital	WDC WD30EFRX	100259	9	3,28%
Seagate	ST4000DM000	3196552	234	2,67%
Seagate	ST6000DX000	173720	8	1,68%
Seagate	ST8000DM002	663697	30	1,65%
HGST	Hitachi HDS5C3030ALA630	412752	13	1,15%
HGST	HGST HMS5C4040BLE640	809119	14	0,63%
HGST	Hitachi HDS5C4040ALE630	241665	4	0,60%
HGST	HGST HMS5C4040ALE640	648393	9	0,51%



Slika 4: Grafički prikaz učestalost kvarova modela tvrdih diskova u Q4 2016.

9.1.1 Rezultat analize i zaključak obrađenih podataka

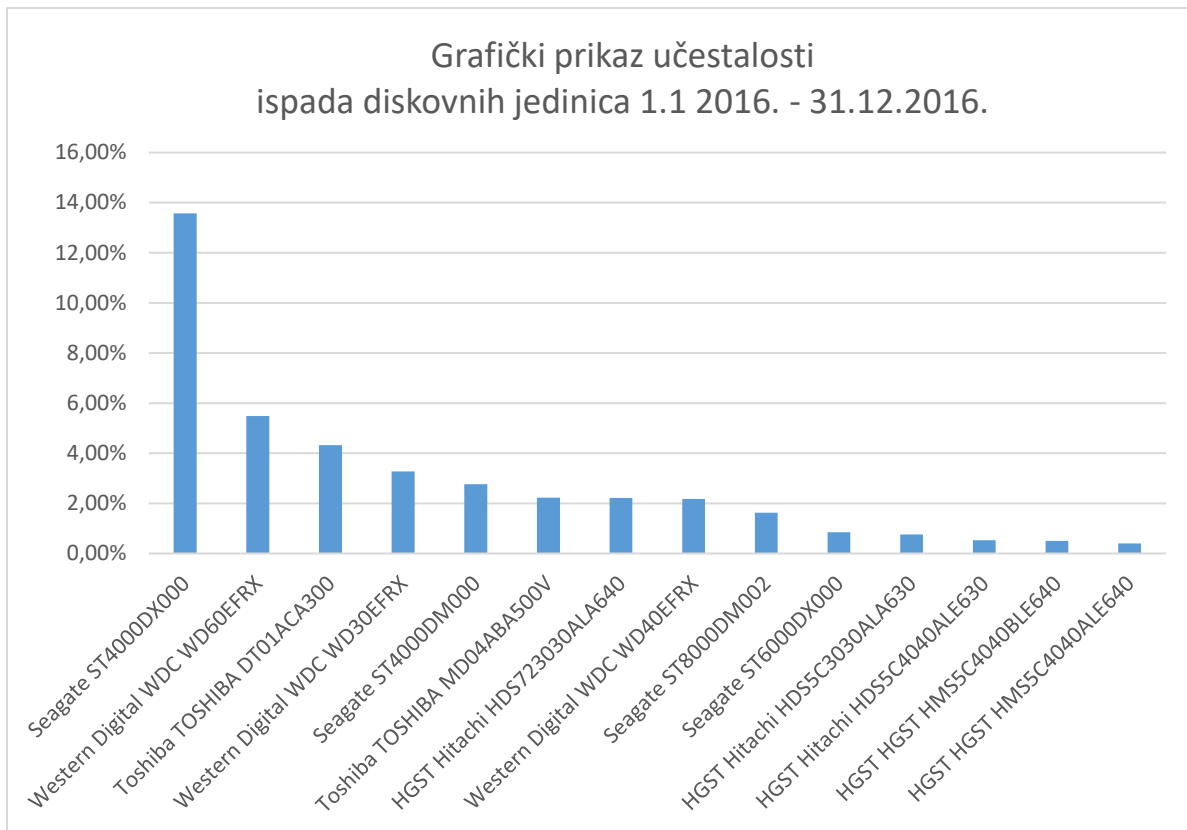
Analizom podataka iz prikupljenih podataka možemo zaključiti da je u promatranom razdoblju (posljednji kvartal 2016. godine) najpouzdaniji proizvođač tvrtka HGST s njihovim modelom HMS5C4040ALE640. Prilikom kupovine, pametno bi bilo izbjeći proizvođače Seagate i Western Digital jer su njihovi diskovi skloni kvarovima (imaju najveću učestalost kvara).

9.2 Dulji vremenski termin – cijela 2016. godina

Kako bi bili sigurni da posljednji kvartal 2016. godine nije previše kratak termin za analizu proširimo opseg podataka na cijelu 2016. godinu. Pripremimo li podatke jednako kao i za posljednji kvartal 2016. godine, izvršimo li prilagođene upite na našem bazom podataka dobijemo sljedeće podatke:

Tablica 3: Učestalost ispada diskovnih jedinica 1.1 2016. - 31.12.201

Proizvođač	Model	disk-dana	broj-kvarova	učestalost kvara
Seagate	ST4000DX000	72615	27	13,57%
Western Digital	WDC WD60EFRX	166152	25	5,49%
Toshiba	TOSHIBA DT01ACA300	16900	2	4,32%
Western Digital	WDC WD30EFRX	390379	35	3,27%
Seagate	ST4000DM000	12359750	938	2,77%
Toshiba	TOSHIBA MD04ABA500V	16425	1	2,22%
HGST	Hitachi HDS723030ALA640	361937	22	2,22%
Western Digital	WDC WD40EFRX	16790	1	2,17%
Seagate	ST8000DM002	1075720	48	1,63%
Seagate	ST6000DX000	684840	16	0,85%
HGST	Hitachi HDS5C3030ALA630	1647137	34	0,75%
HGST	Hitachi HDS5C4040ALE630	960249	14	0,53%
HGST	HGST HMS5C4040BLE640	2436130	34	0,51%
HGST	HGST HMS5C4040ALE640	2579698	28	0,40%



Slika 5: Grafički prikaz učestalosti ispada diskovnih jedinica 1.1 2016. - 31.12.2016.

9.2.1 Rezultat analize i zaključak obrađenih podataka

Analizom podataka iz prikupljenih podataka možemo zaključiti da je u promatranom razdoblju (cijela 2016. godina) najpouzdaniji proizvođač tvrtka HGST s nekoliko njihovih modela. Prilikom kupovine, pametno bi bilo izbjeći proizvođače Segate i Western Digital jer su njihovi diskovi skloni kvarovima (imaju najveću učestalost kvara).

