



**LUNDS**  
UNIVERSITET

# **Maskininlärning & Random Forest: Överträffar traditionella kreditmodeller**

Kandidatuppsats

Nationalekonomiska institutionen, Lunds Universitet

Januari, 2023

**Författare:** Simon Johansson

**Handledare:** Andreas Johansson

## **Abstract**

The Altman Z-Score model is one of the most famous models for predicting bankruptcy and measuring financial distress for companies. It uses multivariate discriminant analysis to classify companies in three different groups based on their calculated Z-Score. The purpose of this thesis is to analyze how well the Altman Z-Score model for emerging markets performs, to then attempt to create a new binary classification model using machine learning and random forest-algorithms in order to get a more precise model that better predicts bankruptcy for companies within 2 years. This was done using financial statements from 9520 Polish firms along with their respective bankruptcy status after 2 years.

The results show that regardless of how you interpret the Altman model, it was possible to create a random forest-model that outperformed it by all measurements. The random forest models managed to beat the Altman Z-Score models by predicting a slightly higher percentage of both bankruptcies as well as non-bankruptcies, resulting in an overall higher accuracy. This was done by using the same four financial ratios that are used in Altman's model, i.e. no further information was added.

## Innehållsförteckning

<b>Introduktion.....</b>	<b>5</b>
Bakgrund.....	5
Syfte.....	6
Begränsningar.....	6
<b>Tidigare litteratur.....</b>	<b>6</b>
<b>Teori &amp; metod.....</b>	<b>7</b>
Altman Z-Score.....	7
Altman Z-Score originalmodell.....	7
Altman Z-Score för tillväxtmarknader.....	8
Metod för utvärdering av Altman Z-Score.....	10
Maskininlärning: Beslutsträd & Random Forest.....	10
Beslutsträd.....	10
Random Forest.....	12
Metod för utvärdering genom Random Forest.....	13
<b>Data.....</b>	<b>14</b>
<b>Analys.....</b>	<b>15</b>

Analys av Altman Z-Score.....	15
Analys av Random-Forest-modeller.....	17
Jämförelse av modellerna.....	19
<b>Slutsats.....</b>	<b>20</b>
<b>Referenser.....</b>	<b>22</b>

# Introduktion

## Bakgrund

Kreditrisk kan definieras som risken att kredittagare inte uppfyller sina skyldigheter gentemot kreditgivare och att en förlust därmed uppstår. Att bedöma kreditrisken och förutspå eventuell ekonomiskt obestånd hos såväl privatpersoner som företag är och har alltid varit en central del av finansbranschen. För att banker och andra kreditinstitut ska kunna bedriva sin verksamhet på ett hållbart sätt krävs det att denna bedömning sker på ett så pricksäkert sätt som möjligt. En bra modell är förutsättningen dels för att kunna undvika kreditförluster i så stor utsträckning som möjligt, men även för att kunna erbjuda kunder med högre kreditvärdighet konkurrenskraftiga avgifter och räntor. Detta område kan anses vara särskilt aktuellt i en tid där Sverige är på väg in i en lågkonjunktur, samtidigt som andra yttre faktorer såsom stigande el- och bensinpriser, riskerar att öka antalet företag som går i konkurs. (Konjunkturinstitutet, 2022)

För att bedöma kreditrisken konstrueras diverse modeller grundade på historisk data. I detta fall är det den första multivariabla modellen, framtagen av Edward I. Altman år 1968, som kommer att ligga till grund för arbetet. Modellen, som har för avsikt att förutspå konkurs eller ekonomisk nöd hos företag, är baserad på en linjär kombination av 4 eller 5 kvoter av företagets nyckeltal (Heine, 2000). Modellen varierar lite beroende på vilken typ av företag som är aktuellt, men är oavsett inte särskilt avancerad utan kräver endast ett par enkla räknesteg för att kunna dra slutsatser. Detta i kontrast till de mer avancerade och metoderna som numera finns tillgängliga. Dels finns det markant större mängd information och data tillgängligt idag jämfört med år 1968, vilket underlättar vid skapandet av modeller. Vidare så spelar den teknologiska utvecklingen naturligtvis stor roll, det är en avsevärd skillnad på datorkraft samt tekniker och metoder för dessa typer av uppgifter idag jämfört med när Altman utvecklade sin modell. Bland annat är så kallad maskininlärning ett område som utvecklats kraftigt under denna period (Dataversity, 2021). Därför kommer maskininlärningsmodeller tillsammans till Altmans Z-Score-modell ligga till grund för arbetet. Hur väl kan Altmans Z-Score-modell predicera konkurser på en

tillväxtmarknad och kan man vidare överträffa detta resultat med hjälp av maskininlärning, mera specifikt Random Forest-metoder?

