

# The Fundamentals of MongoDB Aggregation

## RDB

Databases are spreadsheets with metadata : For patients' work email, entire new entry would be made  
In RDB, data is organized with Primary keys (as index) in different tables : So, to look up at single patient, we traverse dozens of databases tables

Patient Contact Information									
Patient ID	Last Name	First Name	Personal Email	Work Email	Home Phone	Work Phone	Cell Phone	Emergency Contact Phone	Home Street
1	Doe	Jane	jdoe@acme.com		555-1212		555-5090		800 Table St
2	Deer	James	jdeer@msn.com		555-2121	555-1111			800 Relation Drive
3	Appl	Lenser	jappl@stglobal.net		202-555-0118		202-555-0193		678 George Street
4	Delroy	Trammel	jdelroy@icloud.com		202-555-0107	202-555-0107	202-555-0195		89 Coffee Dr
5	Shandy	Spence	carmend@comcast.net	jshandy@stglobal.net	202-555-0165			202-555-0183	8778 Spruce Ave
6	Cassandra	Yeats	kyeats@mac.com		202-555-0109		202-555-0124		483 Birchleaf Lane
7	Maryjo	Varey	avarey@mac.com		202-555-0182		202-555-0188		42 S. Highland Lane
8	Mira	Platt	plattm@stglobal.net		202-555-0198		202-555-0188		482 Pearl Street
9	Raymon	Ryer	ghoyer@gmail.com		202-555-0150	202-555-0150			21 St Louis Street
10	Kimber	Gravel	part@stglobal.net		202-555-0142				7378 Greenwood St
11	Sumiko	Cullinan	padm@stglobal.net		202-555-0154		202-555-0170		31 Nichols Court
12	Boris	Elzey	schweas@verizon.net	jdelroy@icloud.com	202-555-0182				8224 Eagle Drive
13	Dagmar	Morano	dmorano@yahoo.com		202-555-0126				3 Grove Dr
14	Trilla	Krukowski	trilla@stglobal.net		614-555-0164				272 Big Rock Cove
15	Lenny	Walcott	lwalcott@stglobal.net		614-555-0119				2 Phoenix Avenue
16	Jodie	Manion	manionj@comcast.net		614-555-0125			614-555-0519	980 Sycamore St
17	Erinda	Eisenmenger	teoflower@stglobal.net		614-555-0118	614-555-0118			77 Hillcrest Court
18	Jaquelyn	Deffon	jdeffon@stglobal.net		614-555-0196				84 Charles Rd
19	Barbie	Brandis	carrens@outlook.com		614-555-0187				11 Grove Drive
20	Appl	Lenser	jappl@stglobal.net	appls@comcast.net	202-555-0116				678 George Street

Figure 1 Relational databases : complication

- This complicates data understanding ability to write applications
- Hard to add features
- Fetching data from many sources [tables] is inefficient

## MongoDB

### Document Model

### Structuring with flexible schema is:

- Natural for programmers to read and code
- Easy for computers to process

Redundancy - Scalability - cloud native - **Fundamentally Distributed**

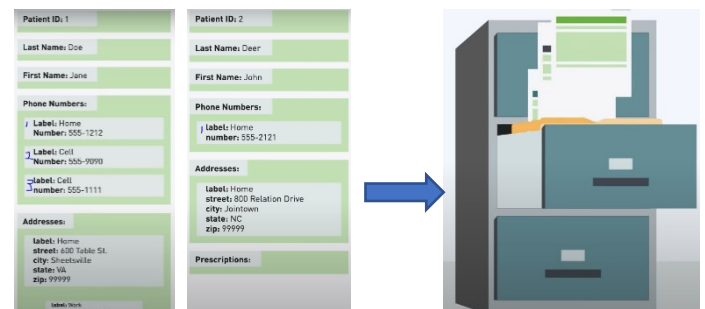


Figure 2 MongoDB document structure : flexible schema

## MongoDB Query vs MongoDB Aggregation

Aggregation framework enables developers to define FUNCTIONAL PIPELINES for data Preparation – Shaping – Analytics. Fundamental components are-

- **Stages** : Grouping-Sorting-Shaping etc.
- **Expressions** : Logic units of functionality

