

# BAN440 - Term Paper Code

Candidate numbers: xx, xx, xx

## Table of contents

<b>Packages used</b>	<b>1</b>
<b>Data</b>	<b>2</b>
Vinmonopolet API . . . . .	2
Dimpostnummer merge . . . . .	4
Vinmonopolet 2024 . . . . .	4
Kommuneendringer 2017 . . . . .	6
Kommuneendringer 2018 . . . . .	7
Kommuneendringer 2020 . . . . .	8
Kommuneendringer 2024 . . . . .	9
Kommune 2025 . . . . .	10
Demography data . . . . .	11
Distance data . . . . .	12
Model variables merge . . . . .	17
<b>Model applications</b>	<b>25</b>
Data preparation . . . . .	25
Model selection and basic regressions . . . . .	25
Demand estimation . . . . .	32
Logit model . . . . .	34

## Packages used

```
# relevant libraries
library(tidyverse)
```

```

library(readxl)
library(fastDummies)
library(knitr)
library(stargazer)
library(caret)
library(here)
library(httr)
library(jsonlite)
library(readr)
library(stringr)
library(tidyr) # Load tidyr for unnesting
library(writexl)
library(geosphere)
library(caret)

```

## Data

### Vinmonopolet API

Detailed descriptions of each Vinmonopolet store per 2024

```

# Define API URL
url <- "https://apis.vinmonopolet.no/stores/v0/details"

# Define your subscription key (replace with your actual key)
subscription_key <- "3b5b02c6793240fe9e6cb6d176e110e0"

# Send GET request with subscription key in header
response <- GET(url,
  add_headers(
    Accept = "application/json",
    `Ocp-Apim-Subscription-Key` = subscription_key # API authentication
  ))

# Check response status
if (status_code(response) == 200) {
  # Convert API response to JSON and store it

```

```

data <- content(response, as = "text", encoding = "UTF-8")
store_data <- fromJSON(data)

# View first few rows
print(head(store_data))
} else {
  print(paste("Error:", status_code(response)))
}

# ----- Combine API with Vinmonopol Data -----

# Ensure store_data_clean is correctly formatted
store_data_clean <- store_data %>%
  unnest_wider(address) %>% # Expands nested address fields
  select(
    storeId,
    storeName,
    status,
    postalCode,
    city,
    gpsCoord
  ) %>%
  rename(
    Store_ID = storeId,
    Store_Name = storeName,
    Store_Status = status,
    Postal_Code = postalCode,
    City = city,
    GPS_Coordinates = gpsCoord
  )

# Transforming to normal characters
store_data_clean$Store_Name <- iconv(store_data_clean$Store_Name, from = "UTF-8", to = "ASCII")
store_data_clean$Store_Name <- trimws(store_data_clean$Store_Name)

```

## Dimpostnummer merge

As stores have postal codes instead of municipality codes, we need a merge data set inbetween

```
# Read the Dimpostnummer data
dimpostnummer_xlsx <- here("Data", "Vinmonopolet", "dimpostnummer.xlsx")

dimpostnummer_data <- read_excel(dimpostnummer_xlsx) %>%
  select("Postnummer", "Poststed", "Fylke", "KommuneKode", "Kommune")

# Merge with store_data_clean
store_data_clean <- store_data_clean %>%
  left_join(dimpostnummer_data, by = c("Postal_Code" = "Postnummer"))
```

## Vinmonopolet 2024

The “foundation” with store names and sales data for 2024

```
# Set locale to UTF-8
Sys.setlocale("LC_ALL", "en_US.UTF-8")

# Use here package to define the working directory
Vinmonopolet_2024 <- here("Data", "Vinmonopolet", "Vinmonopolet_2024.xlsx")

# Get the names of all sheets in the Excel file
sheet_names <- excel_sheets(Vinmonopolet_2024)

# Read each sheet into a list of data frames, skipping the first row
list_of_dfs <- lapply(sheet_names, function(sheet) {
  read_excel(Vinmonopolet_2024, sheet = sheet, skip = 2)
})

# Combine all data frames into a single data frame
combined_data <- bind_rows(list_of_dfs)

# View the combined data frame
print(combined_data)

# Unique values in the first column
```

```

unique_values <- unique(combined_data$...1)

print(unique_values)

# Transforming to normal characters
combined_data$...1 <- iconv(combined_data$...1, from = "UTF-8", to = "ASCII//TRANSLIT")
combined_data$...1 <- trimws(combined_data$...1)

# Define the values to filter out
values_to_exclude <- c(
  "Svakvin", "Rodvin", "Hvitvin", "Musserende vin", "Rosevin",
  "Perlende vin", "Aromatisert vin", "Sider", "Fruktvin",
  "Brennevin", "Vodka", "Likor", "Whisky", "Akevitt",
  "Brennevin, annet", "Gin", "Druebrennevin",
  "Brennevin, noytralt < 37,5 %", "Rom", "Bitter",
  "Fruktbrennevin", "Genever", "Ol", "Alkoholfritt", "Sterkvin", "Totalsum",
  "eLager"
)

# Column names of combined data
colnames(combined_data)

# Filter out the specified values from the first column
filtered_data <- combined_data %>%
  mutate("2024" = as.numeric(`2024`),
         "Store" = as.character(`...1`)) %>%
  filter(!.[[1]] %in% values_to_exclude) %>%
  select("Store", "2024")

# Export the filtered data to an Excel file
#write_xlsx(filtered_data, "filtered_data.xlsx")

# Standardize store names to improve matching
filtered_data <- filtered_data %>%
  mutate(Store = str_trim(str_to_lower(Store))) # Trim spaces and convert to lowercase

store_data_clean <- store_data_clean %>%
  mutate(Store_Name = str_trim(str_to_lower(Store_Name))) # Trim spaces and convert to lowercase

# Remove unwanted characters from store names

```

```

store_data_clean <- store_data_clean %>%
  mutate(Store_Name = case_when(
    Store_Name == "oslo, thereses gate (stengt ja" ~ "oslo, thereses gate",
    Store_Name == "sandnes, sentrum" ~ "sandnes sentrum",
    Store_Name == "buvika" ~ "buvika, apent 24. oktober",
    Store_Name == "sola, tananger" ~ "sola, tananger, apnet 3. oktober",
    Store_Name == "oslo, bjorvika" ~ "oslo, bjorvika, apnet 14. mars 2024",
    Store_Name == "melhus" ~ "melhus, butikken stengt i 2023 pga kranvelt",
    Store_Name == "bergen, valkendorfsgt." ~ "bergen, valkendorfsgate",
    TRUE ~ Store_Name # This keeps all other values unchanged
  ))

# Merge filtered_data (sales) with store_data_clean (store details)
final_data <- filtered_data %>%
  left_join(store_data_clean, by = c("Store" = "Store_Name")) # Match by store name

# Check merged data
head(final_data)

# Write to excel
write_xlsx(final_data, "final_data.xlsx")

