

DIY RISK MANAGEMENT DATA EXPLORATION

Bryce Chamberlain, Sr. Manager, Oliver Wyman Actuarial Consulting

A business of Marsh McLennan



LEARNING OBJECTIVES

At the end of this session, you will:

1

Be familiar with tools commonly used for data exploration.

2

Understand features and differentiators to determine when each tool should be used.

3

Know how to streamline your data exploration and get more value from data.

DATA WHISPERER, MUSIC ENTHUSIAST

- Associate of the Society of **Actuaries**
- Master of Science in **Analytics**, University of Chicago
- **12 years** professional experience
- Lead team at **Oliver Wyman Actuarial Consulting** building business intelligence apps for the web using **R Shiny**
- Passions: data **visualization**, user-friendly **design**, **efficiency** and **flexibility**



TOOLS OF THE TRADE



© Risk and Insurance Management Society, Inc. All rights reserved. Confidential and Proprietary. Do not disclose without written permission from RIMS General Counsel.

© Oliver Wyman

TOOLS OF THE TRADE

Worksheets

Calculator
Excel, Google Sheets

Most people can use it.

Easily manipulate single records.

Easy to make mistakes.

Difficult to automate.
Slow on large data.

Business Intelligence

Get started now
Power BI, Tableau, Salesforce

Lots of options quickly.

Click & drag

Limited functionality.
Difficult to automate.

AutoML

Search for insights
Rapidminer, DataRobot, SageMaker

Find stories across all data.

Limited visualization.
Results are complex.
Expensive if not open source.

Design

Make it pretty
Adobe Illustrator, Inkscape

Lots of features and options for perfecting the visual.

Very time consuming.
Software is complex, difficult to learn.

Free-hand Drawing

Begin with the end in mind
Pen & Paper, Tablet

Fastest method, no interface to slow you down.

Not generated by data.
Not fit for delivery.

Code-based

Automate repetitive tasks
RMarkdown, R Shiny

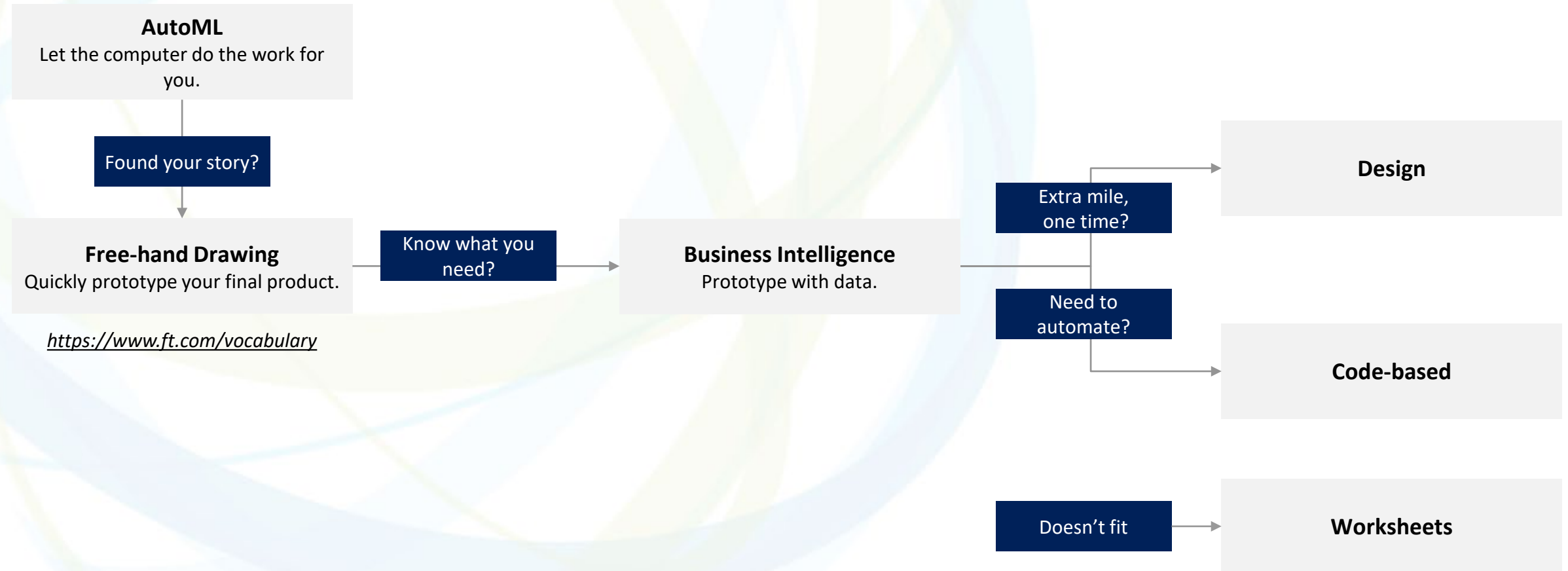
Unlimited functionality.