Dodatno, možemo vidjeti da iako tvrtka Segate ima vrlo lošu reputaciju, to nije slučaj s njihovim modelima ST8000DM002 i ST6000DX000 koji imaju vrlo malenu učestalost kvara.

9.3 Dulji vremenski termin (1. 1. 2015. – 31. 12. 2016.)

Proširimo promatrani vremenski interval na posljednje dvije godine. Pripremimo podatke jednako kao i do sada. Izvršimo li prilagođene upite na našem bazom podataka dobijemo sljedeći rezultat upita:

Tablica 4: Učestalost ispada diskovnih jedinica 1.1 2015. - 31.12.2016

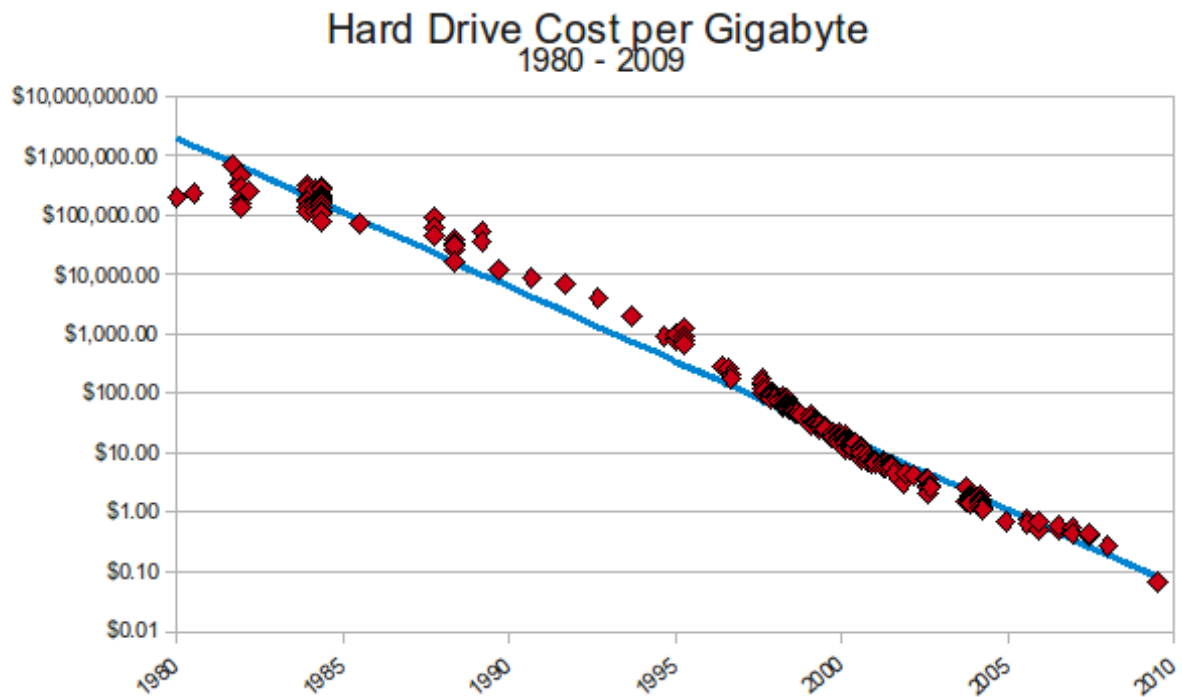
Proizvođač	Model	disk-dana	broj-kvarova	učestalost kvara
Western Digital	WDC WD10EADX	3405	3	32,16%
Western Digital	WDC WD10EADS	98755	34	12,57%
Seagate	ST31500541AS	379791	112	10,76%
Seagate	ST4000DX000	140781	34	8,82%
Western Digital	WDC WD10EACS	13995	3	7,82%
Western Digital	WDC WD60EFRX	328942	52	5,77%
HGST	Hitachi HDS723030BLE640	6557	1	5,57%
Western Digital	WDC WD30EFRX	759657	114	5,48%
Seagate	ST33000651AS	52283	5	3,49%
Western Digital	WDC WD40EFRX	33184	3	3,30%
Toshiba	TOSHIBA DT01ACA300	34104	3	3,21%
Seagate	ST4000DM000	19142515	1524	2,91%
Seagate	HGST HDS5C4040ALE630	38394	3	2,85%
Western Digital	HGST HDS724040ALE640	29476	2	2,48%
Toshiba	TOSHIBA MD04ABA500V	30015	2	2,43%
HGST	HGST HUH728080ALE600	30958	2	2,36%
HGST	Hitachi HDS722020ALA330	2598806	151	2,12%
HGST	Hitachi HDS723030ALA640	727463	40	2,01%
Seagate	ST8000DM002	1075720	48	1,63%
Seagate	ST6000DX000	1195196	47	1,44%
Toshiba	TOSHIBA MD04ABA400V	88972	3	1,23%
HGST	Hitachi HDS5C3030ALA630	3308580	81	0,89%
HGST	Hitachi HDS5C4040ALE630	1919936	35	0,67%
HGST	HGST HMS5C4040ALE640	5153686	72	0,51%
HGST	HGST HMS5C4040BLE640	3561977	46	0,47%

10 Problem dva različita kraja životnog vijeka diskovne jedinice

Tvrtka Backblaze u podacima koje analiziramo objavljuje vrijednosti SMART parametara, svaki dan za svaku pojedinu diskovnu jedinicu u upotrebi u svojim podatkovnim centrima. Podaci se počinju prikupljati prvi dan nakon što je diskovna jedinica uključena u rad podatkovnog centra. Prilikom uvođenja nove diskovne jedinice u rad uključuju se i *smartmon* alati koji u ukupan set podataka dodaju i dnevne vrijednosti SMART parametara za novo dodane jedinice te se ti podaci prikupljaju do kraja životnog vijeka diskovne jedinice koja može biti jedno od sljedećeg:

- Ispad (kvar) diskovne jedinice radi kojeg se diskovna jedinica trajno isključuje iz sustava
- Isključenje ispravne diskovne jedinice iz sustava iz drugih (komercijalnih) razloga

