
LESSON 16

DATABASES & SQL

Michael Twardos

What We Will Do Today

- Databases, data warehouse design and SQL introduction.
- Overview of product adoption and growth analytics.
- Compute retention metrics for a fictional product.
- Apply convolution to the retention curve to project future active users.
- Build a model to predict the retention likelihood of individual customers.
- Think about how a data science model can **actually** be used.

WHAT IS A DATABASE?

Why Databases?

- Databases are used as a repository of information. Allow for efficient storage and access.
- Types
 - Relational (MySQL, PostgreSQL, Redshift)
 - Key Value (Redis)
 - Document (Mongo)
 - Graph (Neo4j)
 - Time Series (Graphite)
 - Search (Elastic, Solr, Splunk)
 - Wide Column (Cassandra, HBase)

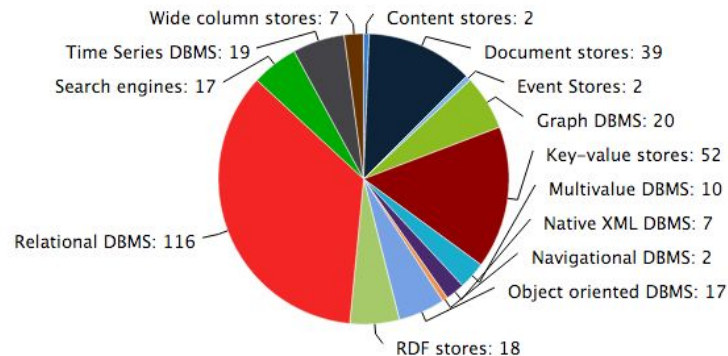
● DB Popularity

The top 5 commercial systems, March 2016

Rank	System	Score	Overall Rank
1.	Oracle	1472	1.
2.	Microsoft SQL Server	1136	3.
3.	DB2	188	6.
4.	Microsoft Access	135	7.
5.	SAP Adaptive Server	77	12.

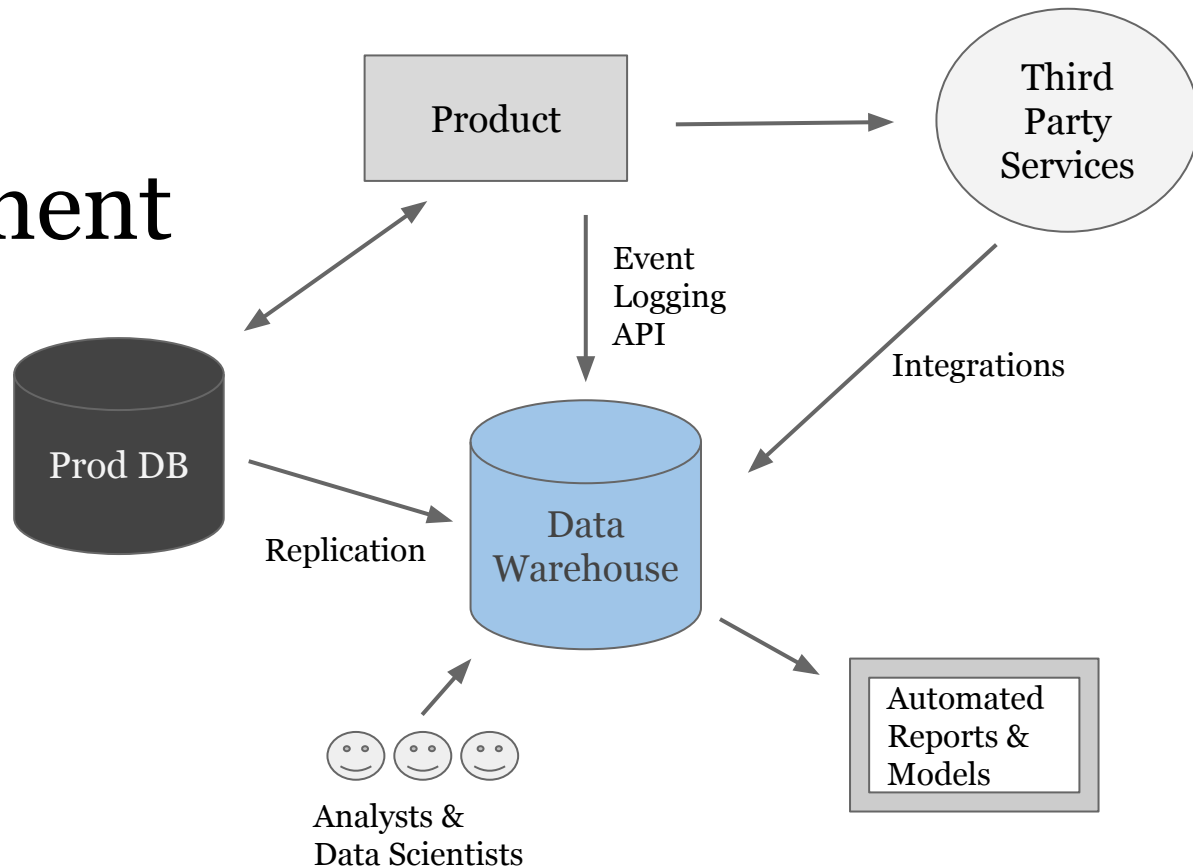
The top 5 open source systems, March 2016