Open source.

Git version control.

Time consuming.
Need to code.

SECRET RECIPE



CASE STUDY

Our Data

- Auto claims dataset
- <https://www.kaggle.com/mykeysid10/insurance-claims-fraud-detection>
- 1,000 claims, 39 columns



GLIMPSE YOUR DATA

<https://bit.ly/rims-dict>
easyr::dict()

| | Name | Class | # Distinct | Top Values | Bottom Values | # Missing | % Missing | % =0 | # <0 | % <0 | Mode | Average |
|----|-----------------------------|-----------|------------|--|--|-----------|-----------|--------|------|--------|-------------|---------|
| 10 | insured_sex | character | 2 | FEMALE, MALE | FEMALE, MALE | 0 | 0 | NA | NA | NA | NA | |
| 37 | fraud_reported | character | 2 | N, Y | N, Y | 0 | 0 | NA | NA | NA | NA | |
| 4 | policy_state | character | 3 | IL, IN, OH | IL, IN, OH | 0 | 0 | NA | NA | NA | NA | |
| 6 | policy_deductable | integer | 3 | 500, 1000, 2000 | 500, 1000, 2000 | 0 | 0 | 0.0000 | 0 | 0.0000 | 500.00 | |
| 26 | property_damage | character | 3 | ?, NO, YES | ?, NO, YES | 0 | 0 | NA | NA | NA | NA | |
| 27 | bodily_injuries | integer | 3 | 0, 1, 2 | 0, 1, 2 | 0 | 0 | 0.3393 | 0 | 0.0000 | 0.00 | |
| 29 | police_report_available | character | 3 | ?, NO, YES | ?, NO, YES | 0 | 0 | NA | NA | NA | NA | |
| 42 | incident_date_month | integer | 3 | 1, 2, 3 | 1, 2, 3 | 0 | 0 | NA | NA | NA | NA | |
| 18 | collision_type | character | 4 | ?, Front Collision, Rear Collision, Side Collision | ?, Front Collision, Rear Collision, Side Collision | 0 | 0 | NA | NA | NA | NA | |
| 19 | incident_severity | character | 4 | Major Damage, Minor Damage, Total Loss, Trivial... | Major Damage, Minor Damage, Total Loss, Trivial... | 0 | 0 | NA | NA | NA | NA | |
| 25 | number_of_vehicles_involved | integer | 4 | 1, 2, 3, 4 | 1, 2, 3, 4 | 0 | 0 | 0.0000 | 0 | 0.0000 | 1.00 | |
| 28 | witnesses | integer | 4 | 0, 1, 2, 3 | 0, 1, 2, 3 | 0 | 0 | 0.2492 | 0 | 0.0000 | 0.00 | |
| 20 | authorities_contacted | character | 5 | Ambulance, Fire, None, Other, Police | Ambulance, Fire, None, Other, Police | 0 | 0 | NA | NA | NA | NA | |
| 21 | incident_state | character | 7 | NC, NY, OH, PA, SC | OH, PA, SC, VA, WV | 0 | 0 | NA | NA | NA | NA | |
| 22 | incident_city | character | 7 | Arlington, Columbus, Hillsdale, Northbend, Nort... | Hillsdale, Northbend, Northbrook, Riverwood, Sp... | 0 | 0 | NA | NA | NA | NA | |
| 22 | incident_city | character | 7 | Arlington, Columbus, Hillsdale, Northbend, Nort... | Hillsdale, Northbend, Northbrook, Riverwood, Sp... | 0 | 0 | NA | NA | NA | NA | |
| 41 | policy_bind_date_dayofweek | integer | 7 | 1, 2, 3, 4, 5 | 3, 4, 5, 6, 7 | 0 | 0 | 0.0000 | 0 | 0.0000 | 1.00 | |
| 44 | incident_date_dayofweek | integer | 7 | 1, 2, 3, 4, 5 | 3, 4, 5, 6, 7 | 0 | 0 | 0.0000 | 0 | 0.0000 | 1.00 | |
| 8 | umbrella_limit | integer | 11 | -1000000, 0, 2000000, 3000000, 4000000 | 6000000, 7000000, 8000000, 9000000, 10000000 | 0 | 0 | 0.7948 | 1 | 0.0010 | -1000000.00 | 1111 |
| 39 | policy_bind_date_month | integer | 12 | 1, 2, 3, 4, 5 | 8, 9, 10, 11, 12 | 0 | 0 | 0.0000 | 0 | 0.0000 | 1.00 | |