## Concept of Functional Pipeline

Assembly line of 1:N stages configurable for transformations using 1:N aggregation operators Or Expressions

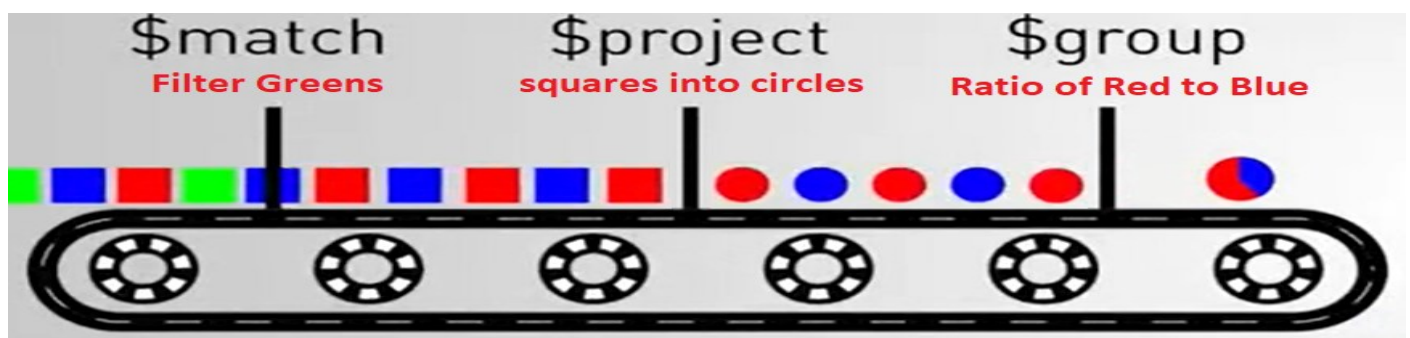
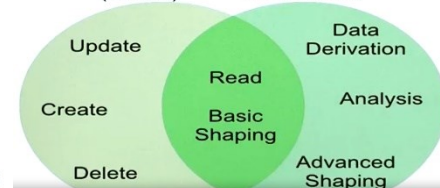


Figure 3 Projection pipeline

### Why we use aggregation?

- Need to prepare data
  - Cleaning
  - Deriving information
- Ask preliminary questions

### Query Framework (CRUD) vs Aggregation Framework



## Aggregation framework: Structure and syntax

### Aggregation Pipeline Quick Reference

- **Operators** typically appear in **Key** position
- **Expressions [act like functions]** typically appear in **Value** position

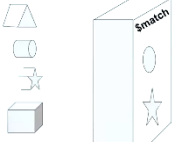

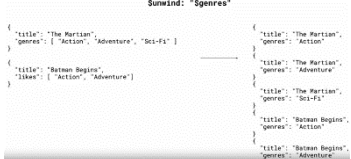
Field Path: "\$fieldName" ("NumberOfMoons")  
Access field in doc

System Variable: "\$\$UPPERCASE" ("\$\$CURRENT")  
System level variable 'CURRENT'

User Variable: "\$\$foo" User variable

```
Cluster0-shard-0:PRIMARY> db.solarSystem.aggregate([
...   $match: {
...     atmosphericComposition: { $in: [/O2/] },
...     meanTemperature: { $gte: -40, $lte: 40 }
...   }, {
...     $project: {
...       _id: 0,
...       name: 1,
...       hasMoons: { $gt: ["NumberOfMoons", 0] }
...     }
...   }, { allowDiskUse: true }
... ])
```