Za razliku od ispada (kvara) diskovne jedinice koji je (iako uobičajen i do neke mjere očekivan) ipak kritičan, neželjen događaj tvrtka Backblaze često isključuje diskovne jedinice različitih modela iz samog sustava iz različitih komercijalnih razloga koji nisu povezani s njihovom pouzdanošću. To su primjerice, situacije kada se radi uštede fizičkog prostora koji ima svoju cijenu namjerno isključuju diskovne jedinice malog kapaciteta te se na njihovo mjesto ugrađuju diskovne jedinice većeg kapaciteta. Svaki slot u „podu“ (odnosno šasiji) zauzima određen fizički prostor koji u sustavu velikog broja pojedinih diskovnih jedinica ima svoju cijenu. Napretkom tehnologije i razvojem sustava diskovne pohrane gustoća zapisa na diskovne jedinice raste te je razumna poslovna odluka da na istom fizičkom prostoru zamjenom diskovnih jedinica malog kapaciteta s jedinicama većeg kapaciteta – smanjujemo ukupnu cijenu koštanja jedinice kapaciteta¹⁵. Isto tako to može biti i manja potrošnja energije novih generacija diskovnih jedinica ili pak preventivna zamjena cijele serije nekih modela.



Slika 7: Graf kretanja cijene 1 gigabayta prostora 1980.-2015.¹⁶

U kontekstu istraživanja ovog rada, važno je razumjeti da se u analiziranim podacima nalaze diskovne jedinice čiji su se SMART podaci prikupljali do trenutka kada je diskovna jedinica isključena iz sustava iz bilo kojeg od ova dva razloga. Potrebno voditi računa da zbog diskovnih jedinica isključenih iz komercijalnih razloga ne dobijemo iskrivljenu ili čak netočnu analizu pouzdanosti i učestalosti kvarova.

¹⁶ <http://www.mkomo.com/cost-per-gigabyte>, prosinac, 2018.

10.1 Kaplan-Meierova krivulja doživljenja - primjena medicinskih statističkih metoda nad promatranim podacima

Kaplan-Meierova krivulja preživljavanja česti je model „vrijeme do događaja“ kojim se najviše u biomedicinskim istraživanjima ali i u drugim primjenama pomoću Kaplan-Meierovih krivulja (grafova) vizualiziraju (primjerice) dva različita liječenja i njihov utjecaj na preživljavanje pacijenata u odnosu na vrijeme do događaja. Događaj može biti primjerice vrijeme do pojave smrti, vrijeme remisije nakon terapije ili pojave primarne bolesti,

U nastavku ovog istraživačkog rada primjeni ćemo isti pristup (principe analize) koji se upotrebljavaju u biomedicina nad našim promatranim podacima. Motivacija za primjenu upravo ovog pristupa je izražena sličnost u modelu ponašanja.

U medicinskim istraživanjima Kaplan-Meierove krivulje koriste se pri vizualizaciji dvaju različita liječenja (ili dvije različite skupine bolesnika) te se želi prikazati vjerojatnost preživljavanja. Kaplan-Meierova princip desne cenzure govori o pacijentima koji su napustili klinička istraživanja iz nekih nepoznatih razloga, odselili su, izgubljen je kontakt i slično... zapravo njihova stvarna i konačna situacija nije poznata te iako su u jednom trenutku pacijenti bili dio istraživanja zbog navedenih razloga nije poznat konačan rezultat liječenja.

U našem slučaju možemo primijetiti vrlo sličan model ponašanja u analizi problema dva različita kraja životna vijeka neke diskovne jedinice koja je isključena iz sustava. Tako za sve diskovne jedinice koje su iz sustava isključene iz komercijalnih razloga nakon trenutka isključenja više nema slogova sa SMART parametrima no znamo da su jedinice uredno radile do tog trenutka.

10.2 Alati za analizu podataka - „R“

R je specifičan programski jezik prilagođen eksploratornoj, statističkoj i dubinskoj analizi podatkovnih skupova. Po svojoj prirodi nalazi se na razmeđu između klasičnih programskih jezika kao što su Python, Java ili C++ i statističkih alata kao što su SAS ili SPSS. Uz interaktivni pristup, ali i mogućnost pisanja složenijih programskih skripti, R se danas nametnuo kao jedan od vodećih analitičkih programskih jezika s kojim se uz pomoć pratećih paketa na vrlo učinkovit način mogu provesti složene analize podatkovnih skupova te stvoriti izvještaji popraćeni kompleksnim vizualizacijama i izračunima. Svladavanje jezika R zahtjeva specifičnu kombinaciju programskih vještina, poznavanja osnova statistike, ali i izvjesnu kreativnost i spremnost na izazove.

10.3 Konstrukcija Kaplan-Meierovih krivulja u alatu „R“

Kako bi konstruirali navedene krivulje na nad našom relacijskom bazom u kojem se nalaze statistički podaci od interesa konstruiramo takav upit koji će nam stvoriti novu relacijsku tablicu s minimalnim nužnim podacima koje ćemo unijeti u programski jezik „R“.

Zbog obima prikupljenih podataka (samo za 2016. godinu u našoj bazi podataka postoji više od 22 milijuna zapisa) nužno je maksimalno smanjiti opseg upita i definirati samo nužno potrebne entitete kako bi u okviru raspoloživih sistemskih resursa mogli pripremiti podatke.