AUTOML

<https://bit.ly/rims-storyteller>

```
dt %<>%  
  clean(run_autotype = FALSE) %>%  
  dropoutliers() %>%  
  convert_date_features() %>%  
  dropnoisecols() %>%  
  correlatedfeatures_find()
```

```
dt %>%  
  correlatedfeatures_address(  
    target = 'fraud_reported'  
  ) %>%  
  fitmodel() %>%  
  summary()
```

```
[1] "Checking for outliers."  
      dropped [10] rows with outlier at [policy_annual_premium] > 1852.745  
      dropped [4] rows with outlier at [capital-gains] > 91516  
      dropped [4] rows with outlier at [capital-loss] > 90918  
      dropped [8] rows with outlier at [total_claim_amount] > 104064  
      dropped [3] rows with outlier at [injury_claim] > 18695  
      dropped [7] rows with outlier at [property_claim] > 20250.3  
      dropped [9] rows with outlier at [vehicle_claim] > 71508  
[1] "Adding features."  
[1] "Removing columns that are not useful."  
      Dropped column [incident_date_year] for reason [singleval]: 2015.  
      Searching for high correlation among 946 feature combinations. May take some time.  
      3% of feature combinations were correlated.
```

```
Column [age] dropped and [months_as_customer] kept.  
Column [incident_state] dropped and [insured_hobbies] kept.  
Column [incident_type] dropped and [collision_type] kept.  
Column [collision_type] dropped and [incident_severity] kept.  
Column [authorities_contacted] dropped and [incident_severity] kept.  
Column [property_damage] dropped and [incident_severity] kept.  
Column [injury_claim] dropped and [total_claim_amount] kept.  
Column [property_claim] dropped and [total_claim_amount] kept.  
Column [total_claim_amount] dropped and [vehicle_claim] kept.  
Column [auto_make] dropped and [auto_model] kept.
```





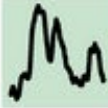



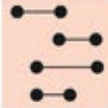











AUTOML

```
dt %>%
  correlatedfeatures_address(
    target = 'fraud_reported'
  ) %>%
  fitmodel() %>%
  summary()
```

| | Estimate | Std. Error | t value | Pr(> t) | |
|-------------------------------------|-----------------|----------------|---------|----------------------|-----|
| (Intercept) | 0.164960445036 | 0.040171694548 | 4.106 | 0.0000437 | *** |
| umbrella_limit | 0.000000010115 | 0.000000004778 | 2.117 | 0.03451 | * |
| policy_bind_date_month | -0.006863287511 | 0.003146919976 | -2.181 | 0.02943 | * |
| policy_bind_date_dayofweek | -0.008459130645 | 0.005460581538 | -1.549 | 0.12169 | |
| incident_date_dayofmonth | -0.001376083338 | 0.001274074058 | -1.080 | 0.28039 | |
| insured_education_level.JD | 0.062367733637 | 0.029905924386 | 2.085 | 0.03730 | * |
| insured_occupation.exec.managerial | 0.117831435862 | 0.042060325484 | 2.801 | 0.00519 | ** |
| insured_hobbies.exercise | -0.086613205931 | 0.047815216374 | -1.811 | 0.07040 | . |
| insured_hobbies.camping | -0.154213077138 | 0.048208930769 | -3.199 | 0.00143 | ** |
| insured_hobbies.chess | 0.608808153878 | 0.051961628124 | 11.716 | < 0.0000000000000002 | *** |
| insured_hobbies.sleeping | -0.155645240007 | 0.056503966915 | -2.755 | 0.00599 | ** |
| insured_relationship.other.relative | 0.051752761650 | 0.029651950080 | 1.745 | 0.08125 | . |
| insured_relationship.husband | -0.045365410736 | 0.029791805908 | -1.523 | 0.12816 | |
| incident_severity.Major.Damage | 0.495195916228 | 0.024774335205 | 19.988 | < 0.0000000000000002 | *** |
| auto_model.Wrangler | -0.094068769716 | 0.054932547630 | -1.712 | 0.08715 | . |
| auto_model.Civic | 0.143865911001 | 0.074785073348 | 1.924 | 0.05469 | . |
| auto_model.X6 | 0.178297750872 | 0.085898700399 | 2.076 | 0.03820 | * |