Rank	System	Score	Overall Rank
1.	MySQL	1348	2.
2.	MongoDB	305	4.
3.	PostgreSQL	300	5.
4.	Cassandra	130	8.
5.	Redis	106	9.



© 2016, DB-Engines.com

Data Management



Master Fact Table

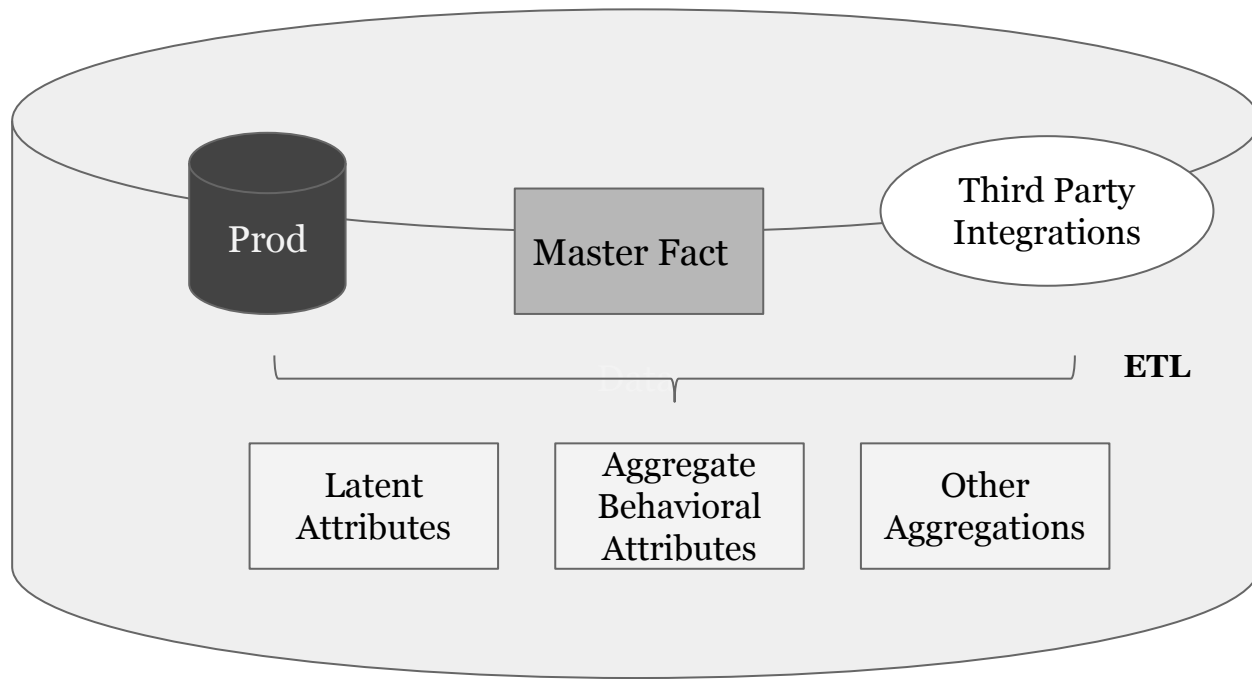
- ▶ A comprehensive historical record of everything that happened. (Think who, what, where, when, how, why)
- ▶ Required fields:

customer_id, timestamp, event_type

- ▶ Recommended fields

marketing campaign, location, device, browser, os...

Data Warehouse Design



Stories from Industry

What have you experienced and what was good or bad?

How was the data stored? How is it accessed?