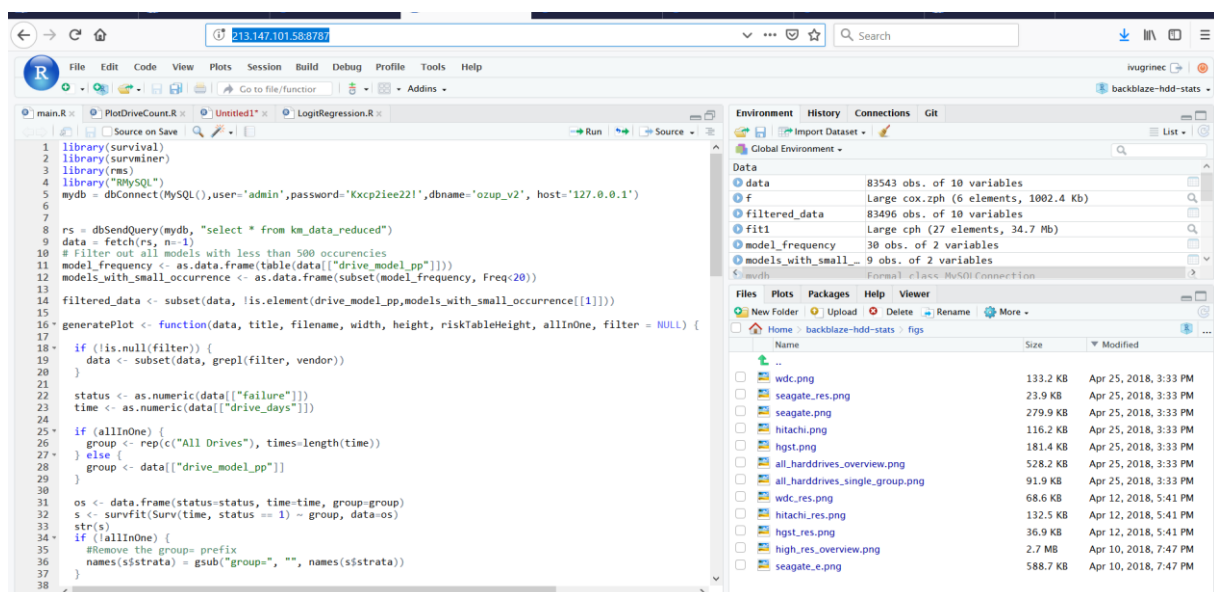
Sama priprema podataka poseban je izazov koji zahtjeva pripremu upita nad malim (reprezentativnim) skupom podataka te njihovo testiranje. Tek nakon što se utvrdi da upit na dovoljno dobar način konstruira tablicu s podacima za pripremu grafa ima smisla proširiti upit nad mnogo većim setom podatka (više godina)

Upit koji će nam pripremiti podatke za konstrukciju Kaplan Meierovih krivulja doživljenja za 2016. godinu glasi:

```
-- Upit za pripremu podataka nužnih za konstrukciju Kaplan-Meierove
-- krivulje preživljavanja (podaci za 2016. godinu)

SELECT serial_number, COUNT(*) drive_days, model, failure FROM
ozup_v2.drive_stats_2016 GROUP BY serial_number
```

Upotrebom programskog alata „R“ kojeg smo implementirali na našu tehnološku osnovicu i javno objavili na adresi: <http://213.147.101.58:8787>



Slika 8: "R" studio na linux poslužitelju za potrebe stvaranja R upita nad promatranim podacima

Tablica 5: Tablica modula korištenih u stvaranju upita u programskom jeziku „R“

Ime modula	Opis modula:
Survminer Drawing Survival Curves using 'ggplot2'	Contains the function 'ggsurvplot()' for drawing easily beautiful and 'ready-to-publish' survival curves with the 'number at risk' table and 'censoring count plot'. Other functions are also available to plot adjusted curves for 'Cox' model and to visually examine 'Cox' model assumptions. https://cran.r-project.org/web/packages/survminer/index.html
survival: Survival Analysis	Contains the core survival analysis routines, including definition of Surv objects, Kaplan-Meier and Aalen-Johansen (multi-state) curves, Cox models, and parametric accelerated failure time models. https://cran.r-project.org/web/packages/survival/index.html

10.3.1 Izvorni kod izrade Kaplan-meirovih krivulja doživljenja u programskom jeziku „R“

```

1 library(survival)
2 library(survminer)
3 library(rms)
4 library("RMySQL")
5 mydb = dbConnect(MySQL(),user='admin',password='Kxcp2iee22!',dbname='ozup_v2', host='127.0.0.1')
6
7
8 rs = dbSendQuery(mydb, "select * from km_data_reduced")
9 data = fetch(rs, n=-1)
10 # Filter out all models with less than 500 occurrences
11 model_frequency <- as.data.frame(table(data[["drive_model_pp"]]))
12 models_with_small_occurrence <- as.data.frame(subset(model_frequency, Freq<20))
13
14 filtered_data <- subset(data, !is.element(drive_model_pp,models_with_small_occurrence[[1]]))
15
16 * generatePlot <- function(data, title, filename, width, height, riskTableHeight, allInOne, filter = NULL) {
17
18   if (!is.null(filter)) {
19     data <- subset(data, grepl(filter, vendor))
20   }
21
22   status <- as.numeric(data[["failure"]])
23   time <- as.numeric(data[["drive_days"]])
24
25 * if (allInOne) {
26   group <- rep(c("All Drives"), times=length(time))
27 * } else {
28   group <- data[["drive_model_pp"]]
29 * }
30
31 os <- data.frame(status=status, time=time, group=group)
32 s <- survfit(Surv(time, status == 1) ~ group, data=os)
33 str(s)
34 * if (!allInOne) {
35   #Remove the group= prefix
36   names(s$strata) = gsub("group=", "", names(s$strata))
37 * }
38

```

Slika 9: Izvorni kod programa za izradu krivulja u programskom jeziku "R"

```

library(survival)
library(survminer)
library(rms)
library("RMySQL")
mydb =
dbConnect(MySQL(), user='admin', password='*****', dbname='ozup_v2',
host='127.0.0.1')

rs = dbSendQuery(mydb, "select * from km_data_reduced")
data = fetch(rs, n=-1)
# Filter out all models with less than 500 occurrences
model_frequency <- as.data.frame(table(data[["drive_model_pp"]]))
models_with_small_occurrence <- as.data.frame(subset(model_frequency,
Freq<20))

filtered_data <- subset(data,
!is.element(drive_model_pp, models_with_small_occurrence[[1]]))

generatePlot <- function(data, title, filename, width, height,
riskTableHeight, allInOne, filter = NULL) {

  if (!is.null(filter)) {
    data <- subset(data, grepl(filter, vendor))
  }
  status <- as.numeric(data[["failure"]])
  time <- as.numeric(data[["drive_days"]])
  if (allInOne) {
    group <- rep(c("All Drives"), times=length(time))
  } else {
    group <- data[["drive_model_pp"]]
  }
  os <- data.frame(status=status, time=time, group=group)
  s <- survfit(Surv(time, status == 1) ~ group, data=os)
  str(s)
  if (!allInOne) {
    #Remove the group= prefix
    names(s$strata) = gsub("group=", "", names(s$strata))
  }
  plot <- ggsurvplot(s,
                    data=os,
                    censor=FALSE,
                    pval=!allInOne,

