# Diplomado de Modelado Predictivo y Machine Learning

# Introduction

## Hamdi Raissi
## José Ruette

# Data Science context

# Data Science Initiatives

# Data Science Initiatives



A Map of Data Science Degree Programs Around the The World

# Data Science Initiatives

# Data Science Initiatives

# Search Metrics

# Search Metrics

# Search Metrics

## Defining data science

I really like the definition quoted above, of data science as *the intersection of software engineering and statistics*. Ofer Mendelevitch goes into more detail, drawing a continuum of professions that ranges from software engineer on the left to pure statistician (or machine learning researcher) on the right.

| Software engineer | — | Data engineer | — | **Data scientist** | — | Data analyst | — | Statistician |

This continuum contains two additional roles, which are often confused with data scientists:

- *Data engineer:* a software engineer that deals with data plumbing (traditional database setup, Hadoop, Spark and all the rest)
- *Data analyst:* a person who digs into data to surface insights, but lacks the skills to do so at scale (e.g., they know how to use Excel, Tableau and SQL but can't build a web app from scratch)

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Befuddment



THE WALL STREET JOURNAL. ☰ | THE CIO REPORT

## CIO Journal
Exclusive reporting and analysis for corporate-technology execut

| CIO REPORT | CONSUMERIZATION | BIG DATA | CLOUD |

11:46 am ET
May 2, 2014    BIG DATA

### Why Do We Need Data Science When We've Had Statistics for Centuries?

ARTICLE    COMMENTS (13)

DATA SCIENCE    DATA SCIENTIST

✉ Email    🖶 Print    f 59    🐦 271

By IRVING WLADAWSKY-BERGER

**Data Science** is emerging as one of the hottest new professions and academic disciplines in these early years of the 21st century. A number of articles **have noted** that the demand for data scientists is racing ahead of supply. People with the necessary skills are scarce, primarily because the discipline is so new. But, the situation is rapidly changing, as universities around the world have started to offer different kinds of graduate programs in data science. This year, for example, New York University is offering two new degrees–a general **Master in Data Science**, and a more domain-specific **Master in Applied Urban Science and Informatics**.

It's very exciting to contemplate the emergence of a **major new discipline**. It reminds me of the advent of **computer science** in the 1960s and 1970s. Like data science, computer science had its roots in a number of related areas, including math, engineering and management. In its early years, the field attracted people from a variety of other disciplines who started out using computers in their work or studies, and eventually switched to computer science from their original field.

## IMS Bulletin online

⋮ HOME    ⋮ LATEST ISSUE PDF    ⋮ ARCHIVE (UNDER CONSTRUCTION)    ⋮ ABOUT    ⋮ ADVERTIS

HADLEY WICKHAM

Sep 4, 2014    👤 Editor    💬 18 Comments

### Data science: how is it different to statistics?

Contributing Editor Hadley Wickham is Chief Scientist at RStudio and Adjunct Professor of Statistics at Rice University. He is interested in building better tools for data science. His work includes R packages for data analysis (ggplot2, plyr, reshape2); packages that make R less frustrating (lubridate for dates, stringr for strings, httr for accessing web APIs); and that make it easier to do good software development in R (roxygen2, testthat, devtools, lineprof, staticdocs). He is also a writer, educator, and frequent contributor to conferences promoting more accessible and more effective data analysis. He writes:

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Befuddment

# Befuddment

## 50 years of Data Science

David Donoho

Sept. 18, 2015
Version 1.00

### Abstract

More than 50 years ago, John Tukey called for a reformation of academic statistics. In 'The Future of Data Analysis', he pointed to the existence of an as-yet unrecognized *science*, whose subject of interest was learning from data, or 'data analysis'. Ten to twenty years ago, John Chambers, Bill Cleveland and Leo Breiman independently once again urged academic statistics to expand its boundaries beyond the classical domain of theoretical statistics; Chambers called for more emphasis on data preparation and presentation rather than statistical modeling; and Breiman called for emphasis on prediction rather than inference. Cleveland even suggested the catchy name "Data Science" for his envisioned field.

A recent and growing phenomenon is the emergence of "Data Science" programs at major universities, including UC Berkeley, NYU, MIT, and most recently the Univ. of Michigan, which on September 8, 2015 announced a $100M "Data Science Initiative" that will hire 35 new faculty. Teaching in these new programs has significant overlap in curricular subject matter with traditional statistics courses; in general, though, the new initiatives steer away from close involvement with academic statistics departments.