SQL Read

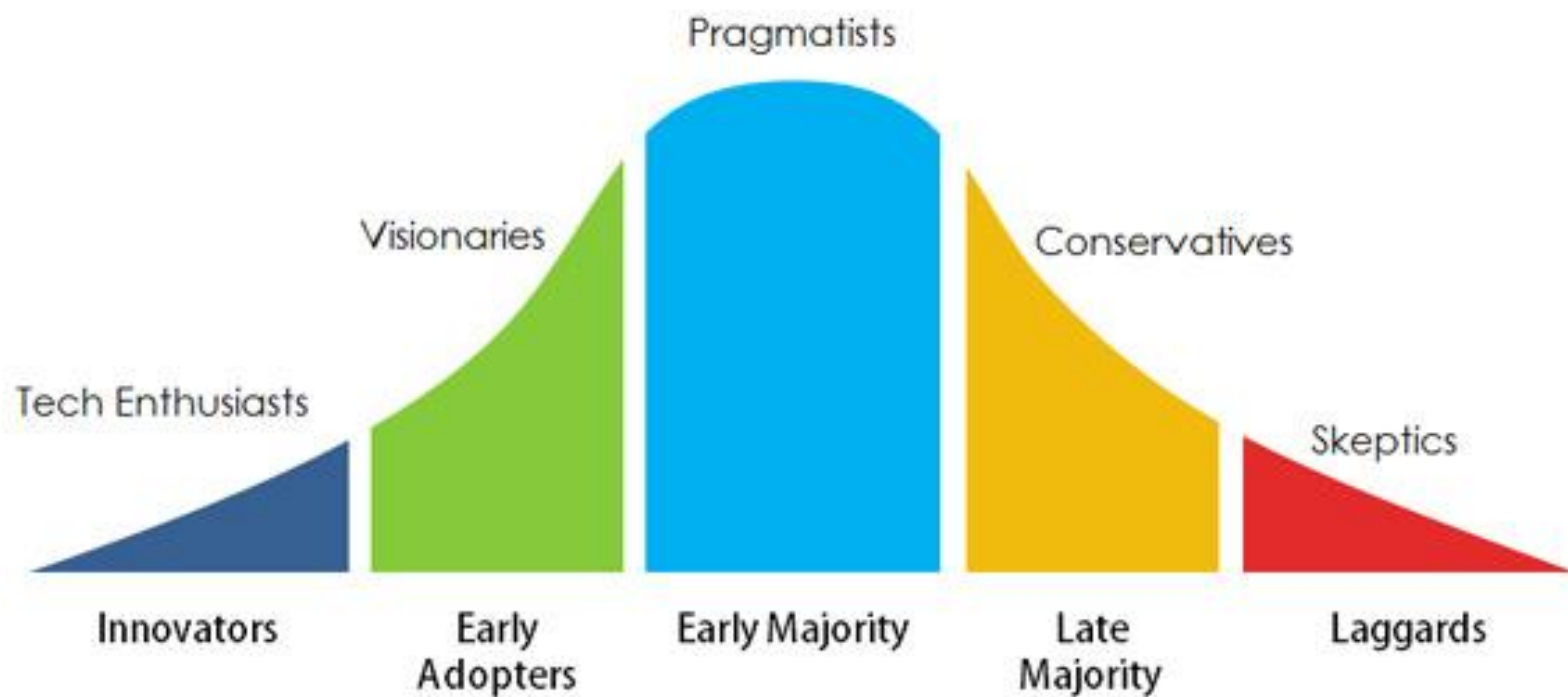
SELECT a
FROM b
WHERE c
GROUP BY d
ORDER BY e
LIMIT f
;