```

## Kommuneendringer 2017

```

### Kommuneendringer 2017 ###

data_df <- final_data %>%
  rename(
    Municipality_Code = KommuneKode,
    Municipality_Name = Kommune
  )

Kommuneendringer_17_xlsx <- here("Data", "Vinmonopolet", "Kommuneendringer_17.xlsx")

kommuneendringer_df <- read_excel(Kommuneendringer_17_xlsx)

# Clean up column names by using the correct column names
colnames(kommuneendringer_df) <- c("New_Codes", "Old_Codes")

```

```

# Split old municipality numbers into separate elements if they are separated by spaces
kommuneendringer_df$Old_Codes <- str_split(kommuneendringer_df$Old_Codes, " ")

# Extract the first four digits from each element in Old_Codes
old_codes_numeric <- lapply(kommuneendringer_df$Old_Codes, function(x) substr(x, 1, 4))

# Create a lookup list that maps old municipality codes to new codes
kommune_mapping <- setNames(rep(kommuneendringer_df$New_Codes, times = sapply(old_codes_numeric,
                                     function(x) length(x))))

# Update Municipality_Code and Municipality_Name in data_df
data_df <- data_df %>%
  rowwise() %>%
  mutate(
    new_val = if (Municipality_Code %in% names(kommune_mapping)) kommune_mapping[[Municipality_Code]]
    Municipality_Code = if (!is.na(new_val)) substr(new_val, 1, 4) else Municipality_Code,
    Municipality_Name = if (!is.na(new_val)) {
      # Remove the municipality number and hyphen from the new value to get the municipality name
      str_trim(str_remove(new_val, "[0-9]{4}-\\s*-\\s*"))
    } else {
      Municipality_Name
    }
  ) %>%
  ungroup() %>%
  select(-new_val)

# Save the updated data to a new Excel file
write_xlsx(data_df, "final_data_17.xlsx")

```

## Kommuneendringer 2018

```

# Read in data from Excel files
Kommuneendringer_18_xlsx <- here("Data", "Vinmonopolet", "Kommuneendringer_18.xlsx")

kommuneendringer_df <- read_excel(Kommuneendringer_18_xlsx)

# Clean up column names by using the correct column names
colnames(kommuneendringer_df) <- c("New_Code", "Old_Codes")

```

```

# Split old municipality numbers (in case multiple old municipalities are separated by space)
kommuneendringer_df$Old_Codes <- str_split(kommuneendringer_df$Old_Codes, " ")

# Create a lookup list for old codes to new codes (one-way mapping)
kommune_mapping <- setNames(rep(kommuneendringer_df$New_Code, times = sapply(kommuneendringer_df$Old_Codes, function(x) strsplit(x, " ")))
                             unlist(kommuneendringer_df$Old_Codes))

# Update both Municipality_Code and Municipality_Name in data_df
data_df <- data_df %>%
  rowwise() %>%
  mutate(
    new_val = if (Municipality_Code %in% names(kommune_mapping)) kommune_mapping[[Municipality_Code]] else Municipality_Code,
    Municipality_Code = if (!is.na(new_val)) substr(new_val, 1, 4) else Municipality_Code,
    Municipality_Name = if (!is.na(new_val)) str_trim(str_remove(new_val, "[0-9]{4}\\s*-\\s*")) else Municipality_Name
  ) %>%
  ungroup() %>%
  select(-new_val)

# Save the updated file
write_xlsx(data_df, "final_data_18.xlsx")

```

## Kommuneendringer 2020

```

# Read in data from Excel files
Kommuneendringer_20_xlsx <- here("Data", "Vinmonopolet", "Kommuneendringer_20.xlsx")

kommuneendringer_df <- read_excel(Kommuneendringer_20_xlsx)

# Clean up column names by using the correct column names
colnames(kommuneendringer_df) <- c("New_Code", "Old_Codes")

# Split old municipality numbers (in case multiple old municipalities are separated by space)
kommuneendringer_df$Old_Codes <- str_split(kommuneendringer_df$Old_Codes, " ")

# Create a lookup list for old codes to new codes (one-way mapping)
kommune_mapping <- setNames(rep(kommuneendringer_df$New_Code, times = sapply(kommuneendringer_df$Old_Codes, function(x) strsplit(x, " ")))
                             unlist(kommuneendringer_df$Old_Codes))

```



```

# Update both Municipality_Code and Municipality_Name in data_df
data_df <- data_df %>%
  rowwise() %>%
  mutate(
    new_val = if (Municipality_Code %in% names(kommune_mapping)) kommune_mapping[[Municipality_Code]] else Municipality_Name,
    Municipality_Code = if (!is.na(new_val)) substr(new_val, 1, 4) else Municipality_Code,
    Municipality_Name = if (!is.na(new_val)) str_trim(str_remove(new_val, "[0-9]{4}\\s*-\\"))
  ) %>%
  ungroup() %>%
  select(-new_val)

# Save the updated file
write_xlsx(data_df, "final_data_20.xlsx")

```

## Kommuneendringer 2024

```

# Read in data from Excel files
Kommuneendringer_24_xlsx <- here("Data", "Vinmonopolet", "Kommuneendringer_24.xlsx")

kommuneendringer_df <- read_excel(Kommuneendringer_24_xlsx)

# Clean up column names by using the correct column names
colnames(kommuneendringer_df) <- c("New_Code", "Old_Codes")

# Split old municipality numbers (in case multiple old municipalities are separated by space)
kommuneendringer_df$Old_Codes <- str_split(kommuneendringer_df$Old_Codes, " ")

# Create a lookup list for old codes to new codes (one-way mapping)
kommune_mapping <- setNames(rep(kommuneendringer_df$New_Code,
                                times = apply(kommuneendringer_df$Old_Codes, length(),
                                                function(x) {
                                                  str_split(x, " ")
                                                })
                                ),
                             unlist(kommuneendringer_df$Old_Codes))

# Update both Municipality_Code and Municipality_Name in data_df
data_df <- data_df %>%
  rowwise() %>%
  mutate(
    new_val = if (Municipality_Code %in% names(kommune_mapping)) kommune_mapping[[Municipality_Code]] else Municipality_Name,
    Municipality_Code = if (!is.na(new_val)) substr(new_val, 1, 4) else Municipality_Code,
    Municipality_Name = if (!is.na(new_val)) str_trim(str_remove(new_val, "[0-9]{4}\\s*-\\"))
  ) %>%
  ungroup() %>%
  select(-new_val)

# Save the updated file
write_xlsx(data_df, "final_data_2024.xlsx")

```

```

    Municipality_Code = if (!is.na(new_val)) substr(new_val, 1, 4) else Municipality_Code,
    Municipality_Name = if (!is.na(new_val)) str_trim(str_remove(new_val, "[0-9]{4}\\s*-\\"))
  ) %>%
  ungroup() %>%
  select(-new_val)

# Hardcode row 121 to set "Municipality_Code" to 1580 and "Municipality_Name" to Haram
data_df[121, "Municipality_Code"] <- "1580"
data_df[121, "Municipality_Name"] <- "Haram"

# Save the updated file
write_xlsx(data_df, "final_data_24.xlsx")