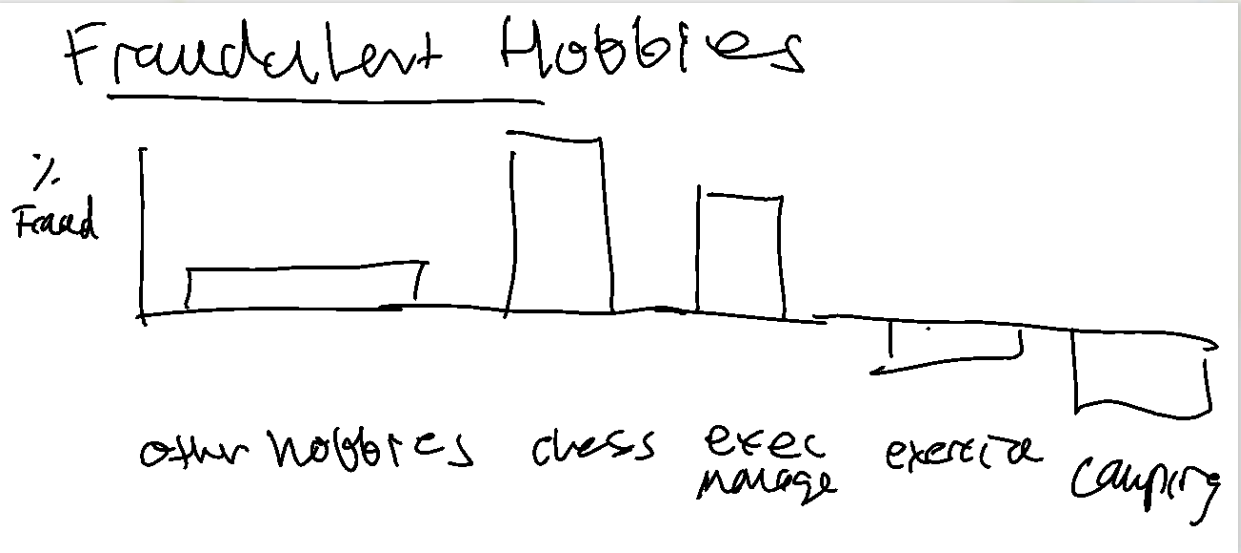
FREEHAND DRAWING

<https://www.ft.com/vocabulary>

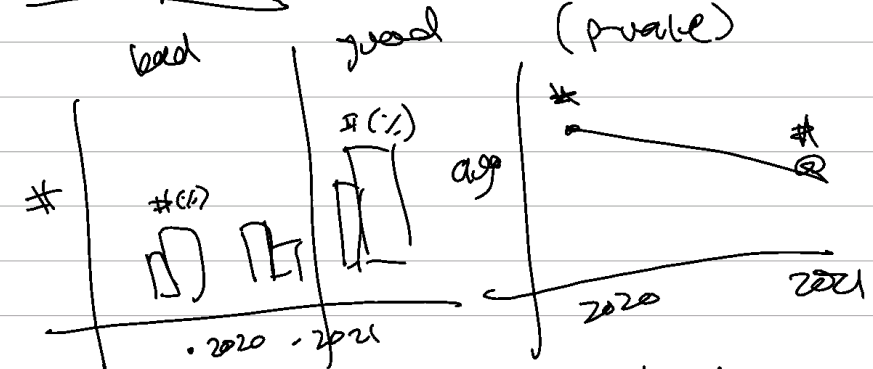
| Deviation | Correlation | Ranking | Distribution | Change over time |
|---|--|--|---|--|
| <p>Emphasise variations (+/-) from a fixed reference point. Typically the reference point is zero but it can also be a target or a long-term average. Can also be used to show sentiment (positive/neutral/negative).</p> <p>Example FT uses Trade surplus/deficit, climate change</p> | <p>Show the relationship between two or more variables. Be mindful that, unless you tell them otherwise, many readers will assume the relationships you show them to be causal (i.e. one causes the other).</p> <p>Example FT uses Inflation and unemployment, income and life expectancy</p> | <p>Use where an item's position in an ordered list is more important than its absolute or relative value. Don't be afraid to highlight the points of interest.</p> <p>Example FT uses Wealth, deprivation, league tables, constituency election results</p> | <p>Show values in a dataset and how often they occur. The shape (or 'skew') of a distribution can be a memorable way of highlighting the lack of uniformity or equality in the data.</p> <p>Example FT uses Income distribution, population (age/sex) distribution, revealing inequality</p> | <p>Give emphasis to changes. These can be short (in movements) or extended (traversing decades or centuries). Choosing the correct format is important to provide a clear picture for the reader.</p> <p>Example FT uses Share price movement series, sectoral change</p> |
| <p>Diverging bar</p>  <p>A simple standard bar chart that can handle both negative and positive magnitude values.</p> | <p>Scatterplot</p>  <p>The standard way to show the relationship between two continuous variables, each of which has its own axis.