```



```

        ##main = title,
        ylab = "Vjerojatnost doživljenja",
        xlab = "Proteklo vrijeme od instalacije (dana)",
        conf.int = FALSE,
        fontsize=3,
        break.time.by = 180,
        pval.size=3,
        size=0.4,
        ylim = c(0.85, 1),

        break.y.by=0.05,

        risk.table = TRUE, # Add risk table
        risk.table.title = "",
        risk.table.height = riskTableHeight,

        ggtheme = theme_light()
    )

plot$table <- plot$table + xlab("Broj diskovnih jedinica u upotrebi nakon
X dana") + ylab(NULL) # Remove the labels from the table
  plot$plot <- plot$plot + theme(legend.title=element_blank()) # Remove
the title from the legends

  png(filename, width = width, height = height, units = "px", res=300,
pointsize=12, bg = "white")

  print(plot)

  dev.off()
}

generateResidualPlot <- function(data, filter, layout, width, height,
filename) {

  if (!is.null(filter)) {
    data <- subset(data, grepl(filter, vendor))
  }
}

```

```

status <- as.numeric(data[["failure"]])
time <- as.numeric(data[["drive_days"]])
group <- data[["model_pp"]]

fit1 = cph(Surv(time, status) ~ group, x=T,y=T)
f <- cox.zph(fit1)

png(filename, width = width, height=height, units = "px", res=30,
pointsize=8, bg = "white")
par(mfrow=layout)
print(plot(f, resid=F))
dev.off()
}

# All Hard Drives in One Group

generatePlot(filtered_data, "All
HardDrives", "figs/all_harddrives_single_group.png", 2200, 1400, 0.3,
TRUE)
generatePlot(filtered_data, "Grouped",
"figs/all_harddrives_overview.png", 4000, 4000, 0.3, FALSE)

# HGST
generatePlot(filtered_data, "Grouped", "figs/hgst.png", 3000, 3300, 0.3,
FALSE, "HGST")
#generateResidualPlot(filtered_data, "HGST", c(1,1),
2000, 2000, "figs/hgst_res.png")

## Hitachi
generatePlot(filtered_data, "Grouped", "figs/hitachi.png", 3500, 2000,
0.3, FALSE, "Hitachi")
#generateResidualPlot(filtered_data, "Hitachi", c(2,2),
4000, 4000, "figs/hitachi_res.png")

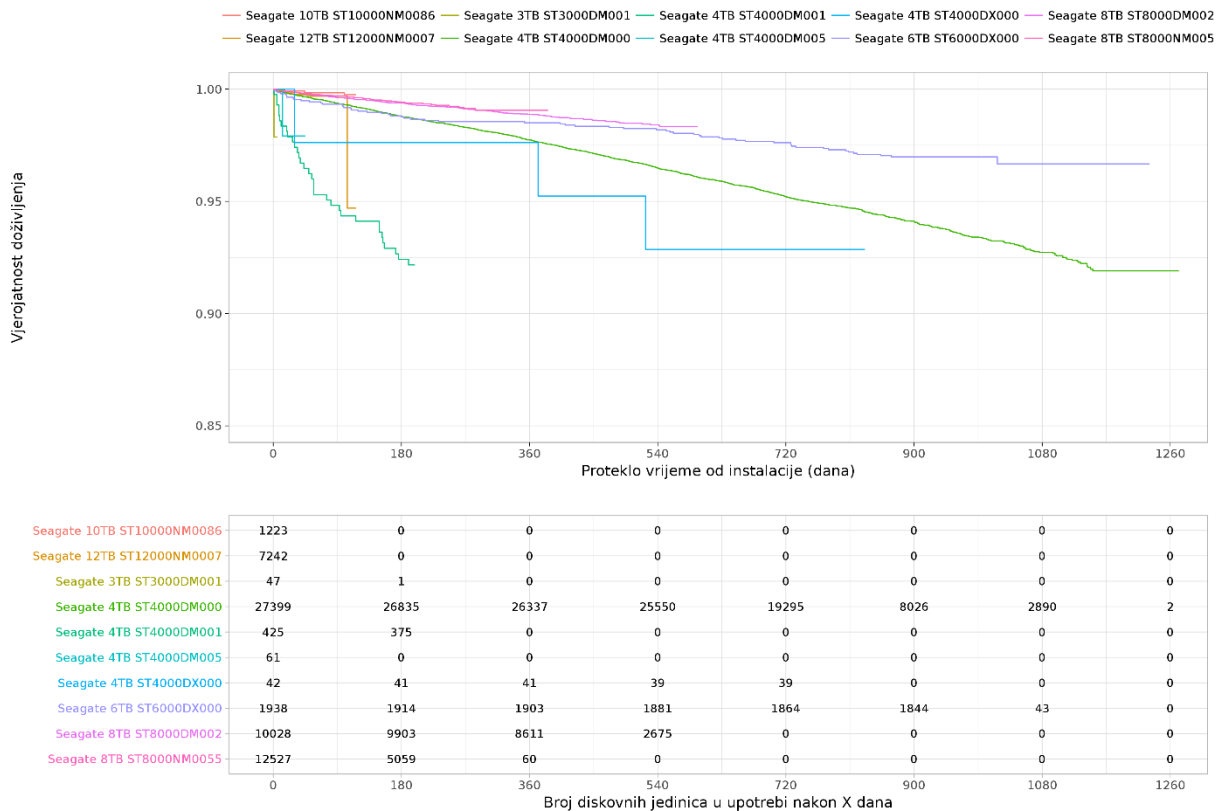
## Seagate
generatePlot(filtered_data, "Grouped", "figs/seagate.png", 3800, 2600,
0.4, FALSE, "Seagate")
generateResidualPlot(filtered_data, "Seagate", c(4,2),
1600, 4000, "figs/seagate_res.png")

```

```
## WDC
generatePlot(filtered_data, "Grouped", "figs/wdc.png", 3500, 2000, 0.3,
FALSE, "Western Digital")
#generateResidualPlot(filtered_data, "WDC", c(1,2),
4000,1600,"figs/wdc_res.png")
```