Aggregation operator  
Query operator  
Expression  
Argument

Operator	Syntax	Description	Details
<b>MATCH</b>	db.<DBname>.aggregate([ \$match : {<query>} ])	<ul style="list-style-type: none"> <li>• Cannot use \$where</li> <li>• Does not have projection</li> <li>• Match should be first operator if \$text is used</li> <li>• First stage match increases query throughput as it can take adv of Indexing</li> </ul>	<b>Filter operation</b> 
<b>PROJECT</b>	db.<DBname>.aggregate([ \$project : { _id:0, <specs>} ])	<ul style="list-style-type: none"> <li>• Select, Remove, Reassign, derive new fields</li> <li>• Similar to “<b>map</b>” function in python</li> <li>• Mention all fields to retain; others are removed auto</li> <li>• Except “_id” : needs explicit removal</li> </ul>	
<b>GROUP</b>	db.<DBname>.aggregate([ \$group: { _id: <expression>, <field1>: { <acc1>: <expr1> },...} ])	<ul style="list-style-type: none"> <li>• Group by a column naturally</li> <li>• Layer of detail using accumulator such as “total_amt_coin_type” : { \$sum : 1 }</li> </ul>	
<b>REDUCE</b>	db.<DBname>.aggregate([ \$reduce: { input: <array>, initialValue: <expr>, in: <expr> } ])	<ul style="list-style-type: none"> <li>• Act on “input” array with “initial” value [acc] with logic “in”; and update “init” at each stage</li> </ul>	
<b>UNWIND</b>	db.<DBname>.aggregate([ \$unwind: { path: <field path>, includeArrayIndex: <string>, preserveNullAndEmptyArrays: <Bool> }) ])	<ul style="list-style-type: none"> <li>• unwind an array field creating a new document for every entry where each field values are now a separate entry.</li> </ul>	
<b>MAP</b>	db.<DBname>.aggregate([ \$map: { input: <expression>, as: <string>, in: <expression> } ])	<ul style="list-style-type: none"> <li>• Applies an expression to each item in an array and returns an array with the applied results.</li> </ul>	
<b>LET</b>	db.<DBname>.aggregate([ \$let: { vars: { <var1>: <expression>, ... }, in: <expression> } ])	<ul style="list-style-type: none"> <li>• Binds variables “<b>var</b>” for use in the specified expression and returns the result of the “in” expression.</li> </ul>	

<b>LOOKUP</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$lookup: {     from: &lt;collection to join&gt;,     localField: &lt;fields input documents&gt;,     foreignField:&lt;field "from" collection&gt;,     as: &lt;output array field&gt; } ])</pre>	<ul style="list-style-type: none"> <li>• <b>Left outer join</b> – strict equality comparison</li> <li>• <b>Collection to “from” CAN NOT be sharded</b></li> </ul>	
<b>GRAPH LOOKUP</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$graphLookup: {     from: &lt;collection&gt;,     startWith: &lt;expression&gt;,     connectFromField: &lt;string&gt;,     connectToField: &lt;string&gt;,     as: &lt;string&gt;,     maxDepth: &lt;number&gt;,     depthField: &lt;string&gt;,     restrictSearchWithMatch: &lt;document&gt;} ])</pre>	<ul style="list-style-type: none"> <li>• Run on primary shard only in sharded environment</li> <li>• <b>Collection to “from” CAN NOT be sharded</b></li> <li>• <b>Final 100MB RAM limit – No effect {allowDiskUsage : True}</b></li> </ul>	
<b>EXPRESSIVE LOOKUP</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$lookup: {     from: &lt;collection to join&gt;,     let: { &lt;var_1&gt;: &lt;expression&gt;, ..., &lt;var_n&gt;: &lt;expression&gt; },     pipeline: [ &lt;pipeline to execute on collection to join&gt; ],     as: &lt;output array field&gt;   } ])</pre>	<ul style="list-style-type: none"> <li>• Regular lookup is essentially join – creates new field DISSASOCIATED WITH OLD SUBFIELDS</li> <li>• Not good for a NoSQL db : hence old subfields are typically remapped</li> <li>• In expressive lookup, Shaping is done before joining : same task within less RAM</li> <li>• Pipelines execute with context of collection specified in “from” field</li> </ul>	
<b>ADD FIELDS</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$addFields : {     &lt;newField&gt;: &lt;expression&gt;,     ... } ])</pre>	<ul style="list-style-type: none"> <li>• appends new fields to existing documents</li> </ul>	
<b>FACET</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$facet:   {     &lt;outputField1&gt;: [ &lt;stg1&gt;, &lt;stg2&gt;,...],     &lt;outputField2&gt;: [ &lt;stg1&gt;, &lt;stg2&gt;,...],     ...   } ])</pre>	<ul style="list-style-type: none"> <li>• categorize and group incoming documents</li> </ul>	
<b>SET DIFFERENCE</b>	<pre>db.&lt;DBname&gt;.aggregate([   { \$setDifference:     [ &lt;expression1&gt;, &lt;expression2&gt; ]   } ])</pre>	<ul style="list-style-type: none"> <li>• Input sets [A,B] and yields values that exist only in A</li> </ul>	
<b>ARRAY ELEMENT AT</b>	<pre>db.&lt;DBname&gt;.aggregate([   \$arrayElemAt: [ &lt;array&gt;, &lt;idx&gt; ] ])</pre>	<ul style="list-style-type: none"> <li>• Returns the element at the specified array index</li> </ul>	