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Tukey's Paper

## THE FUTURE OF DATA ANALYSIS[1]

### By John W. Tukey

*Princeton University and Bell Telephone Laboratories*

# Tukey's Paper

# Tukey's Paper

# Chambers Paper



John Chambers 1992

## Greater or Lesser Statistics: A Choice for Future Research

John M. Chambers
AT&T Bell Laboratories, Murray Hill, New Jersey

**Abstract**

The statistics profession faces a choice in its future research between continuing concentration on traditional topics, based largely on data analysis supported by mathematical statistics, and a broader viewpoint, based on an inclusive concept of learning from data. The latter course presents severe challenges as well as exciting opportunities. The former risks seeing statistics become increasingly marginal in activities to which it can make important contributions.

This paper is one of a set of short position papers on future research directions in statistics invited by the editor of *Statistics and Computation*.

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Cleveland Paper

Bill Cleveland 2002

## Data Science: An Action Plan for Expanding the Technical Areas of the Field of Statistics

William S. Cleveland
Statistics Research, Bell Labs
wsc@bell-labs.com

**Abstract**

An action plan to enlarge the technical areas of statistics focuses on the data analyst. The plan sets out six technical areas of work for a university department and advocates a specific allocation of resources devoted to research in each area and to courses in each area. The value of technical work is judged by the extent to which it benefits the data analyst, either directly or indirectly. The plan is also applicable to government research labs and corporate research organizations.

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Emergence of Meta Analysis

R

## An estimate of the science-wise false discovery rate and application to the top medical literature

LEAH R. JAGER

*Department of Mathematics, United States Naval Academy, Annapolis, MD 21402, USA*

JEFFREY T. LEEK\*

*Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD 21205, USA*

jleek@jhsph.edu

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Emergence of Meta Analysis

Meta Analysis
Systemic Failure
Causes
Solutions

## How Much of Our Published Research Can We Believe?

### Systemic Failures, Their Causes, a Solution

David Donoho

20140530

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Emergence of Reproducible Data Analyses

**Growth in Open Source Software Projects**



Market Realist

Source: Black Duck Management Webinar 2014

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Emergence of Reproducible Data Analyses

# Crisis in Machine Translation mid 1960's

# Crisis in Machine Translation mid 1960's

"We are safe in asserting that speech recognition is attractive to money. The attraction is perhaps similar to the attraction of schemes for turning water into gasoline, extracting gold from the sea, curing cancer, or going to the moon. One doesn't attract thoughtlessly given dollars by means of schemes for cutting the cost of soap by 10%. **To sell suckers, one uses deceit and offers glamor.**"

"It is clear that glamor and any deceit in the field of speech recognition blind the takers of funds as much as they blind the givers of funds. Thus, we may pity workers whom we cannot respect."

JR Pierce, 1969

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Crisis in Machine Translation mid 1960's

"Most recognizers behave, not like scientists, but like mad inventors or untrustworthy engineers. The typical recognizer gets it into his head that he can solve 'the problem.' The basis for this is either individual inspiration (the 'mad inventor' source of knowledge) or acceptance of untested rules, schemes, or information (the untrustworthy engineer approach). . . ."

"The typical recognizer ... builds or programs an elaborate system that either does very little or flops in an obscure way. A lot of money and time are spent. **No simple, clear, sure knowledge is gained.** The work has been an experience, not an experiment."

JR Pierce, 1969

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Common Task Framework

# Common Task Framework





Sebastiao Salgado, Work

# Predictive Modeling fundamentals

# Predictive modeling fundamentals



Predictive model

# Predictive modeling fundamentals

Good rules or formulas



Data source

Predictive model

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Predictive modeling fundamentals

Good rules or formulas

?

Data source

Predictive model

New data source

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Predictive modeling fundamentals

**Objectives**



Predictive model

Fundamentals

Business scenario data

Modeling challenges and solutions

# Predictive modeling fundamentals

## Applications

# Terms and elements of predictive modeling

Population

# Terms and elements of predictive modeling



Population

Inference

Sample

# Terms and elements of predictive modeling



Population

Common tools of inference

- Confidence intervals
- Hypothesis test
- P-values

Inference

Sample

# Terms and elements of predictive modeling

Population

Metrics

# Terms and elements of predictive modeling

Population

Metrics

Generalization

Y

Good rules or formulas

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Terms and elements of predictive modeling