```

## Kommune 2025

Municipality data, including municipality number, population and Area

```

# Kommune data file path
Kommune_data_xlsx <- here("Data", "Vinmonopolet", "Kommune_data.xlsx")

# Read data for total population and area of each municipality
kommune_data <- read_excel(Kommune_data_xlsx, skip = 3) %>%
  rename("Municipality" = "...1",
         "Population" = "2025...2",
         "Area" = "2025...3") %>%
  separate(Municipality, into = c("Mun_num", "Mun_name"), sep = " ", extra = "merge", fill)
  filter(Population != 0,
         Area != 0) %>%
  mutate(Population = as.numeric(Population),
         Area = as.numeric(Area))

# Demographic data file path
Kommune_demo_xlsx <- here("Data", "Vinmonopolet", "Kommune_demo.xlsx")

# Read data for demographic data
demographic_data <- read_excel(Kommune_demo_xlsx, skip = 4) %>%
  rename("Municipality" = "...1",
         "0-17" = "0-17 år",
         "18+" = "18 år eller eldre") %>%

```

```

filter(if_all(everything(), ~ !is.na(.) & . != 0)) %>% # Remove rows with NA or 0 in any
separate(Municipality, into = c("Mun_num", "Mun_name"), sep = " ", extra = "merge", fill
separate(Mun_num, into = c("K", "Mun_num"), sep = "-") %>%
select("-K",
      -"Mun_name")

# Merge the two datasets
kommune_data_final <- kommune_data %>%
  left_join(demographic_data, by = c("Mun_num"))

# Write data to Excel
write_xlsx(kommune_data_final, "Kommune_data_final.xlsx")

```

## Demography data

```

final_data <- data_df

# Transforming to normal characters
final_data$Municipality_Name <- iconv(final_data$Municipality_Name, from = "UTF-8", to = "A

final_data$Municipality_Name <- trimws(final_data$Municipality_Name)

# Standardize store names to improve matching
final_data <- final_data %>%
  mutate(Municipality_Name = str_trim(str_to_lower(Municipality_Name))) # Trim spaces and

# Loading the kommune data
kommune_data <- kommune_data_final

# Standardize the kommune data
kommune_data <- kommune_data %>%
  mutate(Mun_name = iconv(Mun_name, from = "UTF-8", to = "ASCII//TRANSLIT"),
         Mun_name = str_trim(str_to_lower(Mun_name))) # Trim spaces and convert to lowerca

# Perform a full join to include all rows from both datasets
merged_data <- final_data %>%
  full_join(kommune_data, by = c("Municipality_Code" = "Mun_num"))

```

```

# Replace NA values in store-related columns with 0
# Assuming 'Store_Info_Column' is the column in final_data that contains store information
# Replace 'Store_Info_Column' with the actual column names you want to fill with 0
merged_data <- merged_data %>%
  mutate(across(where(is.numeric), ~ replace_na(.x, 0)))

# If you have specific columns to replace NA with 0, you can specify them like this:
# merged_data <- merged_data %>%
#   mutate(Store_Info_Column = replace_na(Store_Info_Column, 0))

# Write the merged data to an Excel file
#write_xlsx(merged_data, "final_data_mun.xlsx")

```

## Distance data

This is just our code for the calculation of `dist_nearest`. As the actual data file is too large to submit, we jump to the next step with the resulting data saved as “final\_data\_mun\_dist.xlsx”

```

# -----
# 1. Load and Prepare Data
# -----
# Load Vinmonopolet + municipality dataset
#data <- read_excel("final_data_mun.xlsx")

# Load pre-cleaned municipality admin center coordinates
#admin_centers_final <- readRDS("admin_centers_final.rds")

# Ensure join columns match in type
#data <- data %>%
#  mutate(Municipality_Code = as.character(Municipality_Code))

#admin_centers_final <- admin_centers_final %>%
#  mutate(kommunennummer = as.character(kommunennummer))

# -----
# 2. Merge Coordinates
# -----
# Merge admin center lat/lon into dataset by municipality

```

```

#data <- left_join(data, admin_centers_final, by = c("Municipality_Code" = "kommunennummer")

# Overwrite old coordinates with admin center coordinates
#data <- data %>%
#  mutate(
#    Latitude = as.numeric(lat),
#    Longitude = as.numeric(lon)
#  )

# -----
# 3. Parse Store GPS Coordinates
# -----
# Split store GPS into separate numeric lat/lon

# -----
# STEP 1: Load dataset
# -----

# Read merged dataset with both Vinmonopolet store info and municipality info
#data <- read_excel("final_data_mun.xlsx")

# -----
# STEP 2: Parse store coordinates
# -----

# GPS_Coordinates column contains both latitude and longitude as a string separated by ";"
# We split this into two separate numeric columns: store_lat and store_lon

#data <- data %>%
#  separate(GPS_Coordinates, into = c("store_lat", "store_lon"), sep = ";", convert = TRUE)
#  mutate(
#    store_lat = as.numeric(store_lat), # ensure store latitude is numeric
#    store_lon = as.numeric(store_lon)  # ensure store longitude is numeric
#  )

# -----
# 4. Build Store Location Matrix
# -----
# Extract distinct (lon, lat) of all Vinmonopolet stores

```

```

#store_locations <- data %>%
#  filter(!is.na(store_lon), !is.na(store_lat)) %>%

# -----
# STEP 3: Ensure municipality center coordinates are numeric
# -----

# These are already separate in the dataset, but stored as characters - we convert them
#data <- data %>%
#  mutate(
#    Longitude = as.numeric(Longitude), # longitude of the municipality center
#    Latitude = as.numeric(Latitude)    # latitude of the municipality center
#  )

# -----
# STEP 4: Extract store coordinates for distance calculation
# -----

# We only want to use valid store locations for calculating distances
# (some rows in the dataset are just municipality data with no store info)
#store_data <- data %>%
#  filter(!is.na(store_lat), !is.na(store_lon))

# Extract a unique matrix of all Vinmonopolet store locations
# Format required by geosphere is matrix of (longitude, latitude)
#store_locations <- store_data %>%

#  select(store_lon, store_lat) %>%
#  distinct() %>%
#  as.matrix()

# -----
# STEP 5: Define function to calculate distance to nearest store
# -----

# For a given municipality center (lon, lat), compute distance to nearest store
# Uses Haversine formula (accounts for Earth's curvature)
#min_distance_to_store <- function(lon, lat) {
#  if (is.na(lon) || is.na(lat)) {
#    return(NA) # return NA if municipality coordinates are missing