<b>SET UNION</b>	db.<DBname>.aggregate([{\$setUnion: [ <expression1>, <expression2>,...] }])	<ul style="list-style-type: none"> <li>• Takes two or more arrays and returns an array containing the elements that appear in any input array</li> </ul>	
<b>SAMPLE</b>	db.<DBname>.aggregate([{\$sample: { size: <positive integer> } }])	<ul style="list-style-type: none"> <li>• Randomly selects the specified number of documents from its input</li> </ul>	
<b>EXPR</b>	{ \$expr: { <expression> } }	<ul style="list-style-type: none"> <li>• Allows the use of aggregation expressions within the query language.</li> </ul>	
<b>FIRST, LAST</b>	{ \$first: < expr > } { \$last: < expr > }	<ul style="list-style-type: none"> <li>• only for <b>GROUP</b></li> <li>• Returns the value that results from applying expression to first/ last document in a group sharing _id</li> </ul>	
<b>MULTIPLY</b>	db.<DBname>.aggregate([{\$multiply : [expr1>, <expr2>,<expr3>, ...] }])	<ul style="list-style-type: none"> <li>• Expr1*Expr2*...ExprN</li> </ul>	
<b>DIVIDE</b>	db.<DBname>.aggregate([{\$divide: [<expr1>,<expr2>] }])	<ul style="list-style-type: none"> <li>• Expr1 / Expr2</li> </ul>	
<b>SUM , MAX, MIN, Average, Standard Deviation</b>	{ \$sum: <expression> } { \$max: < expr > } { \$min: < expr > } { \$avg: < expr > } { \$stdDevPop: <expr> } { \$stdDevSamp: <expr> }	<ul style="list-style-type: none"> <li>• Have memory within scope of document supplied within project. Need unwind and group to operate on all documents</li> </ul>	

## MongoDB expressions

### Boolean

\$and  
\$or  
\$not

### Set

\$setEquals  
\$setIntersection  
\$setUnion  
\$setDifference  
\$setIsSubset  
\$anyElementTrue  
\$allElementsTrue

### Comparison

\$cmp  
\$eq  
\$gt  
\$gte  
\$lt  
\$lte  
\$ne

### Variable

\$let

### Literal

\$literal

### Arithmetic

\$abs  
\$add  
\$ceil  
\$divide  
\$exp  
\$floor  
\$ln  
\$log  
\$log10  
\$mod  
\$multiply  
\$pow  
\$sqrt  
\$subtract  
\$trunc

### String

\$concat  
\$indexOfBytes  
\$indexOfCP  
\$split  
\$strlenBytes  
\$strlenCP  
\$strcasecmp  
\$substrBytes  
\$substrCP  
\$toLower  
\$toUpper

### Text Search

\$meta

### Array

\$arrayElemAt  
\$arrayToObject  
\$concatArrays  
\$filter  
\$in  
\$indexOfArray  
\$isArray  
\$map  
\$objectToArray  
\$range  
\$reduce  
\$reverseArray  
\$size  
\$slice  
\$zip

### Data Type

\$type

### Conditional

\$cond  
\$ifNull  
\$switch

### Date

\$dayOfYear  
\$dayOfMonth  
\$dayOfWeek  
\$year  
\$month  
\$week  
\$hour  
\$minute  
\$second  
\$millisecond  
\$dateToString  
\$isoDayOfWeek  
\$isoWeek  
\$isoWeekYear

### Accumulators

\$sum  
\$avg  
\$first  
\$last  
\$max  
\$min  
\$push  
\$addToSet  
\$stdDevPop  
\$stdDevSamp