New population

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Y

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | ? | ■ | ■ | $\cdots$ | ■ |
| 2 | ? | ■ | ■ | $\cdots$ | ■ |
| 3 | ? | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | ? | ■ | ■ | $\cdots$ | ■ |

Emphasis

• Empirical quality of predictions

• Understanding relationships

# Predictive modeling fundamentals



Input variables

Predictors
Explanatory variables
Inputs
Features

Target variable

Outcome
Response

Cases

Observation examples

$$y \quad X_1 \quad X_2 \quad \cdots \quad X_k$$

$$1 \quad 0 \quad \blacksquare \quad \blacksquare \quad \cdots \quad \blacksquare$$
$$2 \quad 1 \quad \blacksquare \quad \blacksquare \quad \cdots \quad \blacksquare$$
$$3 \quad 0 \quad \blacksquare \quad \blacksquare \quad \cdots \quad \blacksquare$$
$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \quad \vdots$$
$$n \quad 1 \quad \blacksquare \quad \blacksquare \quad \cdots \quad \blacksquare$$

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Build a model on historic data



|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Build a model on historic data

Supervised classification

The goal is to correctly classify

✔ Known discrete variable

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Build a model on historic data

Supervised classification

The goal is to correctly classify

Target variable Y
(binary response)



1 = Response
0 = No response

✓ Known discrete variable

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| ⋮ | ⋮ | ⋮ | ⋮ | | ⋮ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Build a model on historic data

Supervised classification

The goal is to correctly classify



Known discrete variable

Probability that a particular case belongs to one class

Target variable Y (binary response)

1 = Response
0 = No response

Score

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Build a model on historic data

Supervised classification

The goal is to correctly classify

Known discrete variable

Why use a predictive model?

Probability that a particular case belongs to one class

Score

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Y

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

New data arrives



Why use a predictive model?

Unknown discrete variable

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | ? | ■ | ■ | $\cdots$ | ■ |
| 2 | ? | ■ | ■ | $\cdots$ | ■ |
| 3 | ? | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | ? | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

New data arrives



Why use a predictive model?

$\times$ Unknown discrete variable

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | ? | ■ | ■ | $\cdots$ | ■ |
| 2 | ? | ■ | ■ | $\cdots$ | ■ |
| 3 | ? | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | ? | ■ | ■ | $\cdots$ | ■ |

Score

# Basic steps of predictive modeling

Input variables

$$\overbrace{\phantom{X_1 \quad X_2 \quad \cdots \quad X_k}}$$

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

## Basic steps

1, Supervised classification

• Prepare the inputs

• Select the most predictive inputs and fit models

2. Generalization

• Assess the models

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Basic steps of predictive modeling

Most predictive Input variables

$$y \quad X_1 \quad X_2 \quad \cdots \quad X_k$$

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

## Basic steps

1, Supervised classification

• Prepare the inputs

• Select the most predictive inputs and fit models

2. Generalization

• Assess the models

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Basic steps of predictive modeling

Most predictive Input variables



Basic steps

1, Supervised classification

• Prepare the inputs

• Select the most predictive inputs and fit models

2. Generalization

• Assess the models

# Basic steps of predictive modeling

Most predictive Input variables

Fit statistics



## Basic steps

1, Supervised classification

- Prepare the inputs

- Select the most predictive inputs and fit models

2. Generalization

- Assess the models

# Applications of predictive modeling

| | |
|---|---|
| Target marketing | Attrition predicción |
| Credit scoring | Fraud detection |

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Applications of predictive modeling

Target
marketing



Customers

Cases

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Applications of predictive modeling

Target marketing

Customers

Cases

Input variables

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Applications of predictive modeling

Target marketing

Target variable

Input variables

|  | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |  | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Customers

Cases

Target new potential customers

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Applications of predictive modeling

Target marketing

# Applications of predictive modeling

Target marketing

# Applications of predictive modeling



Attrition predicción

# Applications of predictive modeling

# Applications of predictive modeling



Credit scoring

Past applicants

Cases

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Applications of predictive modeling



Credit scoring

Input variables

Past applicants

Cases

$$
\begin{array}{c|c|cccc}
 & y & X_1 & X_2 & \cdots & X_k \\
\hline
1 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
2 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
3 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\vdots & \vdots & \vdots & \vdots & & \vdots \\
n & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\end{array}
$$

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Applications of predictive modeling

**Credit scoring**

Target variable

Input variables

Past applicants

Cases

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Reduce defaults and serious delinquencies

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Applications of predictive modeling