Python SQL Workshop

- ▶ Import data into SQLite
- ▶ Run basic SQL queries
- ▶ Explore SQL functions

PRODUCT ADOPTION & GROWTH ANALYTICS

Product Adoption Curve



PERCENT OF
U.S. HOUSEHOLDS

CONSUMPTION SPREADS FASTER TODAY

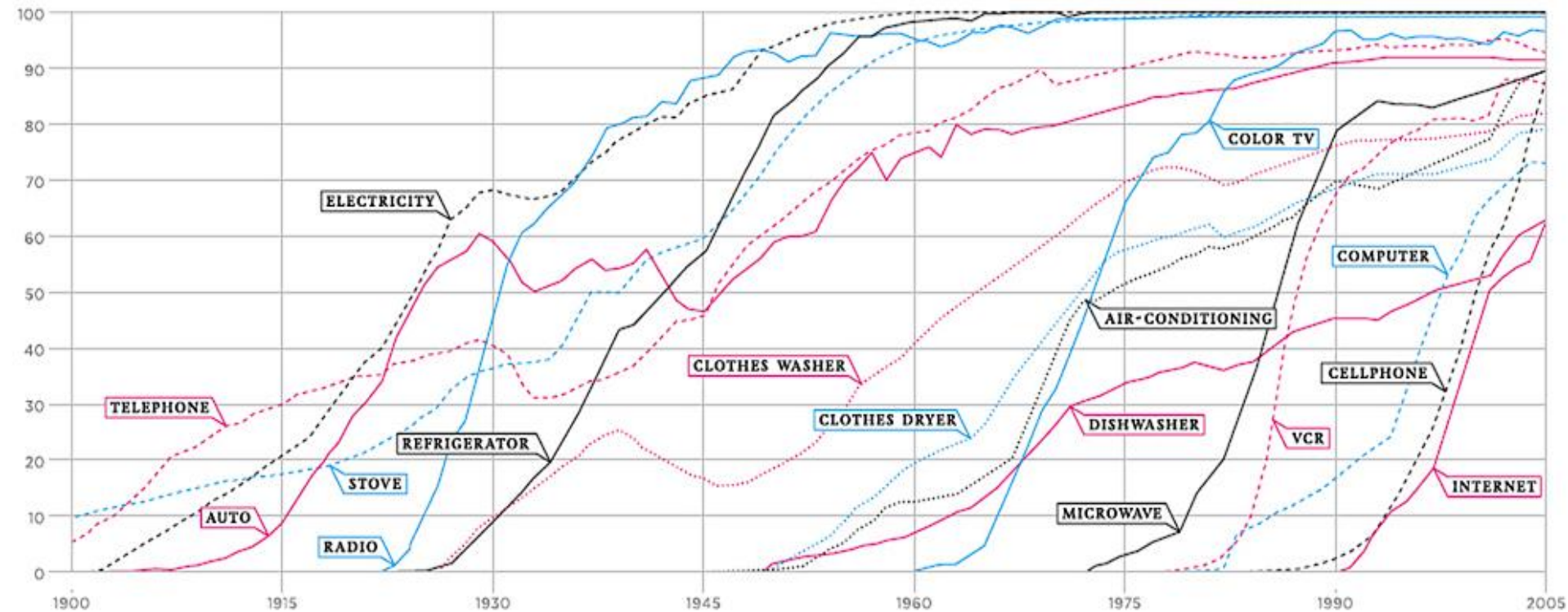
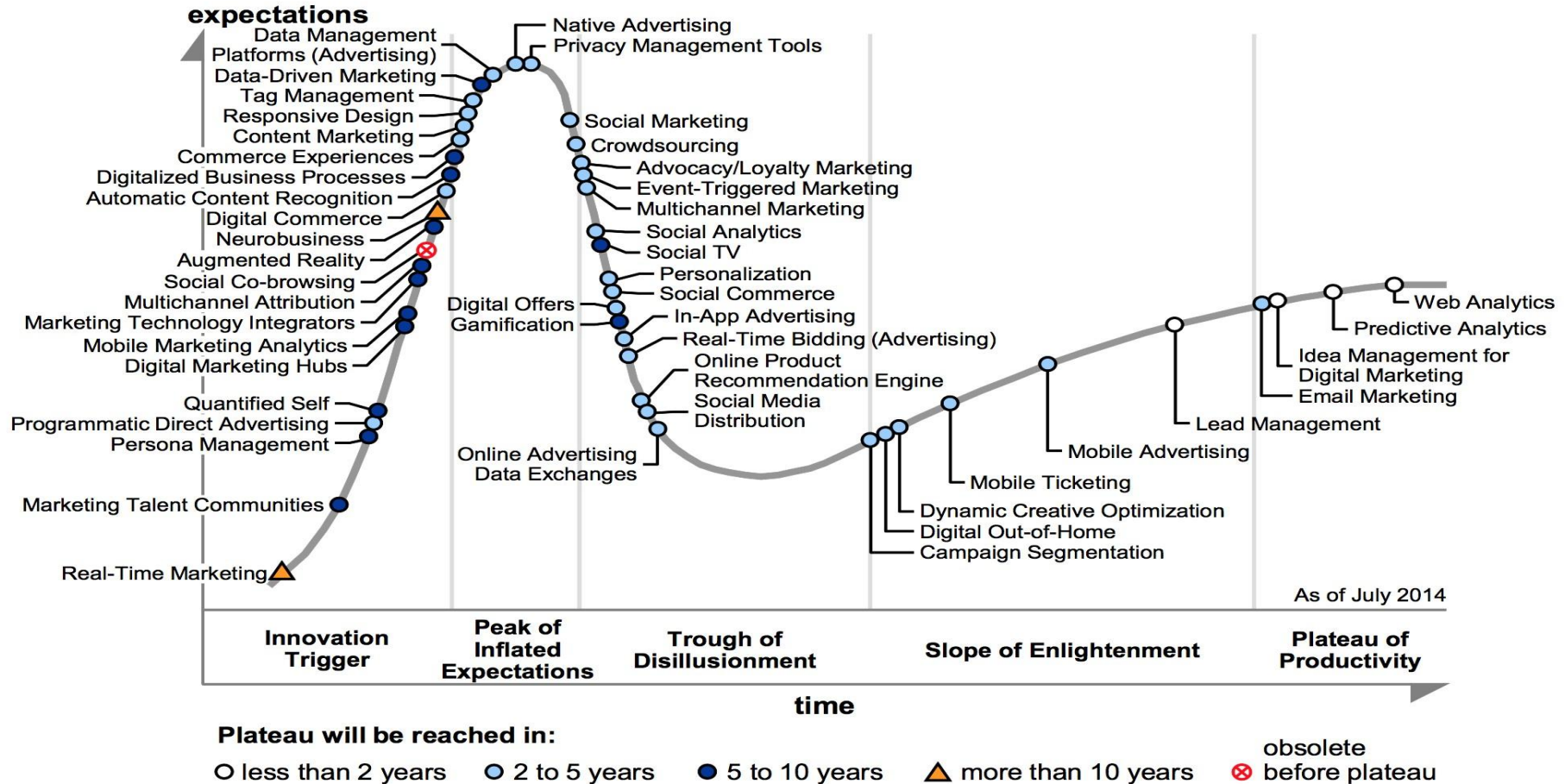