## Cursor like stages

- { \$sort: { <field1>: <sort order>,... } : sort on particular field/s

- Use early in pipeline to take advantage of indexes. If put after project stage, it performs in-memory sort
- Sort has default 100MB RAM limit – to exceed, use aggregation option {allowDiskUsage : True}
- { **\$skip**: <positive integer> } : return all skipping first N → [N:]
- { **\$limit**: <positive integer> } : return first N → operation [:N]
- { **\$count**: <string – name of output field> } : return count of documents returned

**Links:**

[Query and Projection Operators](#)

# Leveraging MongoDB's Flexible Schema

## Importance of Schema Design

Operator	Syntax	Description	Details
<b>MONGO IMPORT</b>	mongo "mongodb+srv://cluster0.jotdp.mongodb.net/<dbname>" --username rasikMooc  Mongoimport --uri mongodb+srv://rasikMooc:rasikMooc@cluster0.jotdp.mongodb.net/aggregation-framework-mooc --collection retail --type csv --file retail.csv --headerline --drop	<ul style="list-style-type: none"><li>• Import csv/ uri data</li><li>• --drop for removing earlier data</li><li>• --mode= [ upsert/insert/merge ] -- upsertFields=[key fields]</li><li>• <b>No shape change / derived fields</b></li></ul>	
<b>CREATE VIEW</b>	db.createView(<view>, <source>, <pipeline>, <options>)	<ul style="list-style-type: none"><li>•</li></ul>	
<b>SYSTEM VIEW</b>	db.system.views	<ul style="list-style-type: none"><li>• contains information about each view in the database</li></ul>	
<b>ADD TO SET</b>	{ \$addToSet: { <field1>: <value1>, ... } }	<ul style="list-style-type: none"><li>• Add values to array omitting duplicates</li></ul>	
<b>PUSH</b>	{ \$push: { <field1>: <value1>, ... } }	<ul style="list-style-type: none"><li>• Appends a specified value to an array</li></ul>	

## Schema Exploration

Make schema explicit

- In SQL, schema is about DATA & RELATIONS
- In NoSQL, schema is about **Performance**
  - Access patterns and frequent queries : What is purpose behind query
- **Schema validation is optional**

## Exploration options

- Collection.find\_one() : naïve
- mongoDB Compass

## Data Migration Processes

- Mongoimport [import, change field type] and aggregation [shape change, field derivation] together offer fetch and process capacity
- In-place operations cannot be done with aggregation

## VIEWS

**db.createView(<view>, <source>, <pipeline>, <options>)**

- Views are public
- Non materialized – computed separately for each read operation
- “Aggregation Pipeline as a collection”
- Create slices of DB
  - **Horizontal** [eg. \$match - select reduced number of documents (based on criteria) with **same shape** i.e. cardinality]
  - **Vertical** [eg. \$project – return same no of documents with **different shape** i.e. reduced cardinality (based on criteria)]

- No permission for : write / mapReduce/ index/ \$text/ \$geoNear/ renaming/ find() with \$elemMatch, \$slice, \$meta

Supplementary schema with accumulators : \$ addFields

Use cases of \$graphLookup

Tree like architectures

Entity resolution

## Regular Lookup

```
MongoDB Enterprise > db.air_alliances.aggregate([
... {
...   $lookup: {
...     from: "air_airlines",
...     localField: "airlines",
...     foreignField: "name",
...     as: "airlines"
...   },
... },
... {
...   $addFields: {
...     airlines: {
...       $map: {
...         input: "$airlines",
...         in: {
...           name: "$$this.name",
...           alias: "$$this.alias",
...           id: "$$this.airline"
...         }
...       }
...     }
...   }
... ]
```

## expressive lookup

```
MongoDB Enterprise > db.air_alliances.aggregate([
... {
...   $lookup: {
...     from: "air_airlines",
...     let: { constituents: "$airlines" },
...     pipeline: [
...       {
...         $match: { $expr: { $in: ["$name", "$$constituents"] } }
...       },
...       {
...         $project: {
...           _id: 0,
...           name: 1,
...           alias: 1,
...           id: "$airline"
...         }
...       }
...     ],
...     as: "airlines"
...   }
... ]
```

