# Practical Big Data Analytics

\* Michael.Burgess(@qa.com)

1.30 ~~4.30~~

9.30 — 4.30 $^{(4-4.30)}$
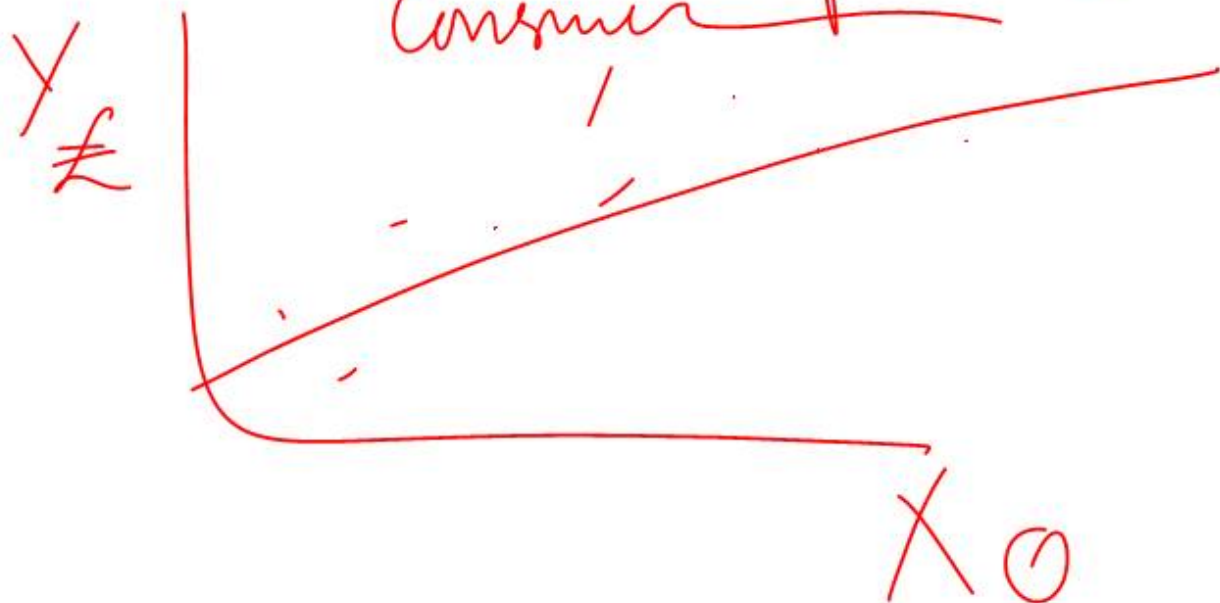
c.12 / 1pm

c. 11am

c. 2 30pm

c

The Skills To Suceed

Traing _ AP

~~Computers~~ Data↑

Consumer

$Y_{£}$

$X_O$

Computational
Statistical

Big Data

Interpreting
Data

Inf

AI

"Just"

NoSQL

Insight vs. Automate

Data Inf Science

Analytic

xs

Inference → Prediction

$V^3$ $V^5$

⊗ Vol.
Variety.
Velocity

$\dfrac{\text{amount}}{\text{structure}}$

⊗ skills &
tools

[flexibility] BD = non-trad method

(hw)

Analytica

Y-axis: 5, 1, 0

X-axis labels: Stats, Maths, algs, (art) Dom. k, Comm & Storytelly, Big Data, Trad. Data, Pry m

Curves labeled: Stats, ML, BA, Stp

# 1 pm

May need QAPYTH3, or eq.

Quick Review

- Wrap up python → fns / import (#2)

- Courseware Review

2→ tinyurl.com/bda-10-02-20

√[∨∧∨>$]

## In python.

$dir(x) \leftarrow$ lists all operations on $x$

Eg dir(dataframe)

$\vdots$

read_csv()

$\vdots$

to_csv()

$\vdots$

m

michael.burgess@
burgess @qq.com

Machine Learning $\longrightarrow$ Self D
$\longrightarrow$ AI

$\longrightarrow$ Maths

Rec.
Alexa
Data $\longrightarrow$ traing

alg/model

$X \longleftarrow$ know

$Y \longleftarrow$ ?

$$\hat{Y}(X) = 3X + 5 \text{ model}$$

hist $(X, Y)$

Movie Co

# Computational Stat

## Inference

AF ?