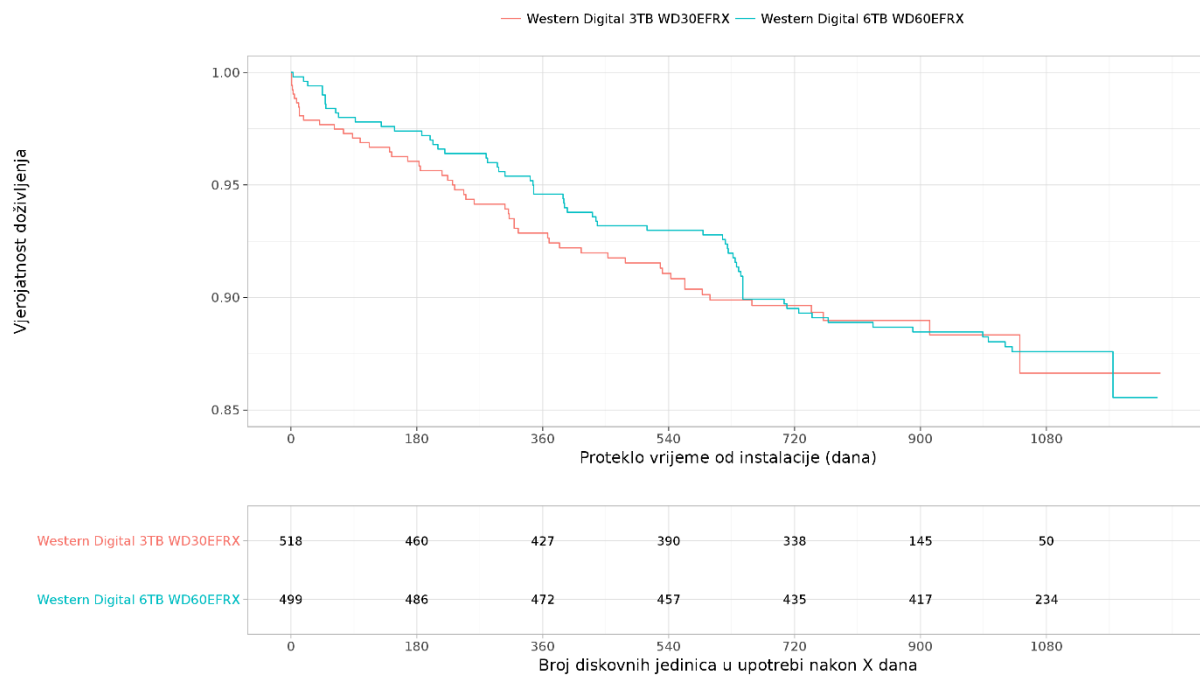
10.4 Rezultati analize podataka u programskom jeziku „R“

Krivulje doživljenja ima smisla prikazati samo za određene skupine podataka. Kada bi željeli prikazati sve modele diskova, takav graf je nepregledan te sama vizualizacije nema smisla. Kako bi mogli prikazati dobivene krivulje diskovne jedinice smo kategorizirali prema proizvođačima i njihovim modelima.

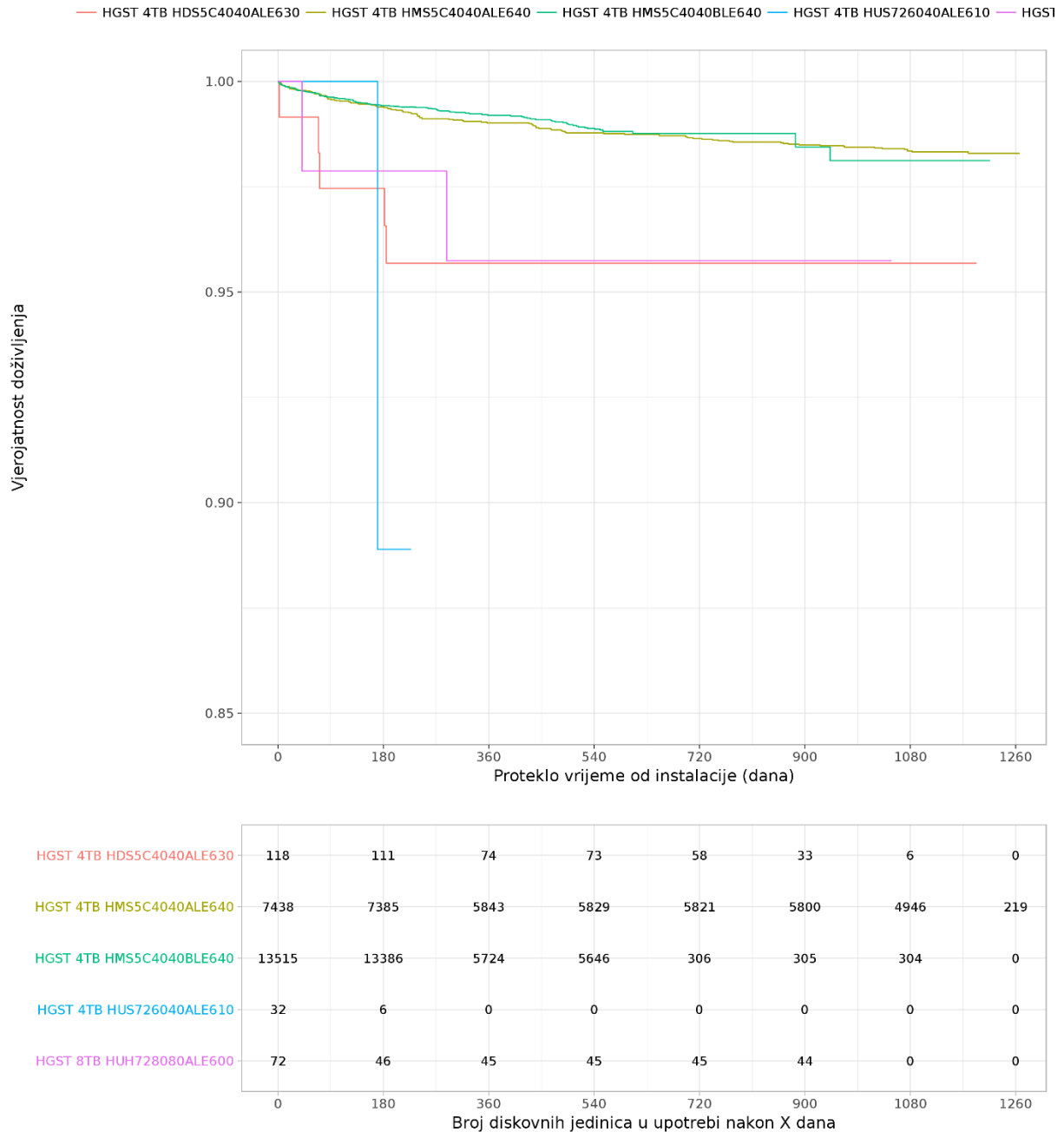


Slika 10: Vjerojatnos doživljenja prema Kaplan-Meierovoj krivulji doživljenja za SEAGATE

Ova vizualizacija podataka jasno prikazuje kako se tvrtka Backblaze odlučila gotovo 27% ukupnog broja diskovnih jedinica u upotrebi temelji na modelu Seagate 4TB ST4000DM000. Velika zastupljenost posljedica je vjerojatno poslovne odluke temeljene na mnogim faktorima (poput cijene, dobavljalivosti, uvjeta i roka isporuke i mnogih drugih) no jedan od njih je relativno visoka pouzdanost. Opravdanost ove odluke na ovom grafu jasno se može prepoznati.