## **Syfte**

Syftet med arbetet är att skapa en binär kreditmodell avsedd för tillväxtmarknader med hjälp av maskininlärning (Random Forest), för att sedan utvärdera om denna kan prestera bättre än Altmans Z-Score-modell från 1968.

## **Begränsningar**

Arbetet är begränsat till att endast hantera och utvärdera Altmans kreditmodell från 1968 samt de modeller som skapas genom Random Forest-metoden inom maskininlärning. Det finns naturligtvis andra kreditmodeller och maskininlärningsmetoder som hade kunnat appliceras, men för att arbetet ska bli konkret läggs fokus endast på ovan nämnda modeller och metoder. Vidare är modellerna begränsade till länder som anses vara tillväxtmarknader. Altmans originalmodell analyseras alltså inte djupare, utan den reviderade modellen avsedd för tillväxtmarknader är den som kommer vara aktuell.

## **Tidigare litteratur**

Altmans Z-Score-modell är en av de mest kända modellerna som kan användas för att predicera ekonomiskt obestånd hos företag. När den först utformades var modellens totala precision 72% enligt det aktuella stickprovet och efter vidare tester fram till och med 1999 visade sig modellen konstant ha en precision mellan 80-90% (Heine, 2000). Detta resultat kan jämföras med andra modeller, bland annat Ohlson O-Score-modellen från 1980 som liknar Altmans modell, men är mer komplex då den tar hänsyn till fler variabler, samtidigt som den är baserad på markant större underlag. Totalt tog Ohlson hänsyn till 2163 bolag, varav 105 av dem var företag som gått i konkurs (Ohlson, 1980). Generellt pekar litteraturen på att Ohlsons modell presterar något bättre överlag. Bland annat visade en studie på som gjordes på små- och medelstora företag i

Indonesien på att Altmans modell lyckades predicera 37,4% av konkurser inom två år, samtidigt som Ohlson lyckades predicera 43,2% (Pramudita, 2020).

När det kommer till maskininlärning tyder litteraturen på att detta är fördelaktigt jämfört med traditionella metoder för att predicera konkurser. Diverse maskininlärningsmodeller visar generellt högre precision än de traditionella, bland annat visade Barboza, Kimura och Altman själv att maskininlärningsmetoder presterar bättre än traditionella metoder med uppskattningsvis 10% (Barboza, Kimura & Altman, 2017). Detta kan styrkas ytterligare när Salim Lahmiri och Stelios Bekiros även de visade på att maskininlärning kan användas för att överträffa de traditionella modellernas resultat (Lahmiri & Bekiros, 2019).

## **Teori och metod**

### **Altman Z-Score**

#### **Altman Z-Score: Originalmodellen**

Altmans Z-Score modell är en kreditmodell som kan användas för att upptäcka ekonomiska problem och förutspå konkurser hos företag inom kommande två år.<sup>5</sup> Det är en linjär modell som i sin originalform är baserad på fem kvoter. (Heine, 2000) De nyckeltal som kvoterna är baserade på definieras enligt följande:

*Rörelsekapital:* Det kapital som används i företagets dagliga verksamhet.

*Totala tillgångar:* Summan av alla företagets tillgångar.

*Balanserad vinst:* Ackumulerad vinst från tidigare verksamhetsår, efter eventuella utdelningar.

*Resultat före räntor och skatt:* Företagets resultat före avdrag för räntor och skatt. Även känt som EBIT (Earnings Before Interest and Tax).

*Marknadsvärde på eget kapital:* Totalt marknadsvärde på företaget enligt den aktuella aktiekursen. Även känt som Market Cap.

*Totala skulder:* Summan av alla företagets skulder.

*Försäljning*: Total försäljning under föregående verksamhetsår.

Dessa nyckeltal kombineras enligt nedan för att framställa de fem kvoterna:

$$X_1 = \frac{\text{Rörelsekapital}}{\text{Totala tillgångar}}$$

$$X_2 = \frac{\text{Balanserad vinst}}{\text{Totala tillgångar}}$$

$$X_3 = \frac{\text{Resultat före räntor och skatt}}{\text{Totala tillgångar}}$$

$$X_4 = \frac{\text{Marknadsvärde på eget kapital}}{\text{Totala skulder}}$$

$$X_5 = \frac{\text{Försäljning}}{\text{Totala tillgångar}}$$

Dessa kvoter viktas sedan med hjälp av en rad konstanter enligt följande:

$$Z = 1,2X_1 + 1,4X_2 + 3,3X_3 + 0,6X_4 + 1,0X_5$$

Ju högre värde på Z desto bättre för företaget, det vill säga lägre risk för ekonomiskt obestånd.

Det är värt att notera att samtliga kvoter har ett positivt samband med Z, vilket innebär ett negativt samband med risken för ekonomiskt obestånd. Z utvärderas sedan enligt följande olikheter:

$Z > 2,99$  – Företagets risk för ekonomiskt obestånd är låg.

