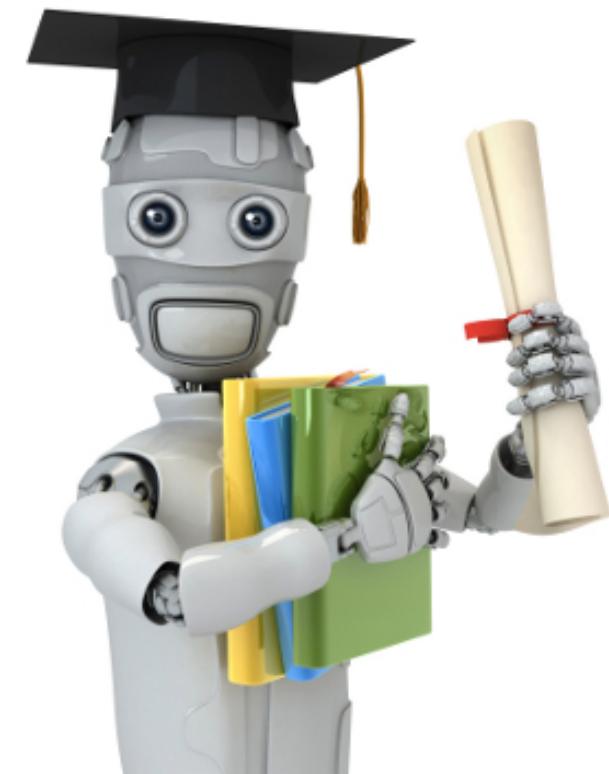


# **VERY BASIC OVERVIEW OF STATISTICS AND MACHINE LEARNING**

INTRODUCTION TO DATA SCIENCE

ELI UPFAL



# MACHINE LEARNING – exciting!



# MACHINE LEARNING – exciting!



# STATISTICS - boring



**MACHINE LEARNING – exciting!**



**STATISTICS - boring**



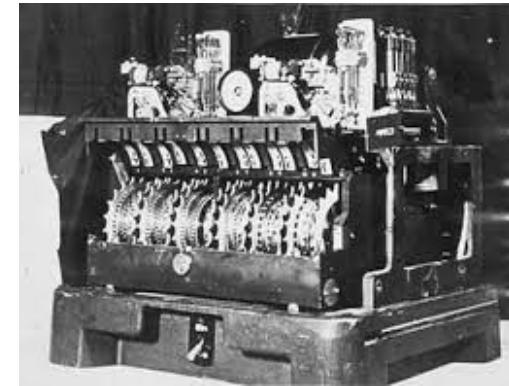
**ACTULLY – not that different**

# EXTRACTING INFORMATION FROM DATA

Data



Analysis



Predictions



Model

**“IT’S DIFFICULT TO MAKE PREDICTIONS,  
ESPECIALLY ABOUT THE FUTURE”**

## **Example of ML predictions:**

**Observation:** Millions of connections to web sites.

If you connect from IP 128.148.31.5, your computer type is ...., your operating system is ...., your mouse movements profile is ...., then your likely age is .... your likely disposable income is .... your likely gender is ...., your likelihood of a purchase is...

Works very well in practice!

**“IT’S DIFFICULT TO MAKE PREDICTIONS,  
ESPECIALLY ABOUT THE FUTURE”**

## **Example of ML predictions:**

**Observation:** Millions of connections to web sites.

If you connect from IP 128.148.31.5, your computer type is ...., your operating system is ...., your mouse movements profile is ...., etc. then your likely age is .... your likely disposable income is .... your gender is ...., your likelihood of a purchase is...

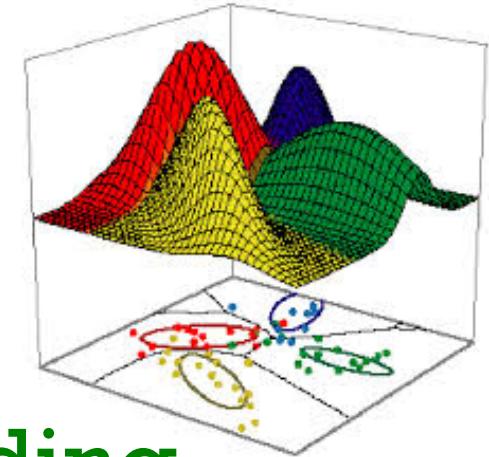
Works very well in practice!

**Observation:** Among flu patients, those who take medicine have a longer recovery

If you take medicine for flu then you have longer recovery.

Obviously wrong!