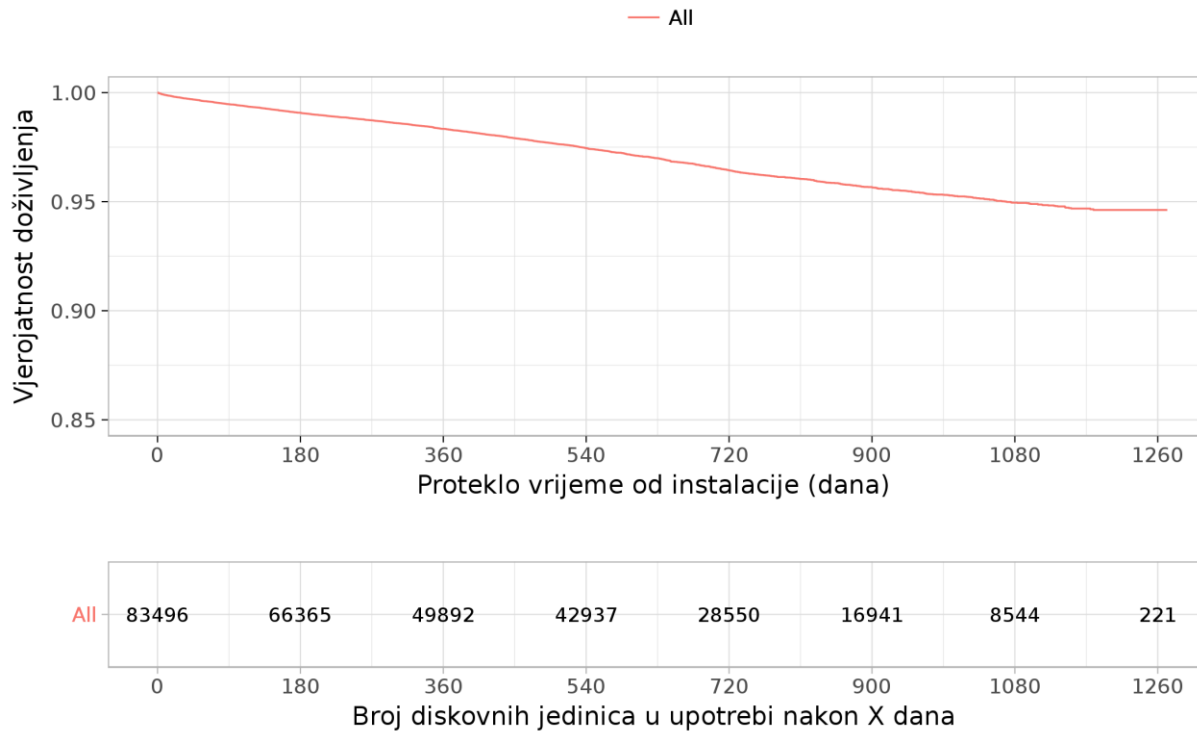
Slika 11: Vjerojatnost doživljenja prema Kaplan-Meierovoj krivulji doživljenja za Western Digital



Slika 12: Vjerojatnost doživljenja prema Kaplan-Meierovoj krivulji doživljenja za HGST

Ova vizualizacija podataka jasno prikazuje kako se oba modela proizvođača HGST – HGST 4TB HMS5C4040ALE640 i HGST 4TB HMS5C4040BLE640 pokazuju izuzetnu pouzdanost.

Potrebno je biti oprezan kako se ne bi donjela kriva odluka. Naime bez dodatne analize drugih uvjeta mogla donjeti pogrešna poslovna odluka da ubuduće sve diskovne jedinice budu od proizvođača tvrtke HGST. Nužno je razmotriti i druge faktore koji utječe na odluku o odabiru a posebno cijenu! Zbog tehnološki vrlo naprednih rješenja (upotreba helija) – njihova cijena je vrlo visoka. Potrebne su zasebne analize i uzimanje u obzir svih dodatnih parametara kako bi se donjela ispravna poslovna odluka.



Slika 13: Vjerojatnost doživljenja - svi modeli svih proizvođača

11 Zaključak istraživanja:

Tvrtka Backblaze donosi odluke na temelju rezultata primjene modela rudarenja podataka. Dodatnim statističkim analizama i dubljim razumijevanjem samog poslovnog procesa, uzimanje u obzir mnogih drugih varijabli koje utječu na poslovanje - kako tehnoloških, vezanih uz samu uslugu koju pružaju tako i utjecaja ostalih elemenata poslovanje moguće je izravno utjecati na poslovanje tvrtke.

Vizualizacija podataka temeljenih na ovom primjeru na jednostavan način omogućava managementu tvrtke donošenje strateških odluka. Ogroman potencijal pružaju otvoreni podaci koje tvrtka Backblaze objavljuje na svojim internet stranicama. Analizom vremenskih serija ponašanja SMART parametara svih diskovnih jedinica koje su ispale, potencijalno može sa određenom statističkom vjerojatnosti omogućiti predviđanje samog kvara. Međuodnos pojedinih SMART parametara i njihov doprinosi statističkoj vjerojatnosti u predviđanju kvara, siguran sam može stvoriti nove modele predviđanja ispada (kvara) pojedinih diskovnih jedinica. Omogućilo bi to izmjenu same diskovne jedinice prije kvara, smanjujući time rizike gubitka podataka na pojedinim poljima, eliminirati dodatne rizike prisutne prilikom obnove diskovnih polja (rebuild) a posredno operativne troškove upravljanja podatkovnim centrom a samim time uspješnijim poslovanjem same tvrtke.