- 40s
- 50s

$OG_1$

Parameters
Reposting

```
if (@ > 18):
    open()
    clse
    CloseUx
```

6 6 6 6 → 6 6

in-sample → out-sample

£ UK ≠ € FR

France

anti

① in ~ out

② in ≠ out ← unsafe

$$\text{ML} \quad 1) \begin{cases} \text{Regression} \leftarrow \hat{y} \in \mathbb{R} \\ \text{Classification} \\ \text{sup learning} \end{cases} \begin{array}{l} \in \{-1, +1\} \quad \underline{\text{Binary Class}} \\ \in \{0, 1, \dots, 5, \dots\} \end{array}$$

$$\text{historical } (X, Y) \xrightarrow[\text{traing}]{\text{learning}} \hat{y}(X) \qquad \min \text{Loss}(\hat{Y}, Y)$$

$$2) \text{ unsup}$$

$$(X_1, X_2, X_3)$$
$$z \quad q_0 \quad t$$

① Simplify
(Dim. Reduction)

② Common Occurance

$$P(X, X, X)$$

# Regression

$$0.1x + 1$$

$$\frac{}{70}$$ — $(X, Y)$

$(0.1, 1)$

$loss(\hat{Y}, Y)$

pref $Y$

3

1

0

20R    63h    $X$

$4^2 = 16$    $-4^2 = 16$

# Classification



days $f$

$5$

$2.5$  $3$

$10$

| $X_1$ | $X_2$ | $X$ |
|---|---|---|
| $100$ | $10$ | ✓ |
| | | ✗ |
| | | ✓ |
| | | ✗ |
| | | ✗ |

Grade/
10

IM

Medicies

New W/ Support

Support

11.05

$X_2$

$X_1$ (units/wk)

1pm

$$(x^i, y^i)$$

$$\hat{y}(x_i, \omega) = \omega_1 x_1 + \omega_2 x_2 + \omega_3$$

var $\longleftrightarrow$ parameter
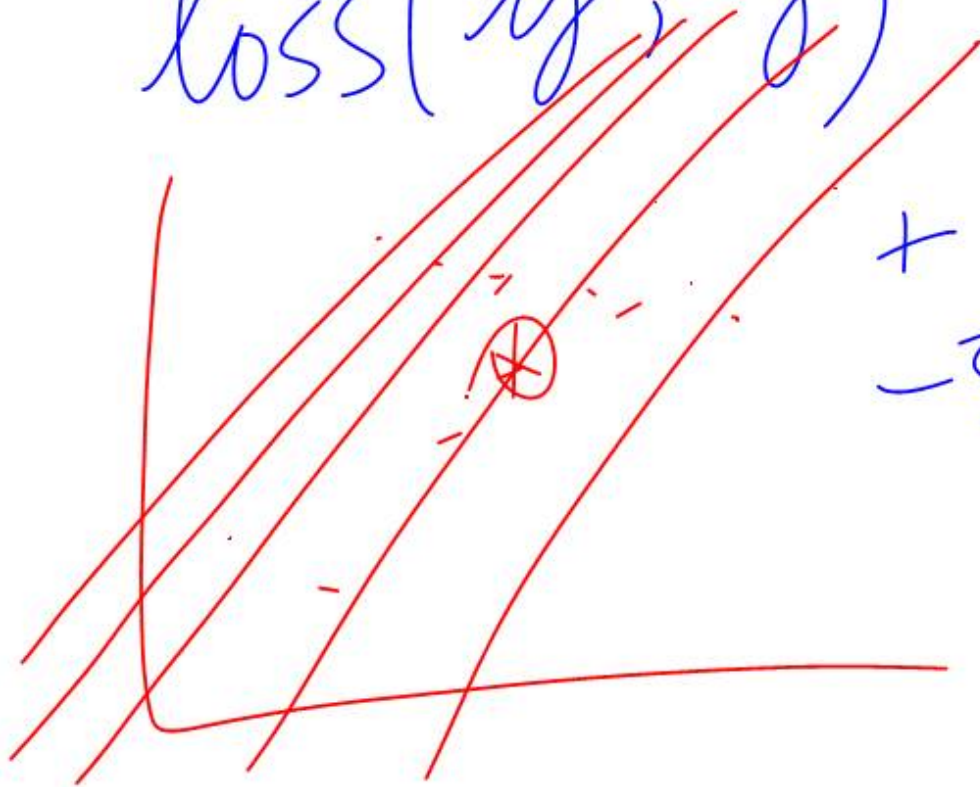
$\downarrow$ a single number