Fraud detection

Transactions

Cases

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

# Applications of predictive modeling

Fraud detection

Input variables

Transactions

Cases

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Applications of predictive modeling

**Fraud detection**

Target variable

Input variables

Transactions

Cases

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0   | ■     | ■     | $\cdots$ | ■     |
| 2 | 1   | ■     | ■     | $\cdots$ | ■     |
| 3 | 0   | ■     | ■     | $\cdots$ | ■     |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■     | ■     | $\cdots$ | ■     |

Anticipate fraud or abuse

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Demonstration Scenario: Target Marketing for a Bank

Marketing campaign

Bank

Customers

Y

Model

# Demonstration Scenario: Target Marketing for a Bank

```
%global inputs;
%let inputs=ACCTAGE DDA DDABAL DEP DEPAMT CASHBK
            CHECKS DIRDEP NSF NSFAMT PHONE TELLER
            SAV SAVBAL ATM ATMAMT POS POSAMT CD
            CDBAL IRA IRABAL LOC LOCBAL INV
            INVBAL ILS ILSBAL MM MMBAL MMCRED MTG
            MTGBAL CC CCBAL CCPURC SDB INCOME
            HMOWN LORES HMVAL AGE CRSCORE MOVED
            INAREA;

proc means data=work.develop n nmiss mean min max;
    var &inputs;
run;

proc freq data=work.develop;
    tables ins branch res;
run;
```

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Demonstration Scenario: Target Marketing for a Bank
## Questions

- What type of variable is Moved?

- How many variables have missing values?

- Look at the percentage of cases that have an Ins variable value of 1 versus those with a value of 0. What could you infer about the selection of cases for this data set?

- How many bank branches are represented in the data? Do you think this is a useful number of levels for the analysis?

- How many area classifications are represented in the data? Which area has the largest number of customers?

# Demonstration Scenario: Target Marketing for a Bank
## Questions

- What type of variable is Moved?

A: Moved is a binary variable.

- How many variables have missing values?

A: A total of 15 variables have missing values. Missing values are an issue. You learn how to handle missing values later in the course.

- Look at the percentage of cases that have an Ins variable value of 1 versus those with a value of 0. What could you infer about the selection of cases for this data set?

A: The results of PROC FREQ show that 34.6% of the customers in the develop data set purchased the insurance product. You might think that this percentage seems artificially high. In fact, the target event (buying the insurance product) is rare—only 2% of the population. To build the develop data set, the bank included all cases that have an Ins variable value of 1 and a representative sample of cases that have an Ins variable value of 0. This oversampling of the events increases the efficiency of the analysis because you are using a smaller sample and therefore have fewer cases to process. However, this oversampling also biases the results. You learn more about oversampling events, and how to adjust the model for it, later in the course.

- How many bank branches are represented in the data? Do you think this is a useful number of levels for the analysis?

A: The Branch of Bank table (the frequency table for Branch) indicates that the customers represented in the data do their banking in 19 different branches. When you determine that a categorical input variable has too many levels to be useful, you can collapse the levels. You learn to do this later in the course.

- How many area classifications are represented in the data? Which area has the largest number of customers?

A: The Area Classification table indicates that Res has three levels: R (rural), S (suburban), and U (urban). The largest number of customers live in urban areas, followed by suburban areas, and then rural areas.

# Predictive Modeling Challenges

# Predictive Modeling Challenges

Data challenges

Analytical challenges

Objectives

| Describe challenges that predictive modelers commonly encounter | Identify solutions to some of these challenges |
|---|---|
| Define honest assessment | Split the data |

# Data challenges

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Redundant variab

Operation

Opportunisti

$$
\begin{array}{c c c c c c}
 & y & X_1 & X_2 & \cdots & X_k \\
1 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
2 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
3 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\vdots & \vdots & \vdots & \vdots & & \vdots \\
n & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\end{array}
$$

Missing value

Millions of case

Hundreds of variabl

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

**Intervals**    0, 0,85 , 2000, 50.15 , 10000,....

**Discrete**    Male, Female

**Count** 0, 1, 2, 3,...

**Nomina**    D, B, C, E, A

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Account type

Dummy variable

| Value | Label | D1 | D2 |
|---|---|---|---|
| 1 | Checking | 1 | 0 |
| 2 | Savings | 0 | 1 |
| 3 | Other | 0 | 0 |