```

```

#   }
#   muni_coord <- matrix(c(lon, lat), nrow = 1) # convert to matrix format for geosphere
#   dists <- distHaversine(muni_coord, store_locations) # distances in meters
#   return(min(dists) / 1000) # convert to kilometers
#}

# -----
# STEP 6: Apply distance function to each municipality
# -----

# For each row (i.e., each municipality center), calculate distance to closest Vinmonopolet
# Note: This includes all rows (even ones without a store)

#data$dist_nearest_store <- mapply(
#  min_distance_to_store,
#  data$Longitude,
#  data$Latitude
#)

# -----
# STEP 7: Quick check (optional)
# -----

# Check that coordinates are numeric
#str(data$Longitude)
#str(data$Latitude)

# -----
# 7. Optional: Drop Redundant Columns
# -----

#data <- data %>%
#  select(
#    -lat, -lon, -multikurve, -kommunenavn
#  )

# -----
# 8. Final Checks (Optional)
# -----

```

```

#str(data$dist_nearest_store)
#summary(data$dist_nearest_store)

# -----
# 9. Does Vinmonopolets 30 km threshold 97% goal work based on our data
# -----
# 1. Total population (all municipalities)
#total_pop <- sum(data$Population, na.rm = TRUE)

# 2. Population in municipalities with distance > 30 km
#pop_far_away <- data %>%
#  filter(dist_nearest_store > 30) %>%
#  summarise(total = sum(Population, na.rm = TRUE)) %>%
#  pull(total)

# 3. Share of population far away
#share_far_away <- pop_far_away / total_pop

# 4. Share WITH access (within 30 km)
#share_within_30km <- 1 - share_far_away

# 5. Print results
#cat(sprintf("Share of population within 30 km of a Vinmonopolet: %.2f%%\n", #share_within_
#cat(sprintf("Target (Vinmonopolet): 97%%\n"))

#underserved <- data %>%
#  filter(dist_nearest_store > 30) %>%
#  select(,Mun_name, Population, dist_nearest_store) %>%
#  arrange(desc(dist_nearest_store))

#print(underserved, n = 50)

# -----
# 10. Export the final data to an Excel file
# -----

#library(writexl)
#write_xlsx(data, "final_data_mun_dist.xlsx")
# -----

```



## Model variables merge

```
### Independent variables merge ###

final_data_mun_dist <- here("Data", "Vinmonopolet", "final_data_mun_dist.xlsx")

# Load data
Vinmonopolet <- read_excel(final_data_mun_dist) %>%
  select(-c(Store_ID, Store_Status, Postal_Code, Poststed,
            PostnummerKategoriKode, PostnummerKategori, Region_Code,
            Municipality_Name)) %>%
  mutate(
    Municipality_Name = Mun_name,
    Region_Name = case_when(
      Region_Name == "AUST-AGDER" ~ "Agder",
      Region_Name == "VEST-AGDER" ~ "Agder",
      Region_Name == "AKERSHUS" ~ "Akershus",
      Region_Name == "OPPLAND" ~ "Innlandet",
      Region_Name == "BUSKERUD" ~ "Buskerud",
      Region_Name == "VESTFOLD" ~ "Vestfold",
      Region_Name == "FINNMARK" ~ "Finnmark",
      Region_Name == "HEDMARK" ~ "Innlandet",
      Region_Name == "MØRE OG ROMSDAL" ~ "Møre og Romsdal",
      Region_Name == "NORDLAND" ~ "Nordland",
      Region_Name == "OSLO" ~ "Oslo",
      Region_Name == "ROGALAND" ~ "Rogaland",
      Region_Name == "TELEMARK" ~ "Telemark",
      Region_Name == "TROMS" ~ "Troms",
      Region_Name == "SØR-TRØNDELAG" ~ "Trøndelag",
      Region_Name == "NORD-TRØNDELAG" ~ "Trøndelag",
      Region_Name == "SOGN OG FJORDANE" ~ "Vestland",
      Region_Name == "HORDALAND" ~ "Vestland",
      Region_Name == "ØSTFOLD" ~ "Østfold",
      is.na(Region_Name) & str_starts(Municipality_Code, "03") ~ "Oslo",
      is.na(Region_Name) & str_starts(Municipality_Code, "11") ~ "Rogaland",
      is.na(Region_Name) & str_starts(Municipality_Code, "15") ~ "Møre og Romsdal",
      is.na(Region_Name) & str_starts(Municipality_Code, "18") ~ "Nordland",
      is.na(Region_Name) & str_starts(Municipality_Code, "31") ~ "Østfold",
      is.na(Region_Name) & str_starts(Municipality_Code, "32") ~ "Akershus",
      is.na(Region_Name) & str_starts(Municipality_Code, "33") ~ "Buskerud",
```

```

    is.na(Region_Name) & str_starts(Municipality_Code, "34") ~ "Innlandet",
    is.na(Region_Name) & str_starts(Municipality_Code, "39") ~ "Vestfold",
    is.na(Region_Name) & str_starts(Municipality_Code, "40") ~ "Telemark",
    is.na(Region_Name) & str_starts(Municipality_Code, "42") ~ "Agder",
    is.na(Region_Name) & str_starts(Municipality_Code, "46") ~ "Vestland",
    is.na(Region_Name) & str_starts(Municipality_Code, "50") ~ "Trøndelag",
    is.na(Region_Name) & str_starts(Municipality_Code, "55") ~ "Troms",
    is.na(Region_Name) & str_starts(Municipality_Code, "56") ~ "Finnmark",
    TRUE ~ Region_Name # Keep existing Region_Name if no conditions are met
  )
) %>%
select(-Mun_name)

# Aggregating per municipality data
Vinmonopolet_market <- Vinmonopolet %>%
  group_by(Municipality_Code) %>%
  summarise(
    Mun_name = first(Municipality_Name),
    Region_Name = first(Region_Name),
    Population = first(Population),
    Area = first(Area),
    Number_of_stores = sum(`2024` > 0), # Count non-zero sales
    Sales = sum(`2024`),
    Lat = first(Latitude),
    Lon = first(Longitude),
    Dist_nearest = first(dist_nearest_store),
  )

# Scaling the variables that have not been scaled yet
Vinmonopolet_market <- Vinmonopolet_market %>%
  mutate(Population = Population / 1000,
         Sales = Sales / 1000)

# Now we have loaded and wrangled the main data set, but we can use some
# new variables for our analysis

## Merge 1: Grensehandel #####

Grensehandel_weights <- here("Data", "Vinmonopolet", "Grensehandel_weights.xlsx")

