

Chapter 1

Part I

I

1.1 Beskrivelse af data

Alle boligejere i Danmark betaler en skat, ejendomsværdiskat, som er baseret på værdien af deres ejendom. Dette vil sige hele ejendommen inkl. grunden som boligen ligger på. For at kunne gøre dette laver den danske stat offentlige ejendomsvurderinger som disse skatter bliver baseret på. Det er derfor vigtigt at disse vurderinger er retvisende og ikke mindst forklarbare, således at en borger kan forstå hvilke parametre der ligger til grund for ejendomsvurderingen. Til dette projekt har jeg valgt at arbejde med anonymiseret data fra mit arbejde i udviklings- og forenklingsstyrelsen, hvor jeg til dagligt arbejder med netop dette. Datasættet består af ejendomssalg fra en 6-årig periode. Ud over selve huspriserne består data også af en lang række attributter som beskriver karakteristika ved selve boligen. Det kan f.eks. være tagmateriale, boligens opførelsesår, information om størrelsen af huset og grunden eller bbr koder som dækker over boligens anvendelse. Derudover består data også af en lang række attributter som fortæller noget om hvor boligens beliggenhed. Det kan f.eks. være boligens koordinater eller information om afstanden til kyst og skov eller afstand til motorvej og jernbane. Data kommer fra en række forskellige registre og offentlige styrelser som eks. BBR og Styrelsen for Dataforsyning og Effektivisering.

Til dette projekt vil jeg overordnet set prøve at se hvor godt man kan forudsige ejendomsværdier ud fra salgspriserne fra en 6-årig periode.

Jeg vil med Principal Component Analysis få et overblik over de data der er til rådighed og få et visuelt overblik over attributterne. Herefter vil jeg med en unsupervised learning forsøge at gruppere det data jeg har til at generer yderlige attributer som kan indgå i modellen. Jeg vil her specifikt prøve at se om det er muligt at gruppere salgene i forskellige boligtyper. Jeg vil i samme omgang også forsøge at frasortere outliers i data med anomaly

detection. Herefter vil jeg med regressions model fors ge at kaste lys over projektets overordnede problem ved at fors ge at forudsige huspriserne ud fra salgspriser. I tilf lde af at modellen ikke kan komme med en god pr diktion af en given ejendom vil det v re muligt at denne ejendom bliver manuelt v rdiansat af en sagsbehandler. Jeg vil derfor til slut med en klassifikation fors ge at estimerer om en ejendom skal ud til manuel sagsbehandling baseret p  dens estimerede ejendomsv rdi.

1.2 Detaljeret beskrivelse af data

Det salgsdata som jeg har valgt at arbejde med d kker i udgangspunktet 'r nrow(train)' observationer med 'r ncol(train)' attributter. Inden jeg g r i gang med at kigge p  data har jeg valgt at lave en opr dning i data. Det har jeg gjort fordi mange af attributterne bliver i mit daglige arbejde brugt i forbindelse med im dekommene diverse forretningskrav. Desuden d kker observationerne mange forskellige typer af ejendomssalg. Det er en blanding af parcelhussalg, r kkehussalg, sommerhussalg, salg af ejerlejligheder mm. og ud fra et forretningsm ssigt perspektiv giver det ikke mening at tr ne en model p  alle salg og ejendomsstypen vil p virke salgsprisen. F.eks. vil der p  sommerhuse v re restriktioner p  hvor meget om  ret man m  bo i sommerhuset og der kan v re i sommerhusomr der v re andre regler for hvad man m  bruge sin grund til end der er i et parcelhusomr de. Jeg har derfor ligeledes valgt at reducerer antallet af observationer s ledes at de kun d kker almindelige parcelhus. Dette er gjort ved kun at beholde alle de ejendomssalg, hvor ejendommen i BBR er registreret med enheds- og bygningsanvendelsen 120.