$2,99 < Z < 1,80$  – Gråzon. Företagets risk för ekonomiskt obestånd är måttlig.

$Z < 1,8$  – Företagets risk för ekonomiskt obestånd är hög.

### **Altmans Z-Score för tillväxtmarknader**

Den ursprungliga versionen av Altmans Z-Score är baserad på företag som anses vara tillverkare, det vill säga företag som producerar och monterar produkter från råvaror och/eller inköpta delar.



Vidare är den även baserad på publikt handlade bolag, då den genom  $X_4$  tar hänsyn till det aktuella marknadspriset på utestående aktier. För att modellen skulle bli applicerbar fler företag reviderades modellen och nya versioner skapades, bland annat för icke-tillverkare samt för företag i tillväxtmarknader. (Swalih, Adarsh, & Sulphey, 2021)

En tillväxtmarknad anses vara en marknad som präglas av stor ekonomisk tillväxt samt delvis, men inte fullt ut, besitter de ekonomiska egenskaper som normalt förknippas med utvecklade länder. (Corporate Finance Institute, 2022) Exakt vilka länder och marknader som definieras inom denna ram är godtyckligt.

I Altmans modell för tillväxtmarknader justeras antalet kvoter från fem till fyra, och ett nytt begrepp introduceras:

*Bokfört värde på eget kapital:* Differensen mellan totala tillgångar och totala skulder enligt företagets balansräkning.

De nya kvoterna beräknas enligt följande:

$$X_1 = \frac{\text{Rörelsekapital}}{\text{Totala tillgångar}}$$

$$X_2 = \frac{\text{Balanserad vinst}}{\text{Totala tillgångar}}$$

$$X_3 = \frac{\text{Resultat före räntor och skatt}}{\text{Totala tillgångar}}$$

$$X_4 = \frac{\text{Bokfört värde på eget kapital}}{\text{Bokfört värde på totala skulder}}$$

Sedan viktas de enligt nedan:

$$Z = 3,25 + 6,56X_1 + 3,26X_2 + 6,72X_3 + 1,04X_4$$

Liksom tidigare innebär ett högre Z-värde bättre ekonomiskt tillstånd och därmed lägre risk för ekonomiskt obestånd. Dock utvärderas detta Z-värde något annorlunda jämfört med i den ursprungliga modellen, de olika spannen är reviderade för att bättre passa den nya modellen:

$Z > 2,6$  - Företagets risk för ekonomiskt obestånd är låg.

$2,6 < Z < 1,1$  - Gråzon. Företagets risk för ekonomiskt obestånd är måttlig.

$Z < 1,1$  - Företagets risk för ekonomiskt obestånd är hög.

### **Metod för utvärdering av Altman Z-Score**

När datan skulle analyseras med hjälp av Altman Z-Score användes Microsoft Excel för att utifrån de fyra aktuella nyckeltalen göra de nödvändiga beräkningarna för att utvinna varje företags Z-värde.

## **Maskininlärning: Beslutsträd och Random Forest**

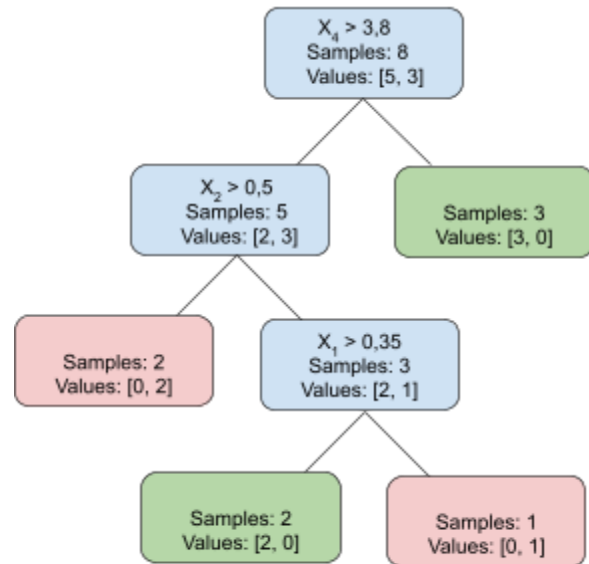
### **Resultsträd**

Beslutsträd (*eng. Decision Tree*) är en form av övervakad maskininlärning som kan användas för att förutspå konkurser hos företag. Beslutsträd där målvariabeln endast kan anta på förhand bestämda diskreta värden kallas för ett klassificeringsträd. Om målvariabler däremot kan anta värden inom ett kontinuerligt spann, exempelvis mellan 0 och 1, kallas det istället för ett regressionsträd. I detta fall är det binära klassificeringsträd som är aktuella, då resultaten ska jämföras med Altmans Z-score modell, vilken man kan se som en binär modell.