```

```

# Load the weights datas
weights <- read_excel(Grensehandel_weights, skip = 3) %>%
  slice(1) %>%
  select(-'...1') %>%
  mutate(
    mean_weight = (as.numeric(`2024K1`) + as.numeric(`2024K2`) + as.numeric(`2024K3`) + as.
  )

weight_grensehandel <- weights$mean_weight / 100

# Load the regional data
Grensehandel_regions <- here("Data", "Vinmonopolet", "Grensehandel_regions.xlsx")

regional <- read_excel(Grensehandel_regions)

total_grensehandel <- sum(regional$"2024")

# Calculate grensehandel per region
regional <- regional %>%
  rename(
    Region = `Fylker`,
    Total_sale = `2024`
  ) %>%
  mutate(
    Grensehandel = Total_sale * weight_grensehandel
  )

# Split the "Vestlandet" region row into three new rows: "Rogaland", "Vestland" and "MC8re
regional <- regional %>%
  rbind(
    regional %>% filter(Region == "Vestlandet") %>% mutate(Region = "Rogaland"),
    regional %>% filter(Region == "Vestlandet") %>% mutate(Region = "Vestland"),
    regional %>% filter(Region == "Vestlandet") %>% mutate(Region = "Møre og Romsdal")
  ) %>%
  filter(Region != "Vestlandet")

# Divide the grensehandel value by three for "Rogaland", "Vestland" and "MC8re og Romsdal"
regional <- regional %>%
  mutate(
    Grensehandel = case_when(

```

```

    Region == "Rogaland" ~ Grensehandel * 0.35,
    Region == "Vestland" ~ Grensehandel * 0.46,
    Region == "Møre og Romsdal" ~ Grensehandel * 0.19,
    TRUE ~ Grensehandel # Keep the original value for other regions
  )
)

# Split the "Nord-Norge" region row into three new rows: "Nordland", "Troms" and "Finnmark"
# And divide the grensehandel value by three
regional <- regional %>%
  mutate(
    Grensehandel = ifelse(Region == "Nord-Norge", Grensehandel / 3, Grensehandel)
  ) %>%
  rbind(
    regional %>% filter(Region == "Nord-Norge") %>% mutate(Region = "Nordland"),
    regional %>% filter(Region == "Nord-Norge") %>% mutate(Region = "Troms"),
    regional %>% filter(Region == "Nord-Norge") %>% mutate(Region = "Finnmark")
  ) %>%
  filter(Region != "Nord-Norge")

# Divide the grensehandel value by three for "Nordland", "Troms" and "Finnmark"
regional <- regional %>%
  mutate(
    Grensehandel = case_when(
      Region == "Nordland" ~ Grensehandel * 0.5,
      Region == "Troms" ~ Grensehandel * 0.35,
      Region == "Finnmark" ~ Grensehandel * 0.15,
      TRUE ~ Grensehandel # Keep the original value for other regions
    )
  )

# Split the "Agder, Telemark, Buskerud og Vestfold" column into four new columns: "Agder",
# And divide the grensehandel value by four
regional <- regional %>%
  mutate(
    Grensehandel = ifelse(Region == "Agder, Telemark, Buskerud og Vestfold", Grensehandel /
  ) %>%
  rbind(
    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region = "Agder"),
    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region = "Telemark"),
    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region = "Buskerud og Vestfold"),
    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region = "Agder, Telemark, Buskerud og Vestfold")
  ) %>%
  filter(Region != "Agder, Telemark, Buskerud og Vestfold")

```

```

    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region
    regional %>% filter(Region == "Agder, Telemark, Buskerud og Vestfold") %>% mutate(Region
  ) %>%
  filter(Region != "Agder, Telemark, Buskerud og Vestfold")

# Divide the grensehandel value by four for "Agder", "Telemark", "Buskerud" and "Vestfold"
regional <- regional %>%
  mutate(
    Grensehandel = case_when(
      Region == "Agder" ~ Grensehandel * 0.31,
      Region == "Telemark" ~ Grensehandel * 0.17,
      Region == "Buskerud" ~ Grensehandel * 0.26,
      Region == "Vestfold" ~ Grensehandel * 0.26,
      TRUE ~ Grensehandel # Keep the original value for other regions
    )
  )

# Removing the "total_sale" column from the regional data set
regional <- regional %>% select(-Total_sale)

# Merge the regional data with the main data set on Region_Name in the Vinmonopolet_market
Vinmonopolet_market <- left_join(Vinmonopolet_market, regional, by = c("Region_Name" = "Reg

# Add a new column "Region_pop" where "Population" is summarized for each region
Vinmonopolet_market <- Vinmonopolet_market %>%
  group_by(Region_Name) %>%
  mutate(Region_pop = sum(Population)) %>%
  ungroup()

Vinmonopolet_market <- Vinmonopolet_market %>%
  mutate(Kommune_share = Population / Region_pop,
    Grensehandel_mun = Grensehandel * Kommune_share) %>%
  select(-c("Region_pop", "Kommune_share", "Grensehandel")) %>%
  rename(Grensehandel = Grensehandel_mun)

## Merge 2: Tourism #####

```

```

Tourism_xlsx <- here("Data", "Vinmonopolet", "Tourism.xlsx")

# Reading tourism data
Tourism <- read_excel(Tourism_xlsx, skip = 4) %>%
  rename(
    Mun = '...1',
    H = 'Hotell og liknande overnattingsbedrifter',
    C = 'Campingplassar, hyttegrender og vandrarheim',
  ) %>%
  select(-'...2') %>%
  mutate_at(vars(H, C), ~as.numeric(str_replace_all(., ":", "0"))) %>%
  mutate(n_stays = H + C) %>%
  separate(Mun, into = c("Municipality_Code", "Municipality_Name"), sep = " ", remove = FALSE) %>%
  select(-c("Mun", "H", "C", "Municipality_Name")) %>%
  filter(!is.na(Municipality_Code))

# Merging the data
Vinmonopolet_market <- left_join(Vinmonopolet_market, Tourism, by = "Municipality_Code") %>%
  mutate(
    n_stays = ifelse(is.na(n_stays), 0, n_stays),
    n_stays = n_stays / 1000
  )

# There is a great deal of missing data, so we do not know the relevance of
# this data yet

## Merge 3: Income #####

# Average monthly salary per inhabitant in the municipality

# Load data
Monthly_Salary <- here("Data", "Vinmonopolet", "Monthly_Salary.xlsx")

data <- read_excel(Monthly_Salary)

# Cleaning data by removing rows with missing values and rows with dots
clean_data <- data %>%
  filter(!apply(., 1, function(row) any(grepl("\\\\.", row)))) %>%
  na.omit()

```

```

# Remove the last two rows from the data, using tidyverse
clean_data <- clean_data %>%
  slice(1:(n() - 2)) %>%
  select(-'...2') %>%
  rename(
    Mun = `12852: Kommunefordelt månedslønn, etter region, statistikk mål, statistikkvariabel`
    Monthly_salary = '...3'
  ) %>%
  separate(Mun, into = c("Municipality_Code", "Municipality_Name"), sep = " ", remove = FALSE)
  select(-c("Municipality_Name", "Mun")) %>%
  mutate(Monthly_salary = as.numeric(Monthly_salary),
    Monthly_salary = Monthly_salary / 1000)

# Merge with the main data set
Vinmonopolet_market <- left_join(Vinmonopolet_market, clean_data, by = "Municipality_Code")

## Merge 4: Concentration #####

# Load data
Concentration_xlsx <- here("Data", "Vinmonopolet", "Concentration.xlsx")

concentration <- read_excel(Concentration_xlsx, skip = 5) %>%
  slice(1:357) %>%
  select('...1',
    'Spredtbygd strøk...3') %>%
  rename(Mun = '...1',
    Spread = 'Spredtbygd strøk...3') %>%
  separate(Mun, into = c("Municipality_Code", "Municipality_Name"), sep = " ", remove = FALSE)
  select(-c("Municipality_Name", "Mun")) %>%
  mutate(Spread = as.numeric(Spread),
    Spread = Spread / 1000)

# Remove the first two characters of each cell in the "Municipality_Code" column
concentration$Municipality_Code <- substr(concentration$Municipality_Code, 3, nchar(concentration$Municipality_Code))

# Merge with the main data set