Figure 1. Hype Cycle for Digital Marketing, 2014

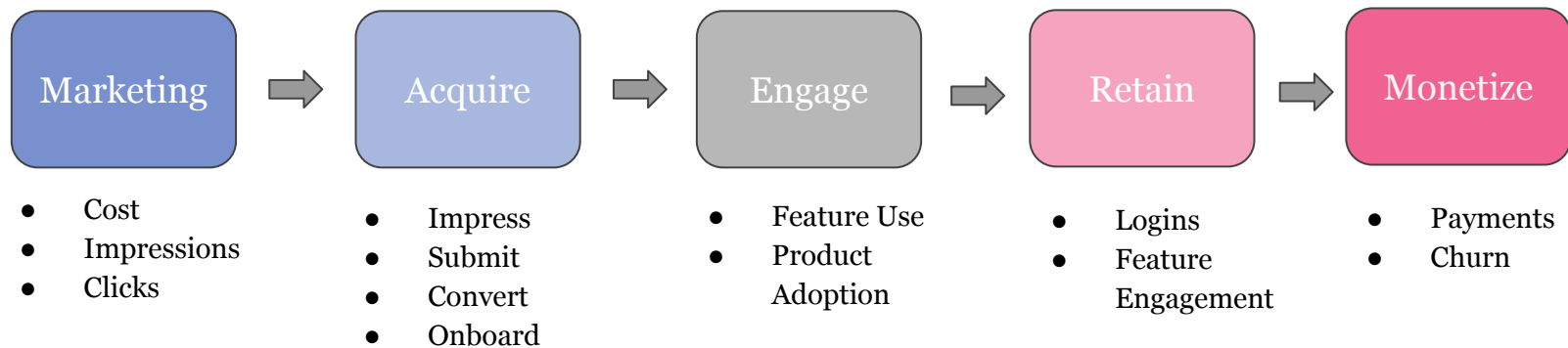


Stories of Product Adoption

Task: Think of a product / service / company you started using/interacting with in the last year and answer these questions.

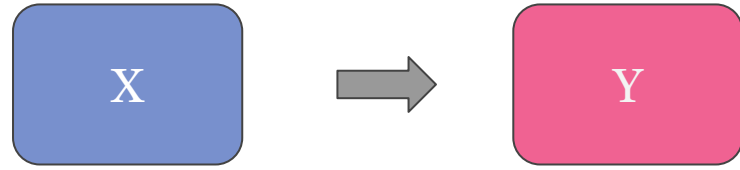
- 1) How did you hear about it?
- 2) What was your first experiences with it? Was it memorable?
- 3) How easy was the product to use? How valuable? How enjoyable?
- 4) How often do you use this product months after you first started to if at all?
- 5) Did you pay? Would you pay again?

A Customer's Lifecycle



- Events are captured through a customer's life.
- The period of time between events can define a user 'state'. A state is the atomic unit of measurement for growth.