## giving same outcome

```
"_id" : ObjectId("5980bef9a39d0ba3c650ae9d"),
"name" : "OneWorld",
"airlines" : [
  {
    "name" : "Canadian Airlines",
    "alias" : "CP",
    "id" : 1615
  },
  {
    "name" : "British Airways",
    "alias" : "BA",
    "id" : 1355
  },
  {
    "name" : "American Airlines",
    "alias" : "AA",
    "id" : 24
  }
],
```



# Machine Learning with MongoDB

**Pearson-Rho** [Pearson correlation coefficient OR bivariate coefficient]

- Strength of linear relation between variables

$$r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

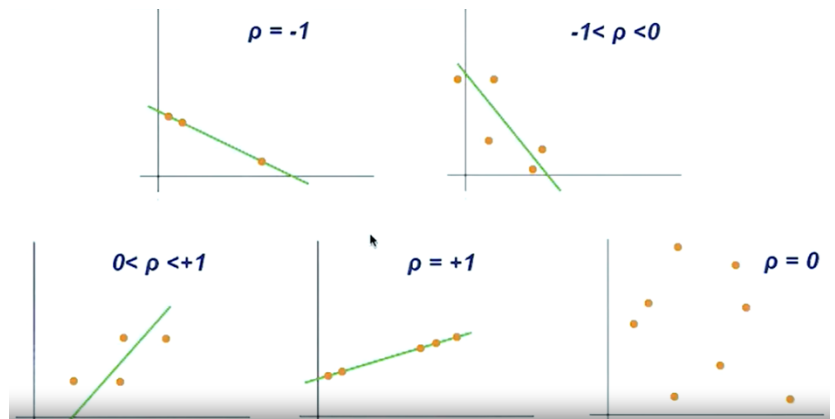


Figure 4 Pearson's Coefficient of correlation

## Market Basket Analysis

**Total Inventory I = {apple, beverage, chips...}**

**Bought Basket b = subset{I}**

### Associative rules

- **Support : likelihood of a basket**
  - $\text{Support}(b) : N(b) / N(\text{Transaction})$
- **Confidence : likelihood of a basket[A,B] given an item[A]**
  - $\text{Confidence}(A \rightarrow B) : \text{Support}(U\{A,B\}) / \text{Support}(A)$
  - Drawback : if consequent[i.e. B] is popular; then implication is indicated, though there isn't any.
- **Lift : Likelihood of another item[B] given an item[A]**
  - $\text{lift}(A \rightarrow B) : \text{Support}(U\{A,B\}) / [\text{Support}(A) * \text{Support}(B)]$
  - lift = 1 NO association
  - lift < 1 - ve association
  - lift > 1 + ve association

## Principal Component Analysis - Reduce noise by reducing dimensionality

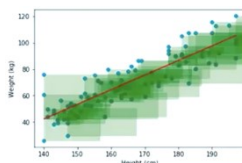
Mathematical steps for PCA:

- Input multidimension data
- Drop target variable
- Calculate covariance matrix of [independent variables]<sup>T</sup>
- Calculate eigenvalue and corresponding eigenvectors
- Sort by eigenvalue
- Take cumulative % sum of eigenvalue
- Choose top N eigenvalues contributing towards n%
- Use eigenvectors of chosen eigenvalues as new Dimensions for further regression/ classification/ clustering

## Linear Regression

- Least square & Mean squared error : minimize area under the squares

MongoDB help transforming large data and pipe it into scikit-learn



## Linear Regression Methods

- |  |                                  |
|--|----------------------------------|
| • Ordinary least squares               | • Ridge regression               |
| • Generalized least squares            | • Least absolute deviations      |
| • Percentage least squares             | • Bayesian linear regression     |
| • Iteratively reweighted least squares | • Quantile regression            |
| • Total least squares                  | • Principal component regression |
| • Maximum likelihood estimation        | • Least-angle regression         |



## Decision Tree

**Precision** : True +ve vs false +ve [imprison most notorious ones, but might miss few less important ones]

**Recall** : True + against false -ve [general model : imprison all rather than missing few]

**F1 score** : harmonic mean of Precision and recall

Random forest : Ensemble of decision trees [averages results from all models]

