# A Brain-Friendly Guide

# Head First
# Data Analysis

Predict
your raise
with linear
regression

A learner's guide to
big numbers, statistics,
and good decisions

Sell more toys by
optimizing your
business model

Experiment to
discover who your
customers *really* are

Overcome
your
cognitive
biases

Load important
statistical concepts
directly into your brain

Clean messy data
for efficient analysis

Michael Milton

# Table of Contents (Summary)

# Table of Contents (the real thing)

## Intro

**Your brain on data analysis.** Here *you* are trying to *learn* something, while here your *brain* is doing you a favor by making sure the learning doesn't *stick*. Your brain's thinking, "Better leave room for more important things, like which wild animals to avoid and whether naked snowboarding is a bad idea." So how *do* you trick your brain into thinking that your life depends on knowing data analysis?
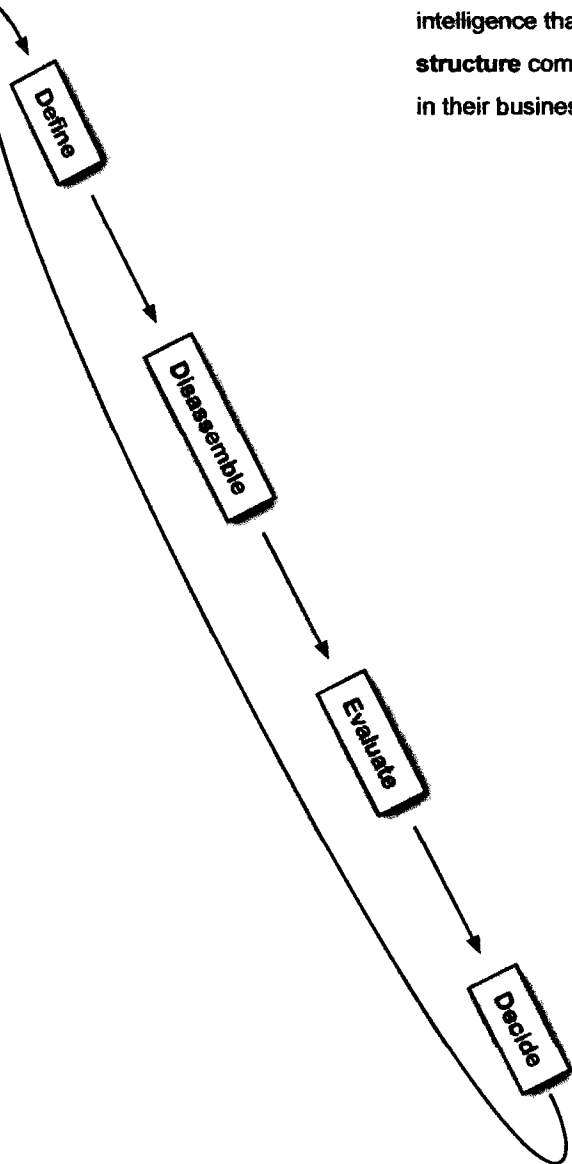
# introduction to data analysis

# Break it down

**1**

### Data is everywhere.

Nowadays, everyone has to deal with mounds of data, whether they call themselves "data analysts" or not. But people who possess a toolbox of data analysis skills have a **massive edge** on everyone else, because they understand what to *do* with all that stuff. They know how to translate raw numbers into intelligence that **drives real-world action**. They know how to **break down and structure** complex problems and data sets to get right to the heart of the problems in their business.
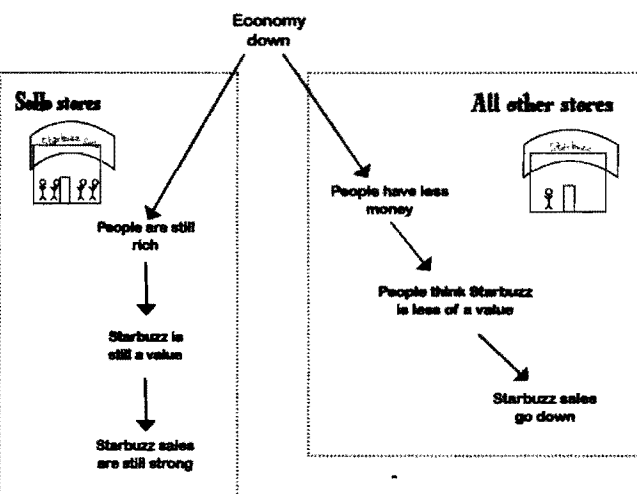
Define