</p> | <p>Ordered bar</p>  <p>Standard bar charts display the ranks of values much more easily when sorted into order.</p> | <p>Histogram</p>  <p>The standard way to show a statistical distribution - keep the gaps between columns small to highlight the 'shape' of the data.</p> | <p>Line</p>  <p>The standard way to show a series of data points over time.</p> |
| <p>Diverging stacked bar</p>  <p>Perfect for presenting survey results which involve sentiment (eg disagree/neutral/agree).</p> | <p>Column + line timeline</p>  <p>A good way of showing the relationship between an amount (columns) and a rate (line).</p> | <p>Ordered column</p>  <p>See above.</p> | <p>Dot plot</p>  <p>A simple way of showing the change or range (min/max) of data across multiple categories.</p> | <p>Column</p>  <p>Column chart showing time with data.</p> |
| <p>Spine</p>  <p>Splits a single value into two contrasting components (eg male/female).</p> | <p>Connected scatterplot</p>  <p>Usually used to show how the relationship between 2 variables has changed over time.</p> | <p>Ordered proportional symbol</p>  <p>Use when there are big variations between values and/or seeing fine differences between data is not so important.</p> | <p>Dot strip plot</p>  <p>Good for showing individual values in a distribution, can be a problem when too many dots have the same value.</p> | <p>Column + line timeline</p>  <p>A good way of showing the relationship between a column and a line over time.</p> |
| <p>Surplus/deficit filled line</p>  <p>The shaded area of these charts allows a balance to be shown -</p> | <p>Bubble</p>  <p>Like a scatterplot, but adds additional detail by sizing the circles.</p> | <p>Dot strip plot</p>  <p>Dots placed in order on a strip are a space-efficient way of displaying data.</p> | <p>Barcode plot</p>  <p>Like dot strip plots, good for displaying all the data in a table.</p> | <p>Slope</p>  <p>Good character as a line.</p> |