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Zip code

Collapse label

| Zip Code | City | State |
|---|---|---|
| 90620 | Buena Park | California |
| 90621 | Buena Park | California |
| 90622 | Buena Park | California |
| 90623 | La Palma | California |
| 90624 | Buena Park | California |
| 90630 | Cypress | California |
| 90631 | La Habra | California |
| 90632 | La Habra | California |
| 90633 | La Habra | California |
| 90680 | Stanton | California |
| 90720 | Los Alamitos | California |
| 90721 | Los Alamitos | California |
| 90740 | Seal Beach | California |
| 90742 | Sunset Beach | California |
| 90743 | Surfside | California |
| 92602 | Irvine | California |
| 92603 | Irvine | California |

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Dimensions

$$
\begin{array}{cc|cccc}
 & y & X_1 & X_2 & \cdots & X_k \\
1 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
2 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
3 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\vdots & \vdots & \vdots & \vdots & & \vdots \\
n & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare
\end{array}
$$

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |

Dimensions

$$\overbrace{\quad\quad\quad\quad\quad}$$

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ | $X_{k+1}$ | $X_{k+2}$ | $X_{k+3}$ | $\cdots$ | $X_{n+w}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |

Really sparse da

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Dimensions

Curse of dimensiona

$$y \quad X_1 \quad X_2 \quad \cdots \quad X_k \quad X_{k+1} \quad X_{k+2} \quad X_{k+3} \quad \cdots \quad X_{n+w}$$

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ | $X_{k+1}$ | $X_{k+2}$ | $X_{k+3}$ | $\cdots$ | $X_{n+w}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ | ■ | ■ | ■ | $\cdots$ | ■ |

$$\boldsymbol{n < n + w}$$

Hard to asses variable relation

Diplomado de Modelado Predictivo y Machine Learning

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

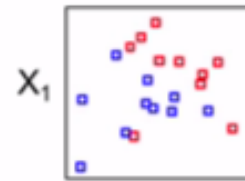| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

$$
\begin{array}{c|cccc}
 & y & X_1 & X_2 & X_3 \\
1 & 0 & \blacksquare & \blacksquare & \blacksquare \\
2 & 1 & \blacksquare & \blacksquare & \blacksquare \\
3 & 0 & \blacksquare & \blacksquare & \blacksquare \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
n & 1 & \blacksquare & \blacksquare & \blacksquare
\end{array}
$$

Reduce dimensiona

Take the most important vari

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

No - Even

Event

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

No - Event

Event

<1%

**Response**
Fraud
Churn
Default
Buy

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

No - Event

Event

<1%

**Response**
Fraud
Churn
Default
Buy

**No Response**
Legitimate
Stay
Pay
Not Buy

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Use all the data is time consuming

No - Even

Event

# Data challenges
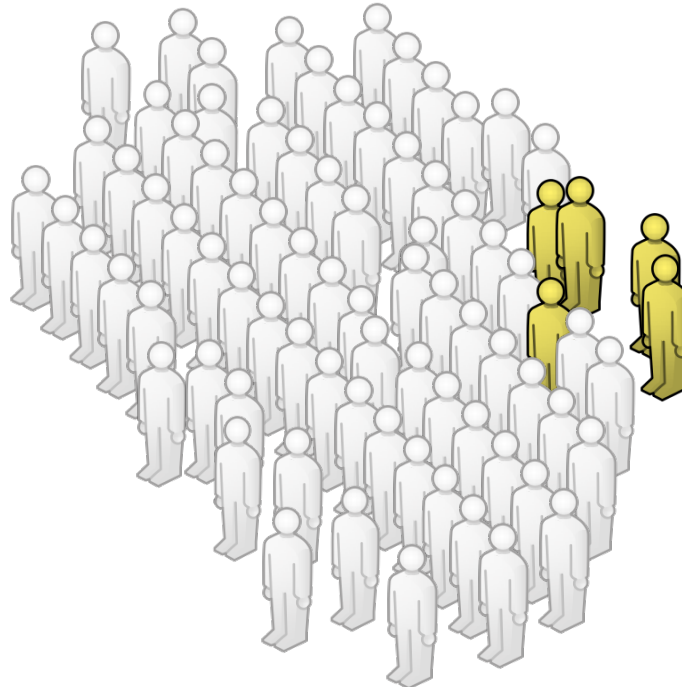
| Observational data | Mixed measurement scales | High dimensionality | Rare target events |

Sample

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
|---|---|---|---|

Representative Samp

# Data challenges



Observational data

Mixed measurement scales

High dimensionality

Rare target events

Representative Sample

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Data challenges

| Observational data | Mixed measurement scales | High dimensionality | Rare target events |
| --- | --- | --- | --- |

x the number of event case

Oversampling

$\approx Predictive\ power$

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Analytical challenges

# Analytical challenges

Non linearity and interactions

Model selection



E(y)