Customer Transitions

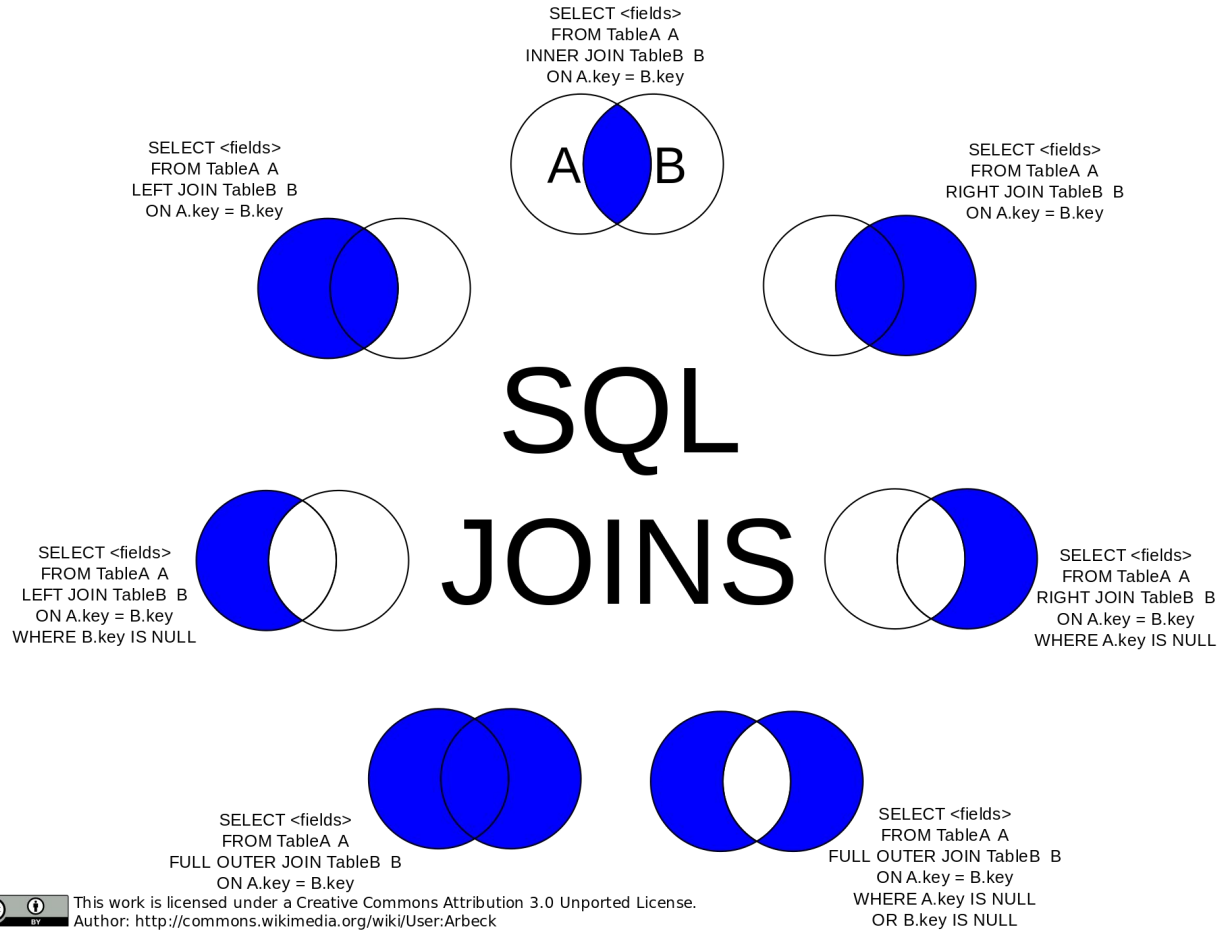


$$Y(t) = XG(t)$$

← The time dependent
'susceptibility' of Y to
X.

- X and Y are behaviors captured as events.
- We measure the transitions of customers as defined by these events.
- Y's relationship to X, can be dependent or independent, singular or multiple. Y can also be the same event as X.
- Measuring $G(t)$ tells us about the rate of change or growth between X and Y.

JOINS



This work is licensed under a Creative Commons Attribution 3.0 Unported License.
Author: <http://commons.wikimedia.org/wiki/User:Arbeck>

Join Example

```
SELECT
    count(DISTINCT(y.customer_id))
FROM
    master_fact x
JOIN
    master_fact y
ON
    x.customer_id = y.customer_id
WHERE
    x.event_type = 'X' AND
    y.event_type = 'Y'
```

Growth Susceptibility Query

```
SELECT
    floor((y.date_created - x.date_created) / (30*86400)) AS "Age",
    count(DISTINCT( y.customerid )) AS customers
FROM
    master_fact x
left join
    master_fact y
ON
    x.customer_id = y.customer_id
WHERE
    x.eventtype = 'X' AND
    y.eventtype = 'Y'
GROUP BY 1
```

Time Dependent Susceptibility



██████████



Python SQL Workshop #2

- ▶ Import master_fact table from eliflo.
- ▶ Explore how to do joins.
- ▶ Measure retention curves of a few cohorts.
- ▶ Look at retention performance over time.