Disassemble

Evaluate

Decide

x

# experiments

## Test your theories

### Can you show what you believe?

In a real **empirical** test? There's nothing like a good experiment to solve your problems and show you the way the world really works. Instead of having to rely exclusively on your **observational data**, a well-executed experiment can often help you make **causal connections**. Strong empirical data will make your analytical judgments all the more powerful.

**2**

Economy
down

Solio stores

All other stores

People are still rich

People have less money

People think Starbuzz is less of a value

Starbuzz is still a value

Starbuzz sales go down

Starbuzz sales are still strong
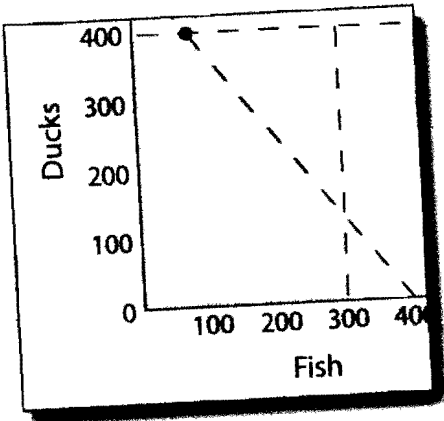
**3**

optimization

# Take it to the max

## We all want more of something.

And we're always trying to figure out how to get it. *If* the things we want more of—profit, money, efficiency, speed—can be represented numerically, then chances are, there's an tool of data analysis to help us tweak our *decision variables*, which will help us find the **solution** or *optimal point* where we get the most of what we want. In this chapter, you'll be using one of those tools and the powerful spreadsheet **Solver** package that implements it.
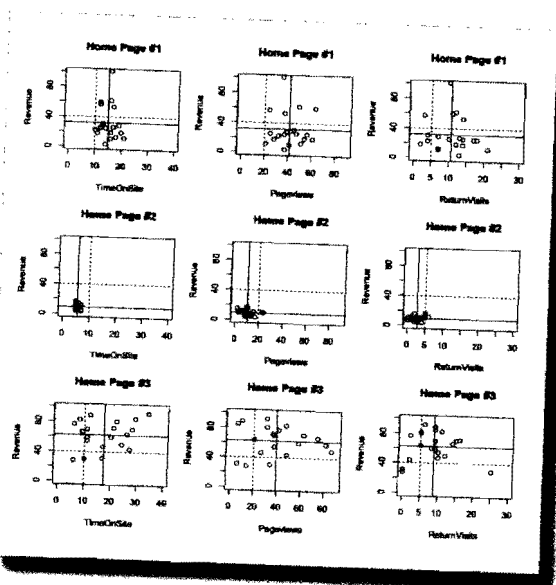
## data visualization

# Pictures make you smarter

## You need more than a table of numbers.

Your data is brilliantly complex, with more variables than you can shake a stick at. Mulling over mounds and mounds of spreadsheets isn't just boring; it can actually be a waste of your time. A clear, highly multivariate visualization can, in a small space, show you the forest that you'd miss for the trees if you were just looking at spreadsheets all the time.

# hypothesis testing
## Say it ain't so

### The world can be tricky to explain.

And it can be fiendishly difficult when you have to deal with complex, heterogeneous data to anticipate future events. This is why analysts don't just take the obvious explanations and assume them to be true: the careful reasoning of data analysis enables you to meticulously evaluate a bunch of options so that you can incorporate all the information you have into your models. You're about to learn about **falsification**, an unintuitive but powerful way to do just that.

# bayesian statistics

## 6

## Get past first base

### You'll always be collecting new data.

And you need to make sure that every analysis you do incorporates the data you have that's relevant to your problem. You've learned how *falsification* can be used to deal with heterogeneous data sources, but what about **straight up probabilities**? The answer involves an extremely handy analytic tool called **Bayes' rule**, which will help you incorporate your **base rates** to uncover not-so-obvious insights with ever-changing data.

*Cough*