$x_1$

$x_2$

linear additive

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Analytical challenges

Non linearity and interactions

Model selection

Theory

Reality



linear additive

nonlinear nonadditive

# Analytical challenges

Non linearity and interactions

Model selection

# Analytical challenges

Non linearity and interactions

Model selection



Overfit

True

# Analytical challenges

# Analytical challenges



Non linearity and interactions

Model selection

Good fit, Generalizes well

True

# Separate sample

Data se

# Separate sample



Data se

Joint sample

Representative sample

**Data set with equal proportion of events to non-events**

# Separate sample

Data set with rare events



Joint sample ✕

Representative sample

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Separate sample



Data set with rare events

Target base sample

# Separate sample

Data set with rare ev[ents]

No event sample

Event sample

Target base sample

# Separate sample



Secondary outcome (Non-e



Primary outcome (eve

# Separate sample



Secondary outcome (Non-e

Some of the cas

Primary outcome (eve

All of the case

# Separate sample



Secondary outcome (Non-e

Primary outcome (eve

Some of the cas

All of the case

Modeling samp

# Separate sample

Modeling samp

Efficient

But generate bi

# Separate sample

a. a data set that consists of 100 events and 5,000 non-events

b. a data set that consists of 50 events and 10,000 non-events

c. a data set that consists of 5,000 events and 25,000 non-events

d. a data set that consists of 1,000 events and 5,000,000 non-events

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Separate sample

Which of the following data scenarios lends itself most to oversampling the target?

a. a data set that consists of 100 events and 5,000 non-events

b. a data set that consists of 50 events and 10,000 non-events

c. a data set that consists of 5,000 events and 25,000 non-events

d. **a data set that consists of 1,000 events and 5,000,000 non-events**

If you have millions of cases but only a thousand events, analyzing all of the non-events is inefficient. In scenario *d*, the ratio of non-events to events is 5000 to 1, which is larger than the ratios for the other scenarios.

Diplomado de Modelado Predictivo y Machine Learning

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Avoiding the Optimism Bias: Honest assessment



Fit

$$y \quad X_1 \quad X_2 \quad \cdots \quad X_k$$

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Y

?

Assess the model with the same

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Avoiding the Optimism Bias: Honest assessment

**10 case data set with two**

Model assessed on development



$x_1$

← gray / black →

$x_2$

# Avoiding the Optimism Bias: Honest assessment



Model assessed on development

Accuracy

10 case data set with two

Generalize well to new data?

# Avoiding the Optimism Bias: Honest assessment



Model assessed on 100 new cases

Accuracy

Generalize well to new data?

# Avoiding the Optimism Bias: Honest assessment



Model assessed on 100 new cases

Accuracy

Accuracy

This indicate that you overfit or underfit the

$X_1$

$X_2$

← gray   black →

# Avoiding the Optimism Bias: Honest assessment



Model assessed on 100 new cases

Fit

Accuracy

Accuracy

This indicate that you overfit or underfit the

Assess

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Avoiding the Optimism Bias: Honest assessment



Model assessed on 100 new cases

Fit

Accuracy

Accuracy

This indicate that you overfit or underfit the

Not a good id

Assess

Optimism bi

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Avoiding the Optimism Bias: Honest assessment

# Avoiding the Optimism Bias: Honest assessment

# Splitting the Data for model training and assessment

$$
\begin{array}{c c c c c c}
 & y & X_1 & X_2 & \cdots & X_k \\
1 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
2 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
3 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
4 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
5 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
6 & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
7 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
8 & 0 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\vdots & \vdots & \vdots & \vdots & & \vdots \\
m & 1 & \blacksquare & \blacksquare & \cdots & \blacksquare \\
\end{array}
$$

**Honest assessme**

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Splitting the Data for model training and assessment



**Honest assessme**

# Splitting the Data for model training and assessment



Training data s     Fit the mode

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Validation data

**Honest assessme**    Holdout portio

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $w$ | 1 | ■ | ■ | $\cdots$ | ■ |

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Splitting the Data for model training and assessment

Training data s

Fit the mode

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |



Train the model

Validation data

**Honest assessme**

Holdout portio

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ |
| $w$ | 1 | ■ | ■ | $\cdots$ | ■ |



Assess and compare models

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Splitting the Data for model training and assessment