PROJECTIONS & PREDICTIVE MODELS

A Thought Experiment

Day 3 of Launch

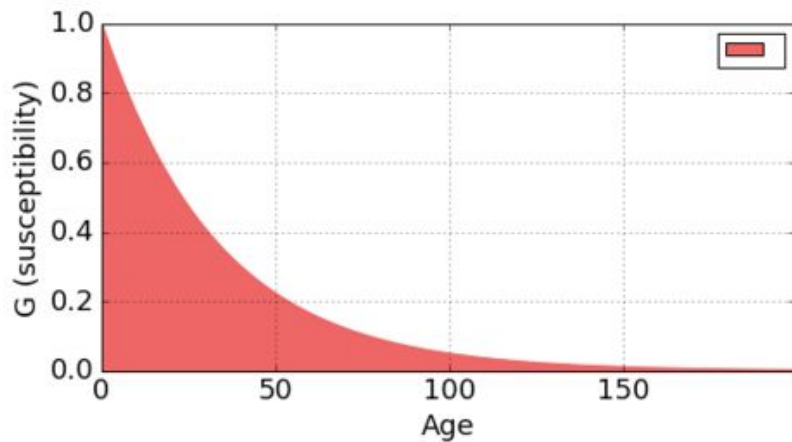
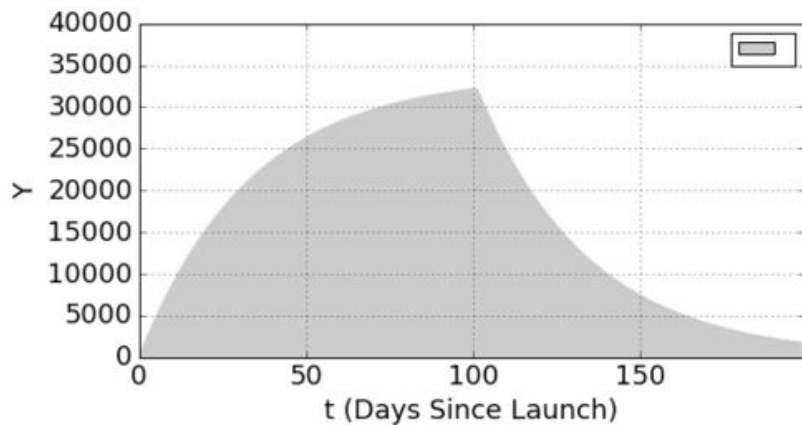
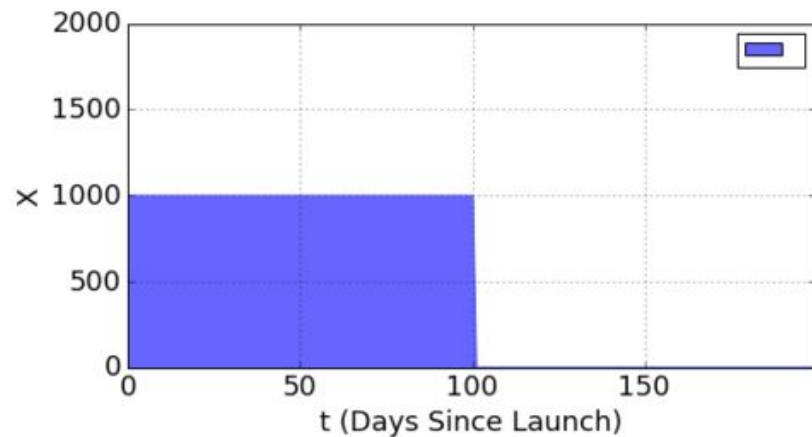
- ▶ 1000 new customers per day
- ▶ 50% return day after signup, 25% return 2 days later, none return after.
- ▶ How many users do we expect on day 5?

Generalize: If $N(t)$ customers join each day and $G(a)$ return how many do we expect to be active T days from now?