FREEHAND DRAWING



Each question



- Notable
- level
 - office
 - tenure

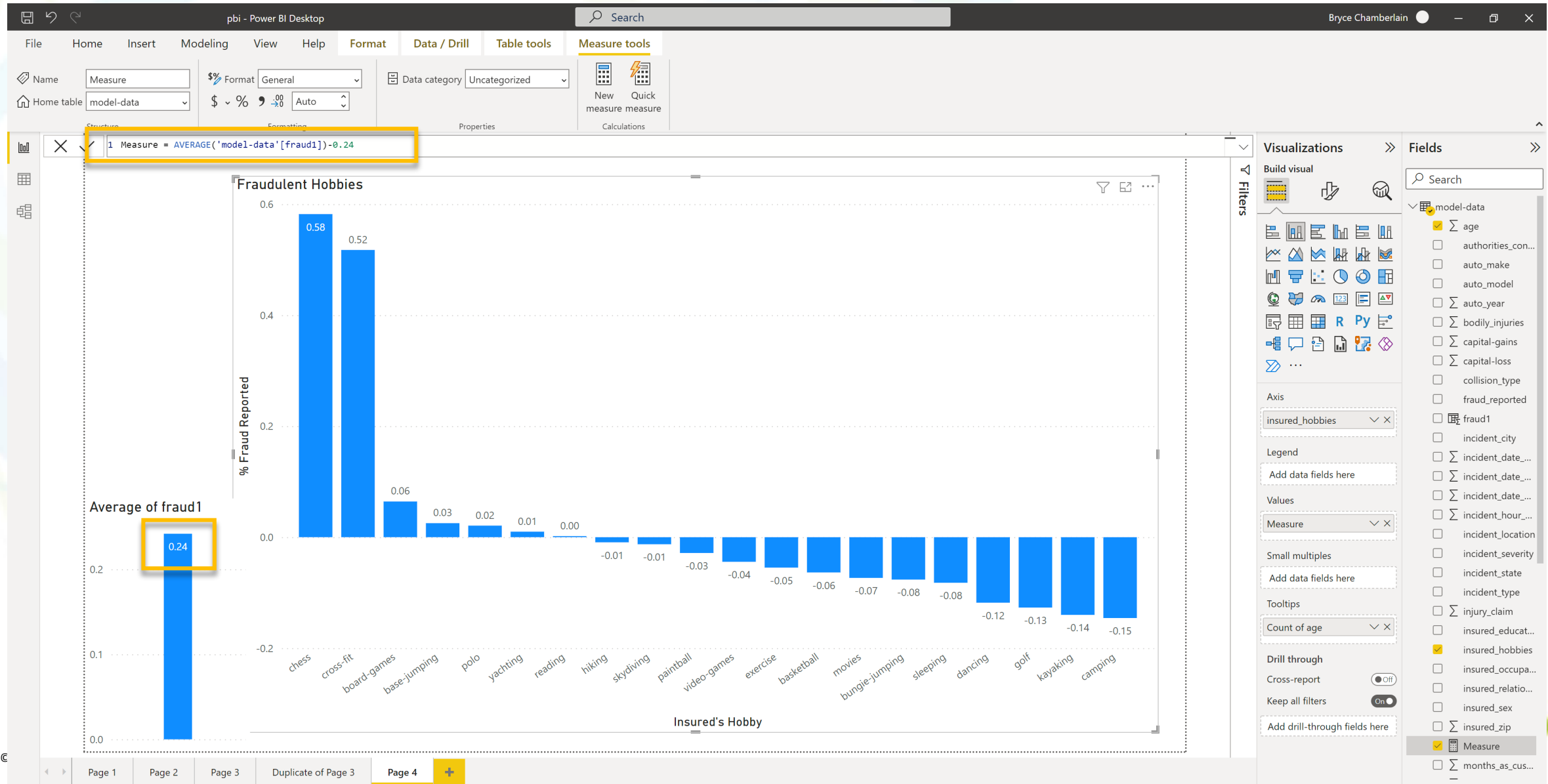
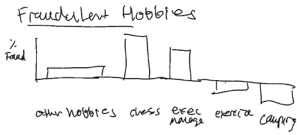
- Catgs of Qs
- higher better?
 - New 2021