## Clustering

- Assign random centroids
- Calculate  $d\{\text{point, centroid}\}$
- Assign point to nearest centroid [dmin]
- Tune centroid to minimize distance from all points in its cluster
- Repeat

- **k-Means clustering**

Data normalization helps improving accuracy

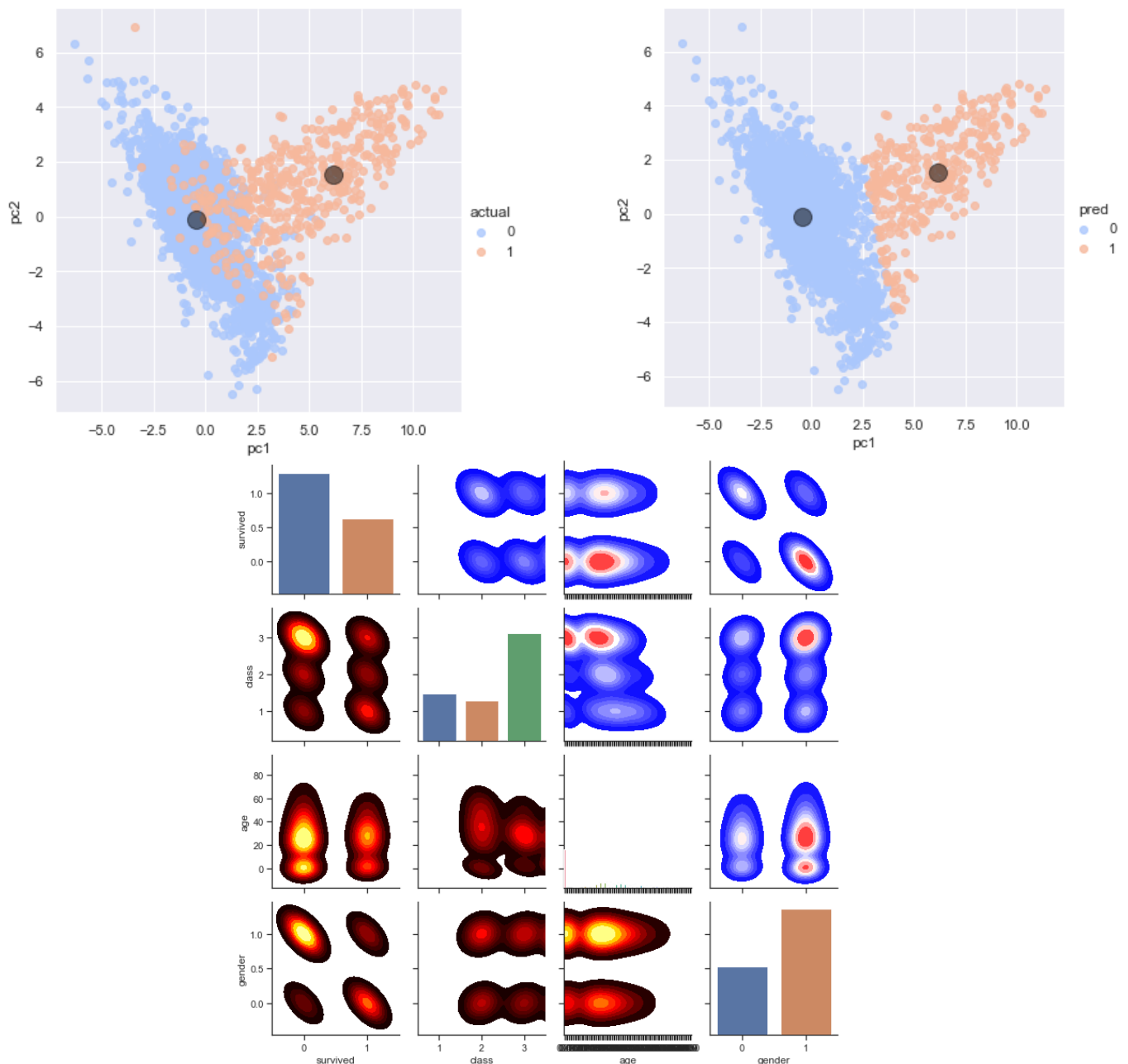


Figure 5 Titanic Dataset correlation distribution