Ett klassificeringsträd skapas genom att en datamängd delas upp upprepade gånger och därmed skapar en struktur som kan liknas vid ett träd. Punkterna i trädet som potentiellt delar upp observationerna kallas för noder och består dels utav beslutsnoder som delar upp observationerna i ytterligare två nya noder, samt bladnoder där observationerna ej delas upp vidare. Som utgångspunkt delas observationerna upp tills endast observationer av samma klass finns kvar i en

nod, men detta kan justeras med hjälp av hyperparametrar som bestäms vid skapandet av modellen, vilket gör att det kan finnas ytterligare villkor som ska uppfyllas för att en nod ska dela upp observationerna. (Du & Zhang, 2002)

Se Figur 1 som visar ett enkelt beslutsträd, där de blåa noderna är beslutsnoder och de gröna respektive röda är bladnoder som motsvarar de två olika klasserna.



Figur 1

Hur bestäms villkoret för hur observationerna delas upp i en beslutsnod? Avsikten med att uppdelningen är att göra datan i de två nya

noderna “renare” och för att uppnå detta kan man använda sig av olika metoder. Två vanliga metoder för detta är att använda sig av entropi- respektive gini, och i detta arbete används gini för att skapa modellerna. Giniindex är ett mått på “orenheten” hos en datamängd, det vill säga hur homogen den är. Indexet ligger för en given datamängd med 2 klasser per definition mellan 0 och 0,5, där ett lägre värde innebär lägre “orenhet”, det vill säga mer homogen data, och vice versa. Indexet definieras enligt följande:

$$GI = 1 - \sum_{i=1}^m p_i^2$$

där  $m$  är definierat som antal klasser och  $p_i$  är definierat som andel av klass  $i$  i den totala datamängden. I varje beslutsnod tillsätts en olikhet för någon av datans förklarande variabler som minimerar summan av de viktade giniindexen hos de två nya noderna. Med andra ord så maximeras “renheten” i dessa nya noder och därmed optimeras uppdelningen.<sup>11</sup>

Efter att trädet har skapats enligt ovan kan man med hjälp av detta klassificera nya observationer genom att låta dem “vandra” genom trädet enligt olikheterna som tillsatts i de olika beslutsnoderna, tills att de slutligen hamnar i en bladnod där de tillsätts en predicerad klass.

## Random forest

Random Forest är en utökad version av ett beslutsträd där man kombinerar ett större antal beslutsträd i samma modell. Dessa beslutsträd skapas med hjälp av mindre delmängder av datan och antalet går att justera utifrån vad som önskas. När en ny observation sedan ska klassificeras, får den vandra genom samtliga träd för att sedan klassificeras utifrån resultatet från dessa. Av olika skäl anses Random Forest-modeller i många fall vara fördelaktiga jämfört beslutsträd. Bland annat har de mindre problem med så kallad *overfitting* (Ying, 2019)

Overfitting innebär att modellen lär sig alltför specifikt av datan och därmed misslyckas generalisera till en rimlig grad. Exempel på detta vore om vi ställer upp ett hypotetiskt scenario där vi med hjälp av bolags föregående års resultat samt eget kapital ska förutspå huruvida bolaget går i konkurs under nästa verksamhetsår. De flesta bolag som gått i konkurs tidigare enligt datan vi tränar vår modell med har haft kraftigt negativa resultat samt lågt eget kapital, och de som inte gått i konkurs har haft bättre resultat samt högre eget kapital. Dock finns det potentiellt bolag i vår data som trots bra resultat föregående verksamhetsår samt högre eget kapital gått i konkurs. I en modell som lider av *overfitting* kan detta komma att skapa problem om vi ska analysera nya bolag vars siffror liknar de ovan beskrivna alltför mycket. Vi önskar naturligtvis att modellen generaliserar och därmed inte förutspår konkurs hos dessa nya bolag med positiva siffror bara för att de råkar likna bolag som avvek från mängden, men desto större problem modellen har med *overfitting* desto större blir risken för att modellen predicerar något som inte är optimalt.

Vidare är precisionen hos Random Forest-modell ofta högre jämfört med ett beslutsträd och en Random Forest-modell kan även hantera avsaknad av värden i datan, vilket ett beslutsträd ej kan. Detta kommer dock på bekostnad av högre krav på datorkraft och högre komplexitet, vilket gör att modellen inte blir lika intuitiv och lätt att visualisera.

## Metod för utvärdering genom Random Forest

Vid analys av datan genom en Random Forest-modell användes programmeringsspråket Python i förening med ett befintligt Python-bibliotek avsett för maskininlärning och dataanalys vid namn *scikit-learn*. Scikit-learn “tränar” modellen med hjälp av data som anges ihop med diverse hyperparametrar som kan justeras för att optimera modellen för sitt syfte.