# Splitting the Data for model training and assessment

**Training data s**

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Percentage of co

Rule of thumb

## 2/3 or 66,66%

**Validation data**

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $w$ | 1 | ■ | ■ | $\cdots$ | ■ |

**Honest assessme**

## 1/3 or 33,33%

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Splitting the Data for model training and assessment

**Training data s**

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

**Percentage of co**

Rule of thumb

## 2/3 or 66,66%

**Random samplin**

**Validation data**

**Honest assessme**

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $w$ | 1 | ■ | ■ | $\cdots$ | ■ |

## 1/3 or 33,33%

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Splitting the Data for model training and assessment

**Training data s**

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0   | ■     | ■     | $\cdots$ | ■     |
| 2 | 1   | ■     | ■     | $\cdots$ | ■     |
| 3 | 0   | ■     | ■     | $\cdots$ | ■     |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

**Percentage of co**

Rule of thumb

2/3 or 66,66%

**Random amplin**

1/3 or 33,33%

**Validation data**

|   | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|-----|-------|-------|----------|-------|
| 1 | 0   | ■     | ■     | $\cdots$ | ■     |
| 2 | 1   | ■     | ■     | $\cdots$ | ■     |
| 3 | 0   | ■     | ■     | $\cdots$ | ■     |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $w$ | 1 | ■ | ■ | $\cdots$ | ■ |

**Honest assessme**

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Splitting the Data for model training and assessment

Stratified random samp

| | Training (66.67%) | Validation (33.33%) | |
|---|---|---|---|
| Event | 7,451 (35%) | 3,724 (35%) | 11,175 (35%) |
| Non-event | 14,061 (65%) | 7,028 (65%) | 21,089 (65%) |
| | 21,512 (100%) | 10,752 (100%) | 32,264 (100%) |

Strata

**Honest assessme**

| | $y$ | $X_1$ | $X_2$ | $\cdots$ | $X_k$ |
|---|---|---|---|---|---|
| 1 | 0 | ■ | ■ | $\cdots$ | ■ |
| 2 | 1 | ■ | ■ | $\cdots$ | ■ |
| 3 | 0 | ■ | ■ | $\cdots$ | ■ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | 1 | ■ | ■ | $\cdots$ | ■ |

Develop data s

32,264 (100%)

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Splitting the Data for model training and assessment

**a.** Data splitting can be used only on data with continuous targets.

**b.** The validation data set is used to calculate the parameter estimates and validate the model.

**c.** Assessing the performance of the model on the data that you used to fit the model usually leads to an optimistically biased assessment.

**d.** Small differences in performance on the training data set versus the validation data set usually indicate overfitting.

PONTIFICIA UNIVERSIDAD CATÓLICA DE VALPARAÍSO

# Splitting the Data for model training and assessment

**a.** Data splitting can be used only on data with continuous targets.

**b.** The validation data set is used to calculate the parameter estimates and validate the model.

**c. Assessing the performance of the model on the data that you used to fit the model usually leads to an optimistically biased assessment.**