data model
Stat model — assoc
Expl

Causal

H

S

$$\text{loss}(\hat{y}, y) = (\hat{y} - y)^{-}$$

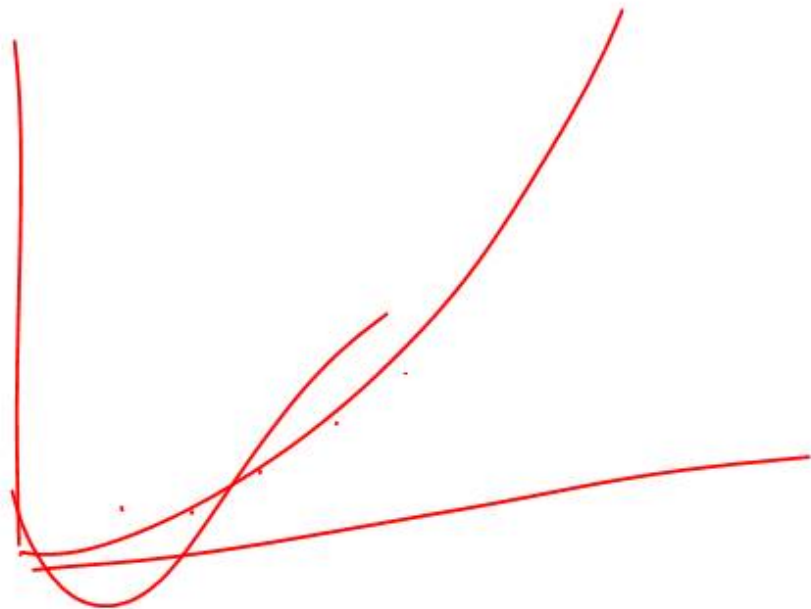$+5$

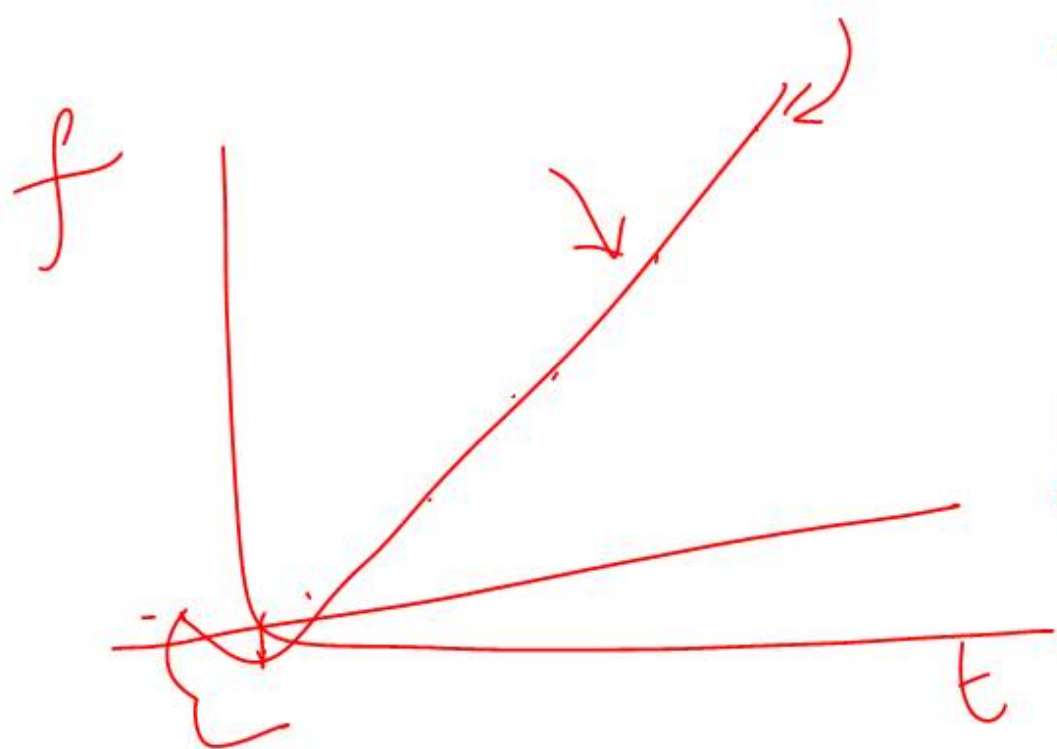$-5$

$$y \in \mathbb{R}$$

$$x_1, x_2 \ldots \in \mathbb{R}$$

$$\underline{find} \quad \hat{y}(x \ldots \mid w \ldots)$$

$$by \; \boxed{\min} \; loss\left(\hat{y} \mid y\right) \; Grad \; Des.$$

How? $\underline{Alg}$

$$\hat{y}(x;w) =$$

$$w_1 x_1 + w_2 x_2 + \cdots$$

$$= a(a(a(wx)))$$

$X$     $\swarrow$     $Y$     $Y$

$18$

$* w_1$   0.1    $a^{(\ )}$   $3.2$

$30$    0.2   $* w_2$

$y = u(u(a\cdots)))$

IK/G

(1430)