Datan som anges vid skapandet av modellen är dock inte hela datamängden som vi har tillgång till, utan en viss andel behöver sparas för att sedan kunna användas för att utvärdera hur väl modellen presterar. Desto större andel av datan man använder för att träna modellen, desto bättre “tränad” blir modellen, men detta på bekostnad av en högre varians när man utvärderar modellens prestation. Det är vanligt att man använder sig av 60-80% för att träna modellen och att man därmed sparar 20-40% till utvärderingen. (V7, 2023) Detta arbete kommer därmed använda sig av 70% av datan för träning respektive 30% för utvärdering. Det är även värt att notera datan delas in helt slumpvis i dessa delmängder. I detta arbete användes den befintliga funktionen i scikit-learn *train\_test\_split* för att utföra uppdelningen.

Som tidigare nämnt skapas modellen med hjälp av en rad hyperparametrar som går att justera för att optimera modellen. Efter att ha justerat de olika parametrarna fram och tillbaks drogs slutsatsen att samtliga hyperparametrar var optimerade genom sina standardvärden, det vill säga de värden som scikit-learn tillsätter om inget annat anges, bortsett från följande två som behövde ändras för att optimera resultatet:

***class\_weight*** (*None* → “*balanced*”): Den kanske viktigaste justeringen som gjordes för att modellen skulle fungera på ett rimligt sätt var att vikta klasserna olika. Eftersom det finns en kraftigt skev fördelning mellan de två klasserna i datan (5,41% konkurser respektive 94,59% icke-konkurser), vilket innebar att modellen predicerade icke-konkurs alltför ofta då de två klasserna var viktade lika. Istället justerades hyperparametern *class\_weight* och sattes istället till “*balanced*”, vilket innebar att klasserna viktades omvänt proportionerligt gentemot frekvensen av respektive klasser i datan. Konkurser viktades därmed mycket högre än icke-konkurser, vilket resulterade i en modell som presterade markant bättre.

*max\_depth* (*None* → 8 eller 10): Som utgångspunkt delar träden upp datan tills antingen tills endast observationer av samma klass finns kvar i noden, eller tills något villkor sätter stopp på uppdelningen. Denna hyperparameter är ett sådant villkor och anger trädets maximala djup, då trädet ej delas upp vidare efter att detta djup har uppnåtts. Den har standardvärde *None*, med andra ord oändligt maximalt djup, vilket innebär att uppdelningen aldrig upphör på grund av trädets djup. Vid justering av denna parameter visade det sig att desto lägre värde den sätts till, desto lägre blir antal *fel av typ 1* samtidigt som *fel av typ 2* ökade. Vid ett värde på 8 blev resultatet jämförbart med Altmans modell där gråzonen definierades som konkurs och vid ett värde på 10 blev modellen istället jämförbar med Altmans modell där gråzonen definierades som icke-konkurs.

## Data

Datan som kommer att användas för att testa modellerna är en befintlig datamängd som kommer från hemsidan Kaggle (<https://www.kaggle.com/datasets/bhadaneeraj/bankruptcy-detection>). Enligt specifikationen är datan ursprungligen hämtad från Emerging Markets Information Services (EMIS), en organisation som är verksam i och rapporterar om diverse tillväxtländer runt om i världen. Den avser information om polska företag vilket är intressant eftersom Polen definieras som en tillväxtmarknad vilket innebär att Altmans Z-score modell för tillväxtmarknader är applicerbar.

Datan är uppdelad i fem datamängder som avser olika tidsperioder och för detta arbetet är det filen som heter "4year.csv" som är aktuell. Detta eftersom denna fil innehåller information om huruvida företag gick i konkurs inom de nästkommande två åren, vilket är samma tidshorisont som Altmans modell också avser. Utöver informationen huruvida ett företag gick i konkurs eller ej finns även 64 stycken olika nyckeltal, inklusive de fyra som Altman använder i sin modell, tillgängliga för respektive företag. Totalt innehåller datamängden 9539 observationer, men 19 av dessa saknade värden på ett eller flera av nyckeltalen. 19 stycken utav 9539 är en minimal andel, så detta hanterades genom att dessa observationer helt enkelt togs bort, vilket resulterade i att endast 9520 observationer användes i analysen.