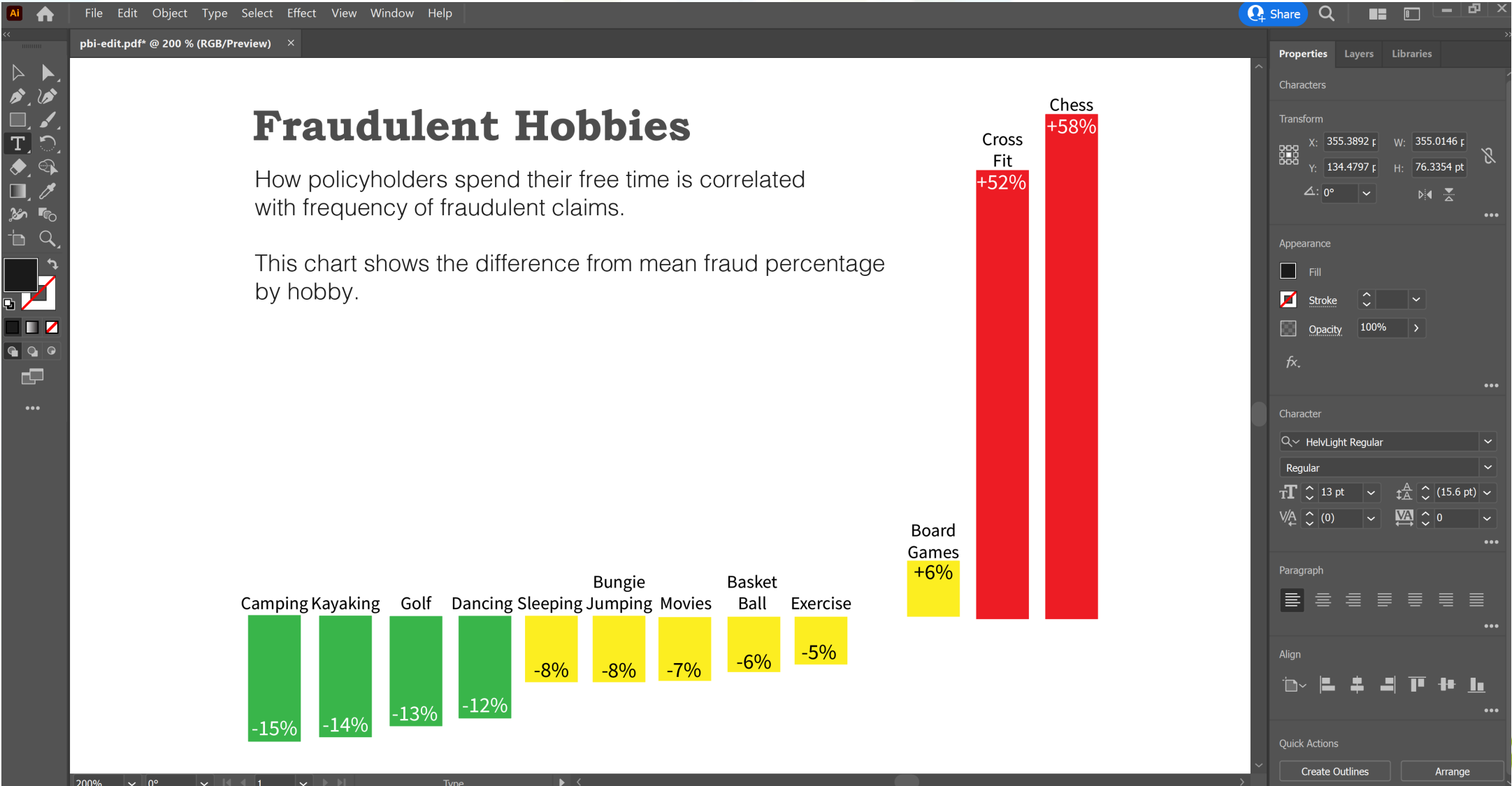
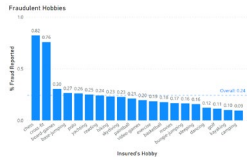


- Overall
- # questions
 - # responses

BUSINESS INTELLIGENCE (BI TOOLS)