A

B

Aging

7  8  9  10  1

Simpson's
paradox

$$\frac{m - samp}{test + train} \quad \Rightarrow outsum$$

100%

m - samp

⇒ outsum

split

⇒70%

Perf.    100%

underfitty

overfitty

avg. loss
on spouts

validation

1    2    3    4    upprovens

# "ML workflow

$$\left( \overset{B_3}{\text{Exploratory}} \atop (x,y) \right)$$

⊗ Can we do you
anythg?

⊗ Problems? ~~Gatekeepr~~

Opions — ~~(HIPPO)~~ not sup.

order ← CEU
Chief D/S

① Understand. Prob ( BA )

② Data Understand
- loss
- Ex. D. A.

age
$x_1, x_2, \ldots$

1)



20  30   40

1000
③
996726

1.
N
N
:
N
y

N
N
N
a
r
)

# ③ Data Prep (why?) ky lm ToDie

Features (x)
Eng
- Clean ← remove missing
- ETL ← Joins (Enriching)
- Transformation ← Changing repr

City × Encoders

| | is LDN | is LOS | is MAN |
|---|---|---|---|
| LOS → 1 → 2 | 1 | 0 | C |
| MAN → 2 ② 1 | 0 | 1 | C |
| LDN → 3 | O | O | 1 |

OHE →

④ Modeling

NN-1
  -2
  -3

CR -1

Data $(X, y)$

test | train

Cross V

$\sqcap\sqcap\sqcap\sqcap \rightarrow 90\%$
A

$\sqcap\sqcap\sqcap\sqcap 9$
B

$\sqcap\sqcap\sqcap 8$
C

# Evaluation

test

Eval

B *

Retrain model on [ test + train data ] stafely

100%

Deployment

Repeat software

① Exp
② Problem
③ Data
④ Modely
⑤ Eval
⑥ Deployment

② 1 mo · T
= $\frac{1}{2}$ mo · T

1 mo · T
$\frac{1}{2}$ mo · T
$\frac{1}{2}$ mo · T
$\frac{1}{2}$ mo (1 mo / 6 mo)

6 mo —— years

# Big Data

100GB
1TB

2005 2015

(non-tr.) Skills

image

Variety ← "Dimensionality"
↳ vary indpt

Velocity  24h·s
10MB/s

Volume  10TB  100KB/s

[Json (K-V
F-C)]?

1000

6000

1

1000

/ 1PB

t

Non/traditional

BI

Relational — SQL

"One Machine"

Velocity

injestion

**BigData**

Analysts | Engineers

1) Data Structure ——→ CAP

2) Query Structure ——→ | Consistency
                         availability
3) "Representation"      partition tolerance
   Mostly Emory
   "Query"

   "Readshru"

Eventual
US
UK | US
US | UK

High Availability

Cloud
AWS DC  11.59

C
2s

A —— P

1 2 3 4 5 6 7
12

UK +20
F +30
BA
US -10
-20
F -20
JE -30
Blocky

ACID
AC | AC
A | B | C

Data 11.05

Data Module

Columnar ‖‖‖

MySQL —— Table (Relation)

*Neo4j —— Graph (Network)

Redis —— Key —— Value

*Mongo —— Documents

Unsr. /Schemaless "U.P.a" {"tag": }

{name
Loans. {1: {
2: {}}}

2pm

# Graphs & Graph Analytics

$$[ \; ( u, v, w ) \quad \rule{2cm}{0pt} \; ]$$

"M"  "A"  O

Graph — heirarchy

KEG
CTO

(Relationships) Communsin

Proless

Symmetric
friendships CEO→CIO

Enemy
directed HS ⇄ StiS
Enemy
DiS

PG

80%

70%

LA

0.8*0.1

10*0. 1

55%

nodes

*

100%

100%

RW

90%

10%

DM

Edges

Lot

Graph + DF.
+ BF.
Traversal.

purchase

weighted ✓

(Vertex = Node)

Queen
Bee
↓ neigh-

"certainty"
(importance)

(+ neighbors
+ height)

middle

out

terriot
freedom group



3

10

4   50%

2

2   2   3   3

5

DL

DL

team?

DL

Travel Map

$\big(\ Subgraph\ \big)$

— Min. Span Tree

: Shortest path

NX

LP

B

N

LP
X

→ Python

↑ NetworkX

Neo4j

DB

Pareto, Networks

pref. Attach.