Forecasting

$$Y(t) = XG(t)$$

$$Y(t) = \sum_a X(t-a)G(a)$$



SQL for Predictive Models

Generalizable to
multiple features.

```
SELECT
    X.customer_id,
    X.feature,
    Y.dependent_variable
FROM
    (
    SELECT
        customer_id,
        count(*) as feature
    FROM
        master_fact
    WHERE
        event_type = 'X' and
        date_created in 'Some Time Period'
    GROUP BY 1
    ) X
LEFT JOIN
    (
    SELECT
        customer_id,
        count(*) as dependent_variable
    FROM
        master_fact
    WHERE
        event_type = 'Y' and
        date_created in 'Some Later Time'
    GROUP BY 1
    ) Y
ON
    X.customer_id = Y.customer_id
GROUP BY 1
```

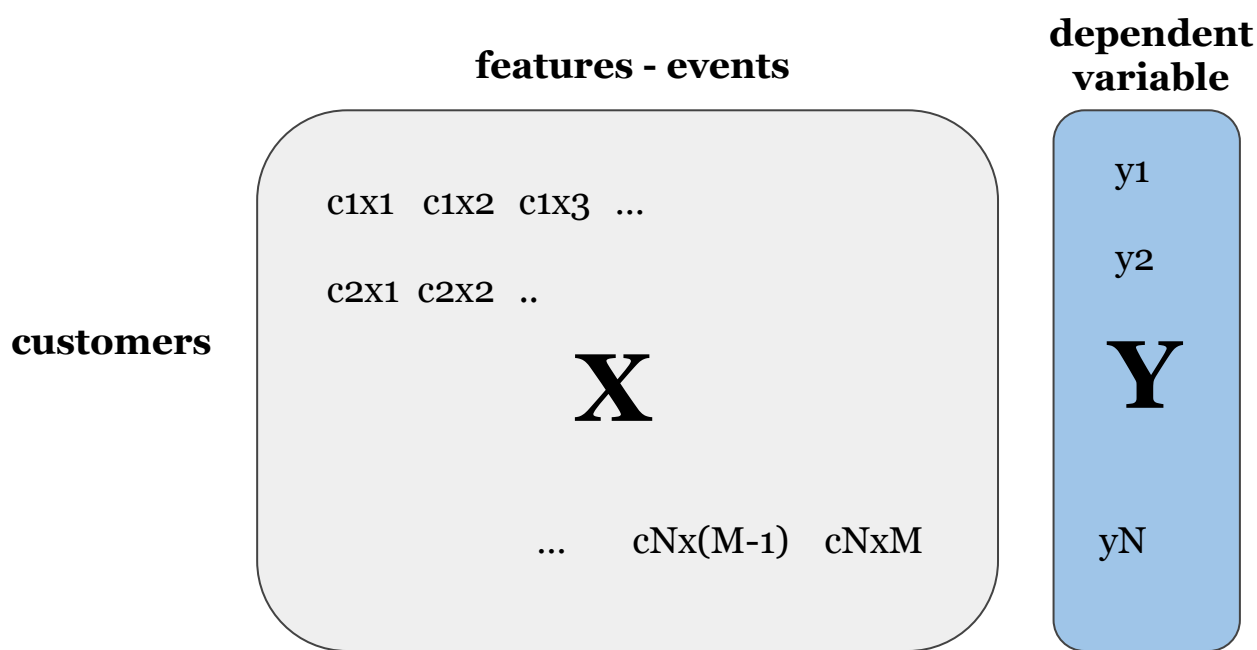
SQL for Predictive Models

Generalizable to multiple features.

Cohort Activity.sql

```
1  SELECT
2      x.user_id
3      x.feature,
4      y.dependent_variable
5  FROM
6      (SELECT
7          user_id,
8          count(*) as feature
9      FROM
10         event_log
11     WHERE
12         event_type = 'X' and
13         date_trunc(date_created, 'month') = "Some Month"
14     GROUP BY 1
15     ORDER BY 1) a
16  LEFT JOIN
17      (SELECT
18          user_id,
19          count(*) as dependent_variable
20      FROM
21         event_log
22     WHERE
23         event_type = 'Y' and
24         date_trunc(date_created, 'month') = "Some Month Later"
25     GROUP BY 1
26     ORDER BY 1) b
27  on a.user_id = b.user_id
28
29
```

Extending to Multiple X - Supervised Learning



- Rows are customers.
- Features are behavioral measures and latent attributes of those customers.
- Dependent variable, **Y**, is the customer's activity we want to predict.

Python SQL Workshop #3

- ▶ Convolve the retention curve with new users to compute future daily active users.
- ▶ Build a prediction model to estimate whether a customer will be active in their 3rd month.
- ▶ Write a query to create a data set you can use to build this model.