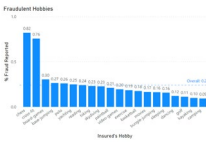
Fraudulent Hobbies





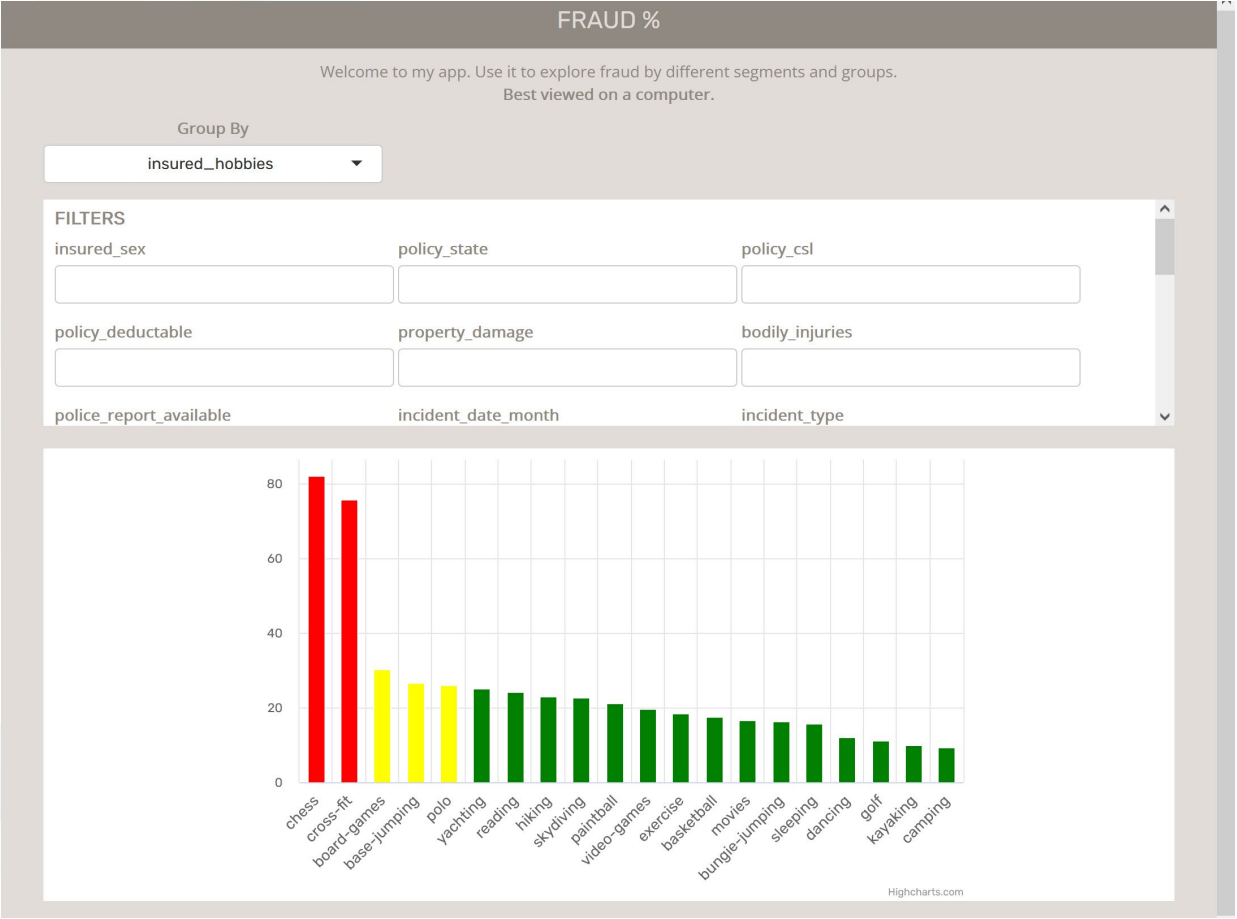
COD-BASED

<https://bit.ly/rims-codebased>



```
File Edit Selection View Go Run Terminal Help mainchart.R - shiny-forrims - Visual Studio Code
EXPLORER
mainchart.R X
mainchart.R app\ser...
SHINY-FORRIMS
app
  global
  rconnect
  server
    bookmark.R
    filters.R
    hcslim.R
    highcharter_replace.R
    inputs.R
    mainchart-data.R
    mainchart.R
    progress.R
    readme.txt
    utils.R
  ui
  www
  external
    fonts
      1-highcharts.js
      2-exporting.js
      fonts.css
      general.css
      setwidth.js
      shiny-override.css
      specific.css
    global.R
  OUTLINE
  TIMELINE
master 0 0 0
Ln 30, Col 41 Spaces: 2 UTF-8 CRLF R
```

```
mainchart.R
1 # https://github.com/superchordate/hcslim
2 output$mainchart = renderUI({
3
4   proginit('')
5
6   idt = chartdt()
7
8   proginc('Build Chart')
9
10  # select x and y.
11  idt$x = idt[[input$groupby]]
12
13  meanfraud = mean(idt$fraud_reported == 'Y') * 100
14
15  idt %>%
16    group_by(x) %>%
17    summarize(y = round(sum(fraud_reported == 'Y')/length(fraud_reported) * 10
18    arrange(y) %>%
19    mutate(x = factor(x, levels = x))
20
21  iqr = quantile(idt$y, c(0.5, 0.25, 0.75))
22
23  idt$color = fifelse(
24    abs(idt$y) <= 5,
25    'yellow',
26    fifelse(idt$y > 5, 'green', 'red')
27  )
28
29  options = list(
30    chart = list(animation = FALSE, marginLeft = 50),
31    plotOption = list(series = list(animation = FALSE)),
32    title = list(text = '', style = list(fontSize = '16pt')),
33    xAxis = list(categories = as.character(idt$x), labels = list(style = lis
34    yAxis = list(gridLineWidth = 1, endOnTick = FALSE, title = list(enabled
35    series = list(list(
```



SECRET RECIPE REFRESHER



REMINDERS, Q&A



Bryce.Chamberlain@oliverwyman.com
[linkedin.com/in/brycechamberlain](https://www.linkedin.com/in/brycechamberlain)



Hand me your card
for a follow-up email
including links



Complete Evaluation Rating
in the RIMS Mobile App



Your feedback is very
important to determine
the success and help
make improvements

Enjoy your next session!