$50\%$

$2.8\%$

$$\hat{y}(x; a,b) = ax + b^5 \qquad (x, y)$$

$$\ell(\hat{y}, y) = (\hat{y} - y)^2 = \ell(a, b; x, y)$$

$$\min_{a,b} \ell(a, b; x, y)$$

T·0/2/3

a := Rand()
b := Rand()

$a := a_{OLD} := \text{change}(\ell, a_{OLD \to NEW})$

$a := a_{NEW}$

$i = [30\ 0\ ..16\ \ .\ \ .\ \ 20]$

$\overset{m}{\phantom{.}} \qquad \overset{|X|}{\phantom{.}}$

$0\ 1\ 2\ 3\ 4$

API

$\boxed{\text{Sum}(images)}$

$$for\ (i=0;\ i < len;\ \boxed{i++}\ )\ \{$$

$[x+1\ for\ x\ in\ I]$

$total\ += images[i];$

$x$

RM

⊗ Streaming (RAM) ———— HDD
  — constant memory                    In

  [ ——————————⟶ ] 1000

  ↘ Streamy API

⊗ Parallynt ———— Manage for _____ dS

  CPU +              ds[i]

Demand

HDFS  <2003

MapReduce

Job Sched
Res.

A
B
C

Commodity → H

why?

A
B
C
C
A
B
C

$f(.jpg)$

$f(J, ps)$

1) Shuffle

map

$f(\downarrow)$ reducer

$lhev$

shuffle

$(inello': 2)$

$\delta$

Spark

hdfs

mr

Big Data

RDD
Resilient
Distr.
Dataset

Eng/SD *bython*

"[          ]"

RDD

(V1) FP f Resilient

(V2) Tabular + Over.

Dataset

SOL

2/5

$r$ . filter(test) . map(f) . collect()

↳ transformation

$a =$
$a =$
$- a =$

r —|— filter —|— m ——
 left        f

$a = (r . map(g) . ) filter(\rightarrow)$

~~filter~~ (~~test~~)

flatMap ( lifting )

~~fold~~ → Reduce
     → aggregate

map ← 1) Projection / "Structure"
    ← 2) Calc

) max()
collect ( )
~~top ( )~~

# Questions

→ **Next Steps**

⎰ Presentations / output ...
↳ psychology
↳ Evidence & argument

(11am) 15-20m Ex
5-10 Break

10.20

# Projects / Pipelines

**Problem /**
**Solution**

BDA

—Data —< Steaming / online
          Batch / DB / offline

— Injection

$f_y(x_1, x_2 \dots)$

Len

Online  Offline

St — Processing Ex. Prep. Analyst
     Preheat | Analysis / Model
              — Model
              — Eval.

This is a hand-drawn whiteboard diagram illustrating a machine learning pipeline for loan approval.

Top labels:
- **Loan Apr.** ① — eg. $x$: Age, $y$: yes/No ② — Sbng. ③
- SEng / DS (green box)
- Dept v.n (green box)

Pipeline stages:
- SEng → **Ingestion** (streaming) (3mo) — ① 
- **Processing** (1mo) — ②
- **Prediction** (1mo – 3mo+) SEng/DS → $y$ — ③

Lower section:
- online / offline
- devops
- BA
- **prep.** (3mo) $x$ → $x$
- **modelling** — ML → $\hat{y}(x)$? ($\frac{1}{2}$mo – 1mo)
- DB $x$ ← features
- **Eval** — BA ($\frac{1}{2}$)
- **tidy** — Sn Dep'l
- DB → $y$
- checking $\hat{y}(x)$ & History
- (6mo)

| Prototype | Injestion | Prep. & Process | Deploy & Predict | |
|---|---|---|---|---|
| • Build App Interface | • Build (Big)Data Pipelines | • Include Analytical Code | • Cloud & DevOps | **Online** |
| 2mo | 2mo | 1mo | • Include Prediction Code 3mo | |
| | x → Data Online | | ŷ → Online Data | |
| 1 | | | | |
| Software Eng. | Software Eng. 2 | S. Eng (& Analyst) | S. Eng, DevOps, Analyst 4 | **#2** |
| Explore & Prep. | Model | Evaluate | Package | **Offline** |
| Gather Data & Determine strategic & predictive importance | find Relationship ŷ (x ...) | Eval. ŷ(x) • against strategic metrics | Rephrase & improve code for use in App | |
| 3mo | ½ mo | ½ MO | ½ MO 4 | **#1 First** |
| BA, Analyst | ML & analyst | BA & analyst 3 | S. Eng & Analyst | |