**d.** Small differences in performance on the training data set versus the validation data set usually indicate overfitting.
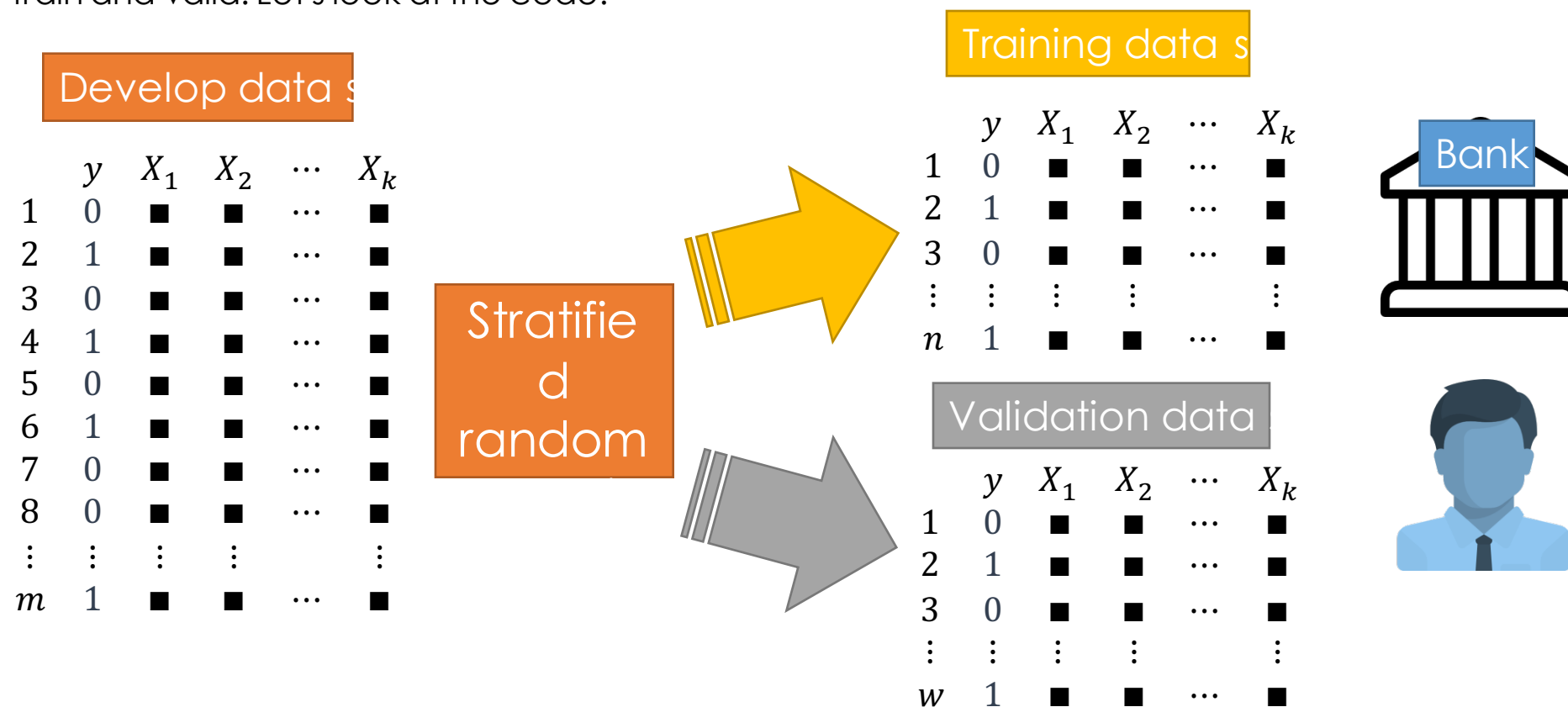
**Answer *a* is <u>incorrect</u>** because data splitting can be used on data with any type of target.

**Answer *b* is <u>incorrect</u>** because the validation data set is used to validate the model. The training data set is used to calculate the parameter estimates.

**Answer *d* is <u>incorrect</u>** because large differences in performance on the training data set versus the validation data set usually indicate overfitting.

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO

# Splitting the Data for model training and assessment

For the target marketing project at the bank, we want to split the develop data set into a training data set and a validation data set. In this demonstration, we do the following: Use a stratified sample to select the records for the training and validation data sets, and create two data sets: train and valid. Let's look at the code.

# Splitting the Data for model training and assessment

```sas
proc sort data=work.develop out=work.develop_sort;
    by ins;
run;

proc surveyselect noprint data=work.develop_sort
                  samprate=.6667 stratumseed=restore
                  out=work.develop_sample
                  seed=44444 outall;
    strata ins;
run;

proc freq data=work.develop_sample;
    tables ins*selected;
run;

data work.train(drop=selected SelectionProb SamplingWeight)
     work.valid(drop=selected SelectionProb SamplingWeight);
    set work.develop_sample;
    if selected then output work.train;
    else output work.valid;
run;
```

# Splitting the Data for model training and assessment

**a.** The results would be the same as in the demonstration.

**b.** The proportion of the SELECTED=1 cases (cases in the training data set) would be different from the corresponding results in the demonstration.

**c.** The proportion of the events in the training data set would probably be different from the proportion of events in the validation data set.

# Splitting the Data for model training and assessment

**a.** The results would be the same as in the demonstration.

**b.** The proportion of the SELECTED=1 cases (cases in the training data set) would be different from the corresponding results in the demonstration.

**c. The proportion of the events in the training data set would probably be different from the proportion of events in the validation data set.**

Unlike a stratified random sample, a simple random sample does not guarantee an equal percentage of events in the training and validation data sets. However, because the sampling rate is the same as in the demonstration (0.6667), the training data set (SELECTED=1) will contain 66.67 percent of the observations regardless of the sampling method.

PONTIFICIA
UNIVERSIDAD
CATÓLICA DE
VALPARAÍSO