## Analys

### Analys med Altmans Z-Score

Altmans modell innefattar som tidigare beskrivits tre olika intervaller för Z-värdet som beskriver olika ekonomiska tillstånd för ett företag. Detta behöver revideras för att kunna jämföras och analyseras med hjälp av datan som används, då målvariabeln hos datan är binär. Z-värden under 1,1 kommer naturligtvis beskrivas som att modellen predicerar en konkurs inom de kommande två åren, eftersom detta intervall enligt modellen innebär att företagets risk för ekonomiskt obestånd är hög. På samma sätt kommer Z-värden över 2,6 beskrivas som att modellen ej predicerar konkurs för företaget inom de kommande två åren, då detta intervall enligt modellen innebär att företagets risk för ekonomiskt obestånd är hög. Gråzonen där emellan,  $1,1 < Z < 2,6$ , som enligt modellen innebär att företagets risk för ekonomiskt obestånd är måttlig, kan däremot tolkas på olika sätt. Därför kommer två separata analyser göras, en där denna gråzon definieras som predicerad konkurs, och en där den definieras som att modellen ej predicerar konkurs.

Begreppet *fel av typ 1* som används nedan definieras som en inkorrekt prediktion genom att modellen förutspått att ett företag ska gå i konkurs, när det i själva verket inte gjorde det. *Fel av typ 2* definieras istället som en inkorrekt prediktion där modellen förutspått icke-konkurs, men företaget visade sig ändå gå i konkurs inom 2 år. Summan av antalet fel av typ 1 och 2 motsvarar alltså det totala antalet inkorrekta prediktioner. Vidare definieras begreppet *precision* som också används nedan som andel observationer där Altmans modell korrekt predicerade huruvida ett företag skulle gå i konkurs eller ej.

Övergripande:

Antal observationer	9520
Antal faktiska konkurser	515
Antal faktiska icke-konkurser	9005

Tabell 1.1 - Gråzon motsvarar icke-konkurs

Antal predicerade konkurser (Z-score < 1,1)	809
Antal predicerade icke-konkurser (Z-Score > 1,1)	8711
Antal korrekt predicerade konkurser	129
Antal korrekt predicerade icke-konkurser	8325
Fel av typ 1	680
Fel av typ 2	386
<b>Precision (%)</b>	<b>88,80</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>25,05</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>92,45</b>

Tabell 1.2 - Gråzon motsvarar konkurs

Antal predicerade konkurser (Z-score < 2,6)	1290
Antal predicerade icke-konkurser (Z-Score > 2,6)	8230
Antal korrekt predicerade konkurser	180
Antal korrekt predicerade icke-konkurser	7895
Fel av typ 1	1110
Fel av typ 2	335
<b>Precision (%)</b>	<b>84,82</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>34,95</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>87,67</b>

Tabellerna 1.1 & 1.2 ovan beskriver resultatet av analysen som gjordes med hjälp av två olika tolkningar av Altmans Z-Score-modell

Då man definierar gråzonen som icke-konkurs lyckas Altmans modell korrekt predicera 25,05% av konkurser respektive 92,45% av icke-konkurser, vilket motsvarar en total *precision* på 88,80%. Då man istället definierar gråzonen som att modellen predicerar konkurs minskar antal *fel av typ 2* och därmed ökar andelen faktiska konkurser som modellen lyckas predicera korrekt. Dock ökar antalet *fel av typ 1* istället markant, vilket resulterar i en lägre andel faktiska icke-konkurser korrekt predicerade därmed en överlag sämre *precision*.

Vilken av dessa två som är mest fördelaktig kan inte bestämmas, utan beror användarens preferenser och i vilket sammanhang den ska användas. Titta exempelvis på ett hypotetiskt scenario där olika kreditgivare ger ut lån till företag och endast företag som kreditgivaren inte



tror kommer gå i konkurs får sina lån beviljade. Anta vidare att kreditgivarna använder sig av modellen ovan för att förutspå detta samt att de alla har olika riskaptit. Under dessa förutsättningar hade kreditgivarna med lägre riskaptit föredragit den version av modellen som har lägre antal *fel av typ 2*, eftersom denna version innebär lägre kreditförluster då man ej ger ut lån till företag i gråzonen. I kontrast till detta hade kreditgivare med högre riskaptit föredragit modellen där gråzonen motsvarar icke-konkurs, eftersom det genom denna version beviljas fler lån och därmed medföljer potentiellt högre intäkter, dock samtidigt som kreditrisken blir högre.

Teoretiskt sett skulle kreditgivare kunna bortse helt från de intervall som Altman satt i sin modell, och anpassa den själv efter egen önskad riskaptit. Desto högre man sätter gränsen desto lägre blir kreditförlusterna, samtidigt som de potentiella intäkterna minskar och vice versa.

### **Analys med Random Forest**

De två resultaten nedan avser två olika Random Forest-modeller som syftar på att jämföras med de två Altman-modellerna som tidigare beskrivit. Hyperparametrarna har justerats för att få fram två modeller som lätt kan jämföras med de två Altman-modellerna. Det endas som skiljer dem åt är hyperparametern *max\_depth* som är satt till 10 respektive 8. Eftersom vi använt 70% av datan för att "träna" modellen, återstår endast 30% att utvärdera modellen genom. Därför är antalet observationer i nedan analys lägre än det var i när Altmans modell utvärderades. Vidare skiljer sig antalet faktiska konkurser respektive icke-konkurser sig åt mellan de två modellerna nedan. Detta på grund av att urvalet till träning- och testdatan som tidigare nämnts sker slumpmässigt, vilket medför en viss varians när det kommer klassfördelningen i datan som används för att utvärdera modellerna.