```

```

Vinmonopolet_market <- left_join(Vinmonopolet_market, concentration, by = "Municipality_Cod

## Merge 5: "Active" stores #####

# Load data
Active_xlsx <- here("Data", "Vinmonopolet", "Active.xlsx")

A1 <- read_excel(Active_xlsx, sheet = 1, skip = 2)

A2 <- read_excel(Active_xlsx, sheet = 2, skip = 2)

# Merge the two data sets
Active <- A1 %>%
  bind_rows(A2) %>%
  select(-c('1', '...3', Fylke))

# Rename columns
names(Active)[1] <- "Mun_name"

# Remove unnecessary spaces and numbers from the "Mun_name" column
Active$Mun_name <- substr(Active$Mun_name, 4, nchar(Active$Mun_name))

Active$Mun_name <- trimws(Active$Mun_name, which = "left")

# Replace norwegian special letters with english ones and make all letters lowercase
Active$Mun_name <- tolower(iconv(Active$Mun_name, from = "UTF-8", to = "ASCII//TRANSLIT"))

# Recode the "Mun_name" column
Active$Mun_name <- case_when(
  Active$Mun_name == "hamaroy" ~ "habmer - hamaroy",
  Active$Mun_name == "hattfjelldal" ~ "aarborte - hattfjelldal",
  Active$Mun_name == "valer (viken)" ~ "valer (ostfold)",
  TRUE ~ Active$Mun_name)

# Merge with the main data set

# Make a dummy variable for active stores
Vinmonopolet_market$Active <- ifelse(Vinmonopolet_market$Mun_name %in% Active$Mun_name, 1,

```



```
## Write to Excel #####

# Write to Excel
# write_xlsx(Vinmonopolet_market, "demand_data.xlsx")
```

## Model applications

### Data preparation

```
### Data preparation #####

# Narrowing down the data to only contain relevant markets
# Excluding the largest cities because they are not representative

# Train and test split, training data all observations with a store
train_data <- Vinmonopolet_market %>%
  filter(Number_of_stores > 0)

# Test data all observations without a store
test_data <- Vinmonopolet_market %>%
  filter(Number_of_stores == 0)
```

### Model selection and basic regressions

```
### Model selection #####

# Forward selection
forward_model <- step(lm(Sales ~ 1, data = train_data),
  scope = ~ Population + Grensehandel + n_stays + Monthly_salary + Area
  direction = "forward")
```

Start: AIC=3281.16

Sales ~ 1

	Df	Sum of Sq	RSS	AIC
--	----	-----------	-----	-----

+ Population	1	238655412	3275767	2263.6
+ Number_of_stores	1	232719071	9212107	2508.6
+ n_stays	1	213396224	28534954	2776.6
+ Grensehandel	1	121971827	119959351	3116.9
+ Monthly_salary	1	46654218	195276960	3232.4
+ Spread	1	4307259	237623919	3278.9
<none>			241931178	3281.2
+ Area	1	1750943	240180236	3281.4

Step: AIC=2263.56

Sales ~ Population

	Df	Sum of Sq	RSS	AIC
+ Number_of_stores	1	655940	2619827	2212.6
+ n_stays	1	611748	2664019	2216.6
+ Grensehandel	1	487930	2787837	2227.3
+ Spread	1	228325	3047442	2248.4
+ Area	1	76705	3199061	2259.9
<none>			3275767	2263.6
+ Monthly_salary	1	2240	3273527	2265.4

Step: AIC=2212.6

Sales ~ Population + Number\_of\_stores

	Df	Sum of Sq	RSS	AIC
+ Grensehandel	1	267494	2352332	2189.1
+ n_stays	1	250512	2369314	2190.8
+ Spread	1	45130	2574697	2210.5
+ Area	1	44973	2574854	2210.5
+ Monthly_salary	1	35351	2584475	2211.4
<none>			2619827	2212.6

Step: AIC=2189.08

Sales ~ Population + Number\_of\_stores + Grensehandel

	Df	Sum of Sq	RSS	AIC
+ n_stays	1	116034	2236298	2179.1
+ Area	1	32495	2319837	2187.8
+ Spread	1	32037	2320296	2187.8
+ Monthly_salary	1	25822	2326510	2188.5
<none>			2352332	2189.1

Step: AIC=2179.09

Sales ~ Population + Number\_of\_stores + Grensehandel + n\_stays

	Df	Sum of Sq	RSS	AIC
+ Monthly_salary	1	94136	2142163	2170.9
+ Spread	1	19793	2216505	2179.0
<none>			2236298	2179.1
+ Area	1	13717	2222581	2179.6

Step: AIC=2170.9

Sales ~ Population + Number\_of\_stores + Grensehandel + n\_stays +  
Monthly\_salary

	Df	Sum of Sq	RSS	AIC
+ Area	1	36877	2105286	2168.8
<none>			2142163	2170.9
+ Spread	1	8522	2133640	2171.9

Step: AIC=2168.78

Sales ~ Population + Number\_of\_stores + Grensehandel + n\_stays +  
Monthly\_salary + Area

	Df	Sum of Sq	RSS	AIC
<none>			2105286	2168.8
+ Spread	1	9654.8	2095631	2169.7

```
summary(forward_model)
```

Call:

```
lm(formula = Sales ~ Population + Number_of_stores + Grensehandel +  
n_stays + Monthly_salary + Area, data = train_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-338.11	-48.12	-6.81	35.54	364.74

Coefficients:

Estimate	Std. Error	t value	Pr(> t )
----------	------------	---------	----------

```

(Intercept)    -4.980e+02  1.249e+02  -3.986  9.03e-05 ***
Population      1.358e+01  6.778e-01  20.040  < 2e-16 ***
Number_of_stores 7.744e+01  1.329e+01   5.825  1.91e-08 ***
Grensehandel   -4.198e+00  1.369e+00  -3.068  0.00242 **
n_stays         1.780e-01  4.175e-02   4.263  2.94e-05 ***
Monthly_salary  7.956e+00  2.222e+00   3.580  0.00042 ***
Area            1.278e-02  6.366e-03   2.007  0.04590 *

```

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 95.67 on 230 degrees of freedom

Multiple R-squared: 0.9913, Adjusted R-squared: 0.9911

F-statistic: 4367 on 6 and 230 DF, p-value: < 2.2e-16