Polja redundantnih diskova unose redundantnost diskovnih jedinica u njihovu upotrebu. Tehnološka rješenja poput RAID10, RAID6, RAID5 ili vrlo naprednih RAID-DP s vrlo visokom vjerojatnošću jamče da neće doći do gubitka podataka ispadom pojedine diskovne jedinice. Napredna tehnološka rješenja redundantnih diskovnih polja uzimaju u obzir da su tvrdi diskovi (diskovne jedinice) - mehanički uređaji čije je životni vijek ograničen (vremenom i načinom upotrebe, okolinom u kojoj se koriste, broj I/O operacija čitanja izapisivanja te mnogim drugim faktorima.

Zbog svega ovog navedenog ispad diskovnih jedinica je očekivan događaj te je on funkcija već spomenutih utjecaja. Usprkos tome, stalna su nastojanja podatkovnih centara da smanje negativne utjecaje ispada diskovnih jedinica. Pouzdanost diskovnih jedinica jedan je od najvažnijih kriterija prilikom donošenja odluka o njihovom odabiru.

I na kraju sva analiza podataka, svi iscrtani grafovi pokazuju da su HGST diskovne jedinice daleko najpouzdanije diskovne jedinice.

Potrebno je biti oprezan kako se ne bi donjela kriva odluka. Naime bez dodatne analize drugih uvjeta mogla donjeti pogrešna poslovna odluka da ubuduće sve diskovne jedinice budu od proizvođača tvrtke HGST. Nužno je razmotriti i druge faktore koji utječe na odluku o odabiru a posebno cijenu! Zbog tehnološki vrlo naprednih rješenja (upotreba helija) – njihova cijena je

vrlo visoka. Potrebne su zasebne analize i uzimanje u obzir svih dodatnih parametara kako bi se donjela ispravna poslovna odluka.

12 Popis literature:

1. Wirth, R., & Hipp, J. (2000, April). CRISP-DM: Towards a standard process model for data mining. In *Proceedings of the 4th international conference on the practical applications of knowledge discovery and data mining* (pp. 29-39). Citeseer.
2. Azevedo, A. I. R. L., & Santos, M. F. (2008). KDD, SEMMA and CRISP-DM: a parallel overview. *IADS-DM.*, (<http://recipp.ipp.pt/bitstream/10400.22/136/3/KDD-CRISP-SEMMA.pdf>, 2018.)
3. Hughes, G. F., Murray, J. F., Kreutz-Delgado, K., & Elkan, C. (2002). Improved disk-drive failure warnings. *IEEE transactions on reliability*, 51(3), 350-357. (<http://www.cs.cmu.edu/~15849g/readings/hughes02.pdf>)
4. Pinheiro, E., Weber, W. D., & Barroso, L. A. (2007). Failure trends in a large disk drive population. (https://www.usenix.org/legacy/events/fast07/tech/full_papers/pinheiro/pinheiro_of_d.pdf)
5. Botezatu, M. M., Giurgiu, I., Bogojeska, J., & Wiesmann, D. (2016, August). Predicting disk replacement towards reliable data centers. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 39-48). ACM. <https://www.kdd.org/kdd2016/papers/files/adf0849-botezatuA.pdf>
6. Aussel, N., Jaulin, S., Gandon, G., Petetin, Y., Fazli, E., & Chabridon, S. (2017, December). Predictive models of hard drive failures based on operational data. In *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)* (pp. 619-625). IEEE. <https://hal.archives-ouvertes.fr/hal-01703140/document>
7. Rosenthal, D. S., Rosenthal, D. C., Miller, E. L., Adams, I. F., Storer, M. W., & Zadok, E. (2012). The economics of long-term digital storage. *Memory of the World in the Digital Age, Vancouver, BC.* <https://www.crss.ucsc.edu/papers/rosenthal-unesco12.pdf>
- Goel, M. K., Khanna, P., & Kishore, J. (2010). Understanding survival analysis: Kaplan-Meier estimate. *International journal of Ayurveda research*, 1(4), 274. (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3059453/>)
8. Knjižnica funkcija programskog jezika R – modul – **survival** - <https://cran.r-project.org/web/packages/survival/index.html>
Knjižnica funkcija programskog jezika R – modul – **survminer** - <https://cran.r-project.org/web/packages/survminer/index.html>

13 Popis slika:

Slika 1: Faze procesa otkrivanja znanja u podacima, (Izvor: dr. Dijana Oreški, FOI, 2017.)	2
Slika 2: Rezultati ankete i usporedba zastupljenosti CRISP-DM metodologije u 2007. i 2014. (N=2000)3	
Slika 3: Udio ispada diskovnih jedinica za sve diskovne jedinice, svih (proizvođača). 2013. - 2015.	12
Slika 4: Grafički prikaz učestalost kvarova modela tvrdih diskova u Q4 2016.	23
Slika 5: Grafički prikaz učestalosti ispada diskovnih jedinica 1.1 2016. - 31.12.2016.	25
Slika 6: Grafički prikaz učestalost ispada diskovnih jedinica 1.1 2015. - 31. 12.2016.	27
Slika 7: Graf kretanja cijene 1 gigabayta prostora 1980.-2015.	29
Slika 8: "R" studio na linux poslužitelju za potrebe stvaranja R upita nad promatranim podacima.....	31
Slika 9: Izvorni kod programa za izradu krivulja u programskom jeziku "R"	32
Slika 10: Vjerojatnos doživljenja prema Kaplan-Meierovoj krivulji doživljenja za SEAGATE.....	37
Slika 11: Vjerojatnost doživljenja prema Kaplan-Meierovoj krivulji doživljenja za Western Digital.....	38
Slika 12: Vjerojatnost doživljenja prema Kaplan-Meierovoj krivulji doživljenja za HGST	39
Slika 13: Vjerojatnost doživljenja - svi modeli svih proizvođača	40

14 Popis tablica:

Tablica 1: Pojašnjenja pojedinih S.M.A.R.T atributa i parametara.....	16
Tablica 2: Učestalost kvarova diskovnih jedinica prema modelu u Q4, 2016.....	23
Tablica 3: Učestalost ispada diskovnih jedinica 1.1 2016. - 31.12.201 24	
Tablica 4: Učestalost ispada diskovnih jedinica 1.1 2015. - 31.12.2016	26
Tablica 5: Tablica modula korištenih u stvaranju upita u programskom jeziku „R“	32