Övergripande:

Antal observationer	2856
---------------------	------

Tabell 2.1 – *max\_depth* = 10

Antal faktiska konkurser	157
Antal faktiska icke-konkurser	2699
Antal predicerade konkurser	193
Antal predicerade icke-konkurser	2663
Antal korrekt predicerade konkurser	41
Antal korrekt predicerade icke-konkurser	2547
Fel av typ 1	152
Fel av typ 2	116
<b>Precision (%)</b>	<b>90,62</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>26,11</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>94,37</b>

Tabell 2.2 – *max\_depth* = 8

Antal faktiska konkurser	155
Antal faktiska icke-konkurser	2701
Antal predicerade konkurser	334
Antal predicerade icke-konkurser	2522
Antal korrekt predicerade konkurser	49
Antal korrekt predicerade icke-konkurser	2416
Fel av typ 1	285
Fel av typ 2	106
<b>Precision (%)</b>	<b>85,96</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>39,19</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>88,52</b>

Tabellerna 2.1 & 2.2 ovan beskriver resultatet av analysen som gjordes med hjälp av två olika *Random Forest*-modeller

Modellen där trädens djup har en maxgräns på 8 lyckas predicera en större andel av konkurserna, men detta på bekostnad av en överlag lägre precision, då en lägre andel av icke-konkurserna lyckas prediceras korrekt. Det motsatta gäller då naturligtvis modellen där trädens djup har en maxgräns på 10, en högre precision uppnås, men detta på bekostnad av en lägre andel av faktiska konkurser som korrekt prediceras. Vilken av ovan som är med fördelaktig liksom med Altmans

modeller inte heller sägas på rak arm, utan beror på användarens preferenser och användningsområde.

### Jämförelse av modellerna

Tabellerna nedan visar en sammanfattning av de två versionerna av Altman Z-Score samt de två versionerna av Random-Forest som testats. Respektive Altman-modell visas i en vars en tabell tillsammans med den Random Forest-modell som skapats för att jämföras med denna.

*Tabell 3 – Jämförelse mellan motsvarande Altman- respektive Random Forest-modell*

	Altman Z-Score (Gråzon = icke-konkurs)	Random Forest (max_depth = 10)
<b>Precision (%)</b>	<b>88,80</b>	<b>90,62</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>25,05</b>	<b>26,11</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>92,45</b>	<b>94,37</b>

*Tabell 4 – Jämförelse mellan motsvarande Altman- respektive Random Forest-modell*

	Altman Z-Score (Gråzon = konkurs)	Random Forest (max_depth = 8)
<b>Precision (%)</b>	<b>84,82</b>	<b>85,96</b>
<b>Andel faktiska konkurser korrekt predicerade (%)</b>	<b>34,95</b>	<b>39,19</b>
<b>Andel faktiska icke-konkurser korrekt predicerade (%)</b>	<b>87,67</b>	<b>88,52</b>

Båda modellerna som skapats med hjälp av maskininlärning lyckas prestera bättre än Altmans modeller på samtliga punkter. Detta kan dock anses vara väntat, eftersom Altmans modell är

förhållandevis gammal (från 1968). Det är ett rimligt antagande att konkursrisk baserat på företags nyckeltal ej är konstant över tid, utan ändras något med tiden, vilket medför att en modell från 1968 potentiellt inte fungerar optimalt på data från 2000-talet. Vidare är en Random Forest-modell markant mer komplicerad och krävande än Altmans linjära modeller, och därmed är det även rimligt att vänta sig en bättre prestation.

## **Slutsats**

Efter att ha utvärderat de olika modellerna med hjälp av den tillgängliga datamängden går det att dra slutsatsen att de skapade Random Forest-modellerna som väntat presterar bättre än Altmans modeller genomgående. Detta stöts även av den befintliga litteratur som existerar, som pekar på att modeller baserade på maskininlärning bör kunna överträffa de traditionella. Dock finns det en del saker att anmärka på:

Det finns en viss varians när man skapar och utvärderar Random Forest-modellerna med hjälp av tränings- respektive testdata. Detta eftersom datan delas upp i dessa delmängder slumpvis, vilket gör att man med största sannolikhet skapar och utvärderar modellen baserat på olika data varje gång, även om den totala datan är den samma. Vidare innehåller själva skapandet av modellen ett slumpmässigt moment (därav namnet), vilket också bidrar till en viss varians.