```
# Backward selection
```

```
backward_model <- step(lm(Sales ~ Population + Grensehandel + n_stays + Monthly_salary + Area +
                           Number_of_stores + Spread,
                           data = train_data),
                       direction = "backward")
```

Start: AIC=2169.69

Sales ~ Population + Grensehandel + n\_stays + Monthly\_salary +  
Area + Number\_of\_stores + Spread

	Df	Sum of Sq	RSS	AIC
- Spread	1	9655	2105286	2168.8
<none>			2095631	2169.7
- Area	1	38010	2133640	2171.9
- Grensehandel	1	86637	2182268	2177.3
- Monthly_salary	1	104814	2200445	2179.3
- n_stays	1	148519	2244150	2183.9
- Number_of_stores	1	253669	2349299	2194.8
- Population	1	3321080	5416711	2392.8

Step: AIC=2168.78

Sales ~ Population + Grensehandel + n\_stays + Monthly\_salary +  
Area + Number\_of\_stores

	Df	Sum of Sq	RSS	AIC
<none>			2105286	2168.8
- Area	1	36877	2142163	2170.9

```
- Grensehandel      1      86135 2191421 2176.3
- Monthly_salary    1      117295 2222581 2179.6
- n_stays            1      166368 2271653 2184.8
- Number_of_stores  1      310561 2415847 2199.4
- Population         1     3676078 5781364 2406.2
```

```
summary(backward_model)
```

Call:

```
lm(formula = Sales ~ Population + Grensehandel + n_stays + Monthly_salary +
    Area + Number_of_stores, data = train_data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-338.11	-48.12	-6.81	35.54	364.74

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-4.980e+02	1.249e+02	-3.986	9.03e-05 ***
Population	1.358e+01	6.778e-01	20.040	< 2e-16 ***
Grensehandel	-4.198e+00	1.369e+00	-3.068	0.00242 **
n_stays	1.780e-01	4.175e-02	4.263	2.94e-05 ***
Monthly_salary	7.956e+00	2.222e+00	3.580	0.00042 ***
Area	1.278e-02	6.366e-03	2.007	0.04590 *
Number_of_stores	7.744e+01	1.329e+01	5.825	1.91e-08 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 95.67 on 230 degrees of freedom

Multiple R-squared: 0.9913, Adjusted R-squared: 0.9911

F-statistic: 4367 on 6 and 230 DF, p-value: < 2.2e-16

```
lm_Area <- lm(Sales ~ Area, data = train_data)
```

```
summary(lm_Area)
```

Call:

```
lm(formula = Sales ~ Area, data = train_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-426.0	-304.2	-213.1	-23.7	12804.5

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	463.91339	88.30921	5.253	3.35e-07 ***
Area	-0.08325	0.06360	-1.309	0.192

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1011 on 235 degrees of freedom

Multiple R-squared: 0.007237, Adjusted R-squared: 0.003013

F-statistic: 1.713 on 1 and 235 DF, p-value: 0.1919

```
lm_pop <- lm(Sales ~ Population, data = Vinmonopolet_market)
summary(lm_pop)
```

Call:

```
lm(formula = Sales ~ Population, data = Vinmonopolet_market)
```

Residuals:

Min	1Q	Median	3Q	Max
-595.79	-25.26	5.15	33.18	417.61

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-27.397	5.466	-5.012	8.52e-07 ***
Population	18.128	0.112	161.791	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 97.81 on 355 degrees of freedom

Multiple R-squared: 0.9866, Adjusted R-squared: 0.9866

F-statistic: 2.618e+04 on 1 and 355 DF, p-value: < 2.2e-16

```

small_data <- Vinmonopolet_market %>%
  filter(Number_of_stores == 1 | 0)

lm_pop_test <- lm(Sales ~ Population, data = small_data)

summary(lm_pop_test)

```

Call:

```
lm(formula = Sales ~ Population, data = small_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-309.98	-36.54	-5.05	38.74	346.01

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	37.038	7.977	4.643	6.31e-06 ***
Population	12.706	0.580	21.905	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 71.84 on 194 degrees of freedom

Multiple R-squared: 0.7121, Adjusted R-squared: 0.7106

F-statistic: 479.8 on 1 and 194 DF, p-value: < 2.2e-16

```

# Linear regression model for predicting sales with all the variables
var_test <- lm(Sales ~ Population + Grensehandel + n_stays + Monthly_salary + Area +
  Number_of_stores + Spread,
  data = Vinmonopolet_market)

stargazer(var_test, type = "text")

```

```

=====
Dependent variable:
-----
Sales
-----

```

Population	14.497*** (0.490)
Grensehandel	-4.610*** (1.117)
n_stays	0.159*** (0.034)
Monthly_salary	4.740*** (1.518)
Area	0.010** (0.005)
Number_of_stores	63.026*** (8.444)
Spread	-6.040*** (2.236)
Constant	-289.696*** (83.210)

```
-----
Observations      357
R2                0.991
Adjusted R2       0.991
Residual Std. Error 79.155 (df = 349)
F Statistic      5,737.855*** (df = 7; 349)
=====
Note:             *p<0.1; **p<0.05; ***p<0.01
```

```
# From these regressions we see that we want to remove the "Area" and "prop_spread" variables
# from the regressions as they are not significant.
```

## Demand estimation



```

### Demand estimation #####

## Linear regression

# Predicting sales using the training data
reg1 <- lm(Sales ~ Population + Grensehandel + n_stays + Monthly_salary,
           data = train_data)

summary(reg1)

```

Call:

```
lm(formula = Sales ~ Population + Grensehandel + n_stays + Monthly_salary,
    data = train_data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-435.51	-46.33	0.49	42.36	402.63

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-270.08466	125.33620	-2.155	0.032201 *
Population	16.44887	0.49049	33.535	< 2e-16 ***
Grensehandel	-5.26373	1.45612	-3.615	0.000369 ***
n_stays	0.24607	0.04315	5.703	3.56e-08 ***
Monthly_salary	4.88463	2.28962	2.133	0.033944 *

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 102.6 on 232 degrees of freedom

Multiple R-squared: 0.9899, Adjusted R-squared: 0.9897

F-statistic: 5683 on 4 and 232 DF, p-value: < 2.2e-16

```

# Applying the model on the test data
test_data$Sales_pred <- predict(reg1, newdata = test_data)

## Merge predicted data into the original data

# Deselect unnecessary columns to merge the data easier

```

```

test_data <- test_data %>%
  select(Municipality_Code, Sales_pred)

# Merge predicted demand (sales) back into the original data
Vinmonopolet_market <- Vinmonopolet_market %>%
  left_join(test_data, by = "Municipality_Code") %>%
  mutate(Sales = ifelse(Sales == 0, Sales_pred, Sales)) %>%
  select(-Sales_pred) %>%
  mutate(Sales = ifelse(Sales < 0, 0, Sales),
         Number_of_stores = as.integer(Number_of_stores)) %>%
  filter(Number_of_stores < 2)

```

## Logit model

```

## Logit regression #####

# Make sure the factor for Number_of_stores has valid R variable names
# that won't cause errors in caret. For instance, rename "0" -> "NoStore"
# and "1" -> "OneStore".
data_for_logit <- Vinmonopolet_market %>%
  mutate(Number_of_stores = as.factor(Number_of_stores))

# Rename factor levels (originally "0" and "1") to "NoStore" and "OneStore"
data_for_logit$Number_of_stores <- factor(
  data_for_logit$Number_of_stores,
  levels = c("0", "1"),
  labels = c("NoStore", "OneStore")
)

# Set up k-fold cross-validation parameters
set.seed(123) # for reproducibility

my_control <- trainControl(
  method = "cv",           # k-fold CV
  number = 5,              # 5 folds
  classProbs = TRUE,       # needed for probability output
  summaryFunction = twoClassSummary
)

```

```
# Train the logistic model with cross-validation
cv_model <- train(
  Number_of_stores ~ Sales,
  data = data_for_logit,
  method = "glm",
  family = binomial,
  trControl = my_control,
  metric = "ROC"          # use AUC (Area Under the Curve) as our metric
)

# Review cross-validation results
print(cv_model)
```

### Generalized Linear Model

```
316 samples
  1 predictor
  2 classes: 'NoStore', 'OneStore'
```

No pre-processing

Resampling: Cross-Validated (5 fold)

Summary of sample sizes: 252, 253, 253, 253, 253

Resampling results:

ROC	Sens	Spec
0.9397703	0.8666667	0.8670513

```
print(cv_model$results)
```

	parameter	ROC	Sens	Spec	ROCSD	SensSD	SpecSD
1	none	0.9397703	0.8666667	0.8670513	0.02518379	0.1078515	0.04987575

```
# Get predicted probabilities from the final trained model
```

```
# caret retrains on the entire dataset after CV by default
```

```
Vinmonopolet_market$prob <- predict(cv_model, newdata = data_for_logit, type = "prob")[, "0"]
```

```
# Use the probabilities for your recommendations
```

```
recommended_stores <- Vinmonopolet_market %>%
```

```
mutate(Number_of_stores = as.integer(as.character(Number_of_stores))) %>%
filter(Number_of_stores == 0, Dist_nearest > 0) %>%
arrange(desc(prob)) %>%
select(Mun_name, prob, Dist_nearest, Sales, Population, Region_Name, Active)

head(recommended_stores, 10) # for example, show top 10
```

# A tibble: 10 x 7

	Mun_name <chr>	prob <dbl>	Dist_nearest <dbl>	Sales <dbl>	Population <dbl>	Region_Name <chr>	Active <dbl>
1	giske	0.995	3.41	156.	8.77	Møre og Romsdal	0
2	lunner	0.993	7.32	150.	9.42	Akershus	1
3	rade	0.967	12.9	118.	7.85	Østfold	0
4	hareid	0.857	9.83	87.2	5.32	Møre og Romsdal	1
5	valer (ostfold)	0.823	10.8	82.1	6.16	Østfold	1
6	aurland	0.817	36.6	81.3	1.84	Vestland	1
7	birkenes	0.752	12.7	73.8	5.41	Agder	0
8	eidskog	0.729	23.3	71.5	6.06	Innlandet	1
9	aure	0.683	28.0	67.0	3.39	Møre og Romsdal	1
10	austrheim	0.658	16.3	64.8	2.92	Vestland	1

```
# Output the top 10 recommended stores as a nice table using kable
# And save it
kable(head(recommended_stores, 10), format = "markdown")
```

Mun_name	prob	Dist_nearest	Sales	Population	Region_Name	Active
giske	0.9950643	3.407605	156.13987	8.773	Møre og Romsdal	0
lunner	0.9932445	7.323172	149.94498	9.420	Akershus	1
rade	0.9672041	12.891027	118.41978	7.850	Østfold	0
hareid	0.8573922	9.828661	87.21263	5.320	Møre og Romsdal	1
valer (ostfold)	0.8226703	10.788931	82.12523	6.162	Østfold	1
aurland	0.8166871	36.608030	81.33082	1.836	Vestland	1
birkenes	0.7517656	12.695963	73.75601	5.413	Agder	0
eidskog	0.7294726	23.342412	71.47769	6.059	Innlandet	1
aure	0.6827221	28.029241	67.04993	3.394	Møre og Romsdal	1

Mun_name	prob	Dist_nearest	Sales	Population	Region_Name	Active
austrheim	0.6575973	16.315532	64.81866	2.915	Vestland	1

```
# Use the probabilities for your recommendations
Active_stores <- Vinmonopolet_market %>%
  mutate(Number_of_stores = as.integer(as.character(Number_of_stores))) %>%
  filter(Active == 1, Dist_nearest > 0) %>%
  arrange(desc(prob)) %>%
  select(Mun_name, prob, Dist_nearest, Sales, Population, Region_Name, Active)

head(Active_stores, 10) # for example, show top 10
```

# A tibble: 10 x 7

	Mun_name	prob	Dist_nearest	Sales	Population	Region_Name	Active
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<dbl>
1	lunner	0.993	7.32	150.	9.42	Akershus	1
2	hareid	0.857	9.83	87.2	5.32	Møre og Romsdal	1
3	valer (ostfold)	0.823	10.8	82.1	6.16	Østfold	1
4	aurland	0.817	36.6	81.3	1.84	Vestland	1
5	eidskog	0.729	23.3	71.5	6.06	Innlandet	1
6	aure	0.683	28.0	67.0	3.39	Møre og Romsdal	1
7	austrheim	0.658	16.3	64.8	2.92	Vestland	1
8	aukra	0.644	14.9	63.6	3.76	Møre og Romsdal	1
9	vaksdal	0.617	19.2	61.4	3.88	Vestland	1
10	sokndal	0.544	20.9	55.5	3.37	Rogaland	1

```
# Output the top 10 recommended stores as a nice table using kable
# And save it
kable(head(Active_stores, 10), format = "markdown")
```

Mun_name	prob	Dist_nearest	Sales	Population	Region_Name	Active
lunner	0.9932445	7.323172	149.94498	9.420	Akershus	1
hareid	0.8573922	9.828661	87.21263	5.320	Møre og Romsdal	1
valer (ostfold)	0.8226703	10.788931	82.12523	6.162	Østfold	1
aurland	0.8166871	36.608030	81.33082	1.836	Vestland	1
eidskog	0.7294726	23.342412	71.47769	6.059	Innlandet	1

Mun_name	prob	Dist_nearest	Sales	Population	Region_Name	Active
aure	0.6827221	28.029241	67.04993	3.394	Møre og Romsdal	1
austrheim	0.6575973	16.315532	64.81866	2.915	Vestland	1
aukra	0.6436248	14.936345	63.61235	3.759	Møre og Romsdal	1
vaksdal	0.6172827	19.244243	61.39293	3.875	Vestland	1
sokndal	0.5442250	20.880302	55.49285	3.371	Rogaland	1