Vidare kan det anses vara problematiskt att använda sig av binära modeller för att förutspå konkurs hos företag. Eftersom det i detta fall skulle jämföras med Altmans modeller var det nödvändigt att skapa binära modeller via maskininlärningen, men annars förutspås konkurser bättre med en sannolikhet istället för en binär målvariabel. Eftersom risk för konkurs är på pass komplext och beror på många faktorer som är mycket svåra eller till och med omöjliga att ta hänsyn till vid skapandet av en kreditmodell, är det inte rimligt att anta att man kan predicera dem binärt med särskilt mycket större framgång än vad som visats i ovan analyser. Istället är det rimligare att skapa en modell som ger en uppskattad sannolikhet för konkurs istället.

Det är även värt att notera att Random Forest-modellernas hyperparametrar potentiellt inte är fullt optimerade. Eftersom det inte finns svart på vitt vilka hyperparametrar är optimala, utan det beror på varje enskilt fall, finns risken att en kombination existerar som hade gett modellerna ännu bättre prestanda än de som tillsattes i detta arbete. Det hade även potentiellt varit möjligt att förbättra modellen genom att utöka de variabler man använder sig av utöver endast de ingår i Altmans modell, då det är fullt möjligt att andra nyckeltal och kvoter existerar som bättre förklarar konkursrisken än just de som Altman tog fram 1968.

## Referenser

- 1) Altman, Edward I., Małgorzata Iwanicz-Drozdowska, Erkki K. Laitinen, Arto Suvas. 2016. *Financial Distress Prediction in an International Context: A Review and Empirical Analysis of Altman's Z-Score Model*. Journal of International Financial Management & Accounting, Volume 28, Issue 2.  
<https://onlinelibrary.wiley.com/doi/full/10.1111/jifm.12053> (Hämtad 2023-01-05)
- 2) Barboza, Flavio, Herbert Kimura, Edward I. Altman. 2017. *Machine learning models and bankruptcy prediction*. Expert Systems with Applications, Volume 83.
- 3) Corporate Finance Institute. 2022. *Emerging Markets*.  
<https://corporatefinanceinstitute.com/resources/economics/emerging-markets/> (Hämtad 2023-01-05)
- 4) Dataversity. 2021. *A Brief History of Machine Learning*.  
<https://www.dataversity.net/a-brief-history-of-machine-learning/> (Hämtad 2023-01-05)
- 5) Du, Wenliang & Zhan, Zhijun. 2002. *Building Decision Tree Classifier on Private Data*. Electrical Engineering and Computer Science. 8.  
<https://surface.syr.edu/eecs/8> (Hämtad 2023-01-05)
- 6) Heine, Max L. 2000. *Predicting Financial Distress of Companies: Revisiting the Z-Score and Zeta® Models*.  
<https://pages.stern.nyu.edu/~ealtman/Zscores.pdf> (Hämtad 2023-01-05)
- 7) Konjunkturinstitutet. 2022. *Lågkonjunkturen står för dörren*.  
<https://www.konj.se/publikationer/konjunkturlaget/konjunkturlaget/2022-09-28-lagkonjunktturen-star-for-dorren.html> (Hämtad 2023-01-05)
- 8) Lahmiri, Salim, Stelios Bekiros. 2019. *Can machine learning approaches predict corporate bankruptcy? Evidence from a qualitative experimental design*. Quantitative Finance.
- 9) Ohlson, James A. 1980. *Financial Ratios and the Probabilistic Prediction of Bankruptcy*. Journal of Accounting Research, Volume 18.
- 10) Pramudita, Adhy. 2020. *The Application of Altman Revised Z-Score Four Variables and Ohlson O-Score as A Bankruptcy Prediction Tool in Small and Medium Enterprise*

*Segments in Indonesia*. Advances in Economics, Business and Management Research, Volume 187.

- 11) Swalih, M., Adarsh, K & Sulphey, M. (2021). *A study on the financial soundness of Indian automobile industries using Altman Z-Score*. Accounting Volume 7, Issue 2.  
<http://m.growingscience.com/beta/ac/4469-a-study-on-the-financial-soundness-of-indian-automobile-industries-using-altman-z-score.html> (Hämtad 2023-01-05)
- 12) V7. 2023. *Train Test Validation Split: How To and Best Practices*.  
<https://www.v7labs.com/blog/train-validation-test-set> (Hämtad 2021-01-05)
- 13) Vitalflux. 2022. *Differences: Decision Tree & Random Forest*.  
<https://vitalflux.com/differences-between-decision-tree-random-forest/> (Hämtad 2023-01-05)
- 14) Xue Ying. 2019. *An Overview of Overfitting and its Solutions*. Journal of Physics: Conference Series 1168 022022.  
<https://iopscience.iop.org/article/10.1088/1742-6596/1168/2/022022/meta> (Hämtad 2023-01-05)