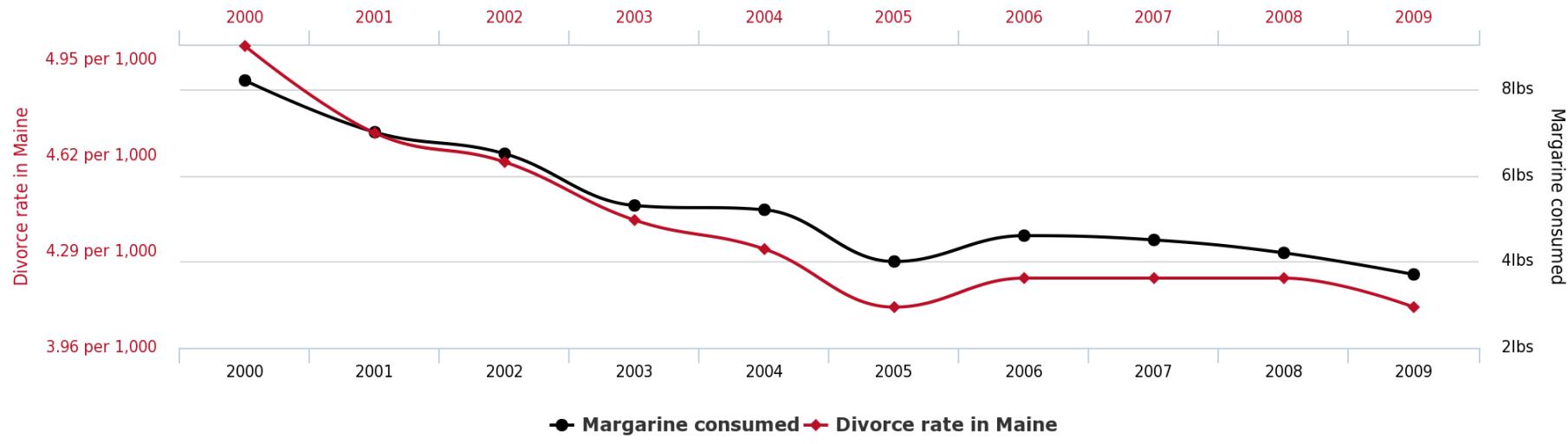
**Just because a machine learning,  
data mining, or data analysis  
application outputs a result - it  
doesn't mean that it's right**



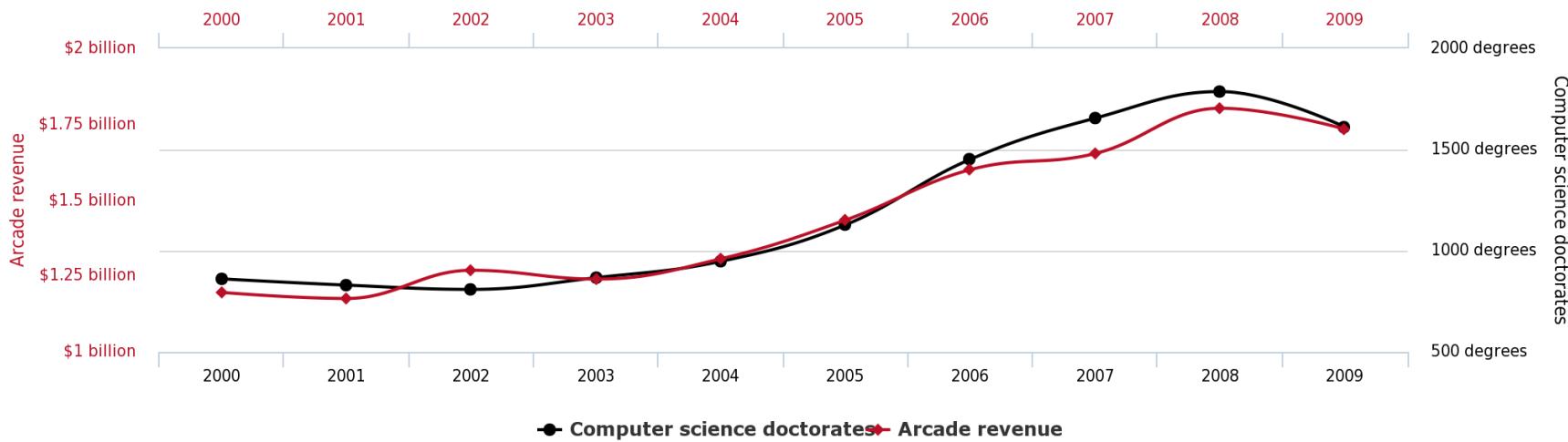
Data analysis is often misleading

Machine learning without statistical  
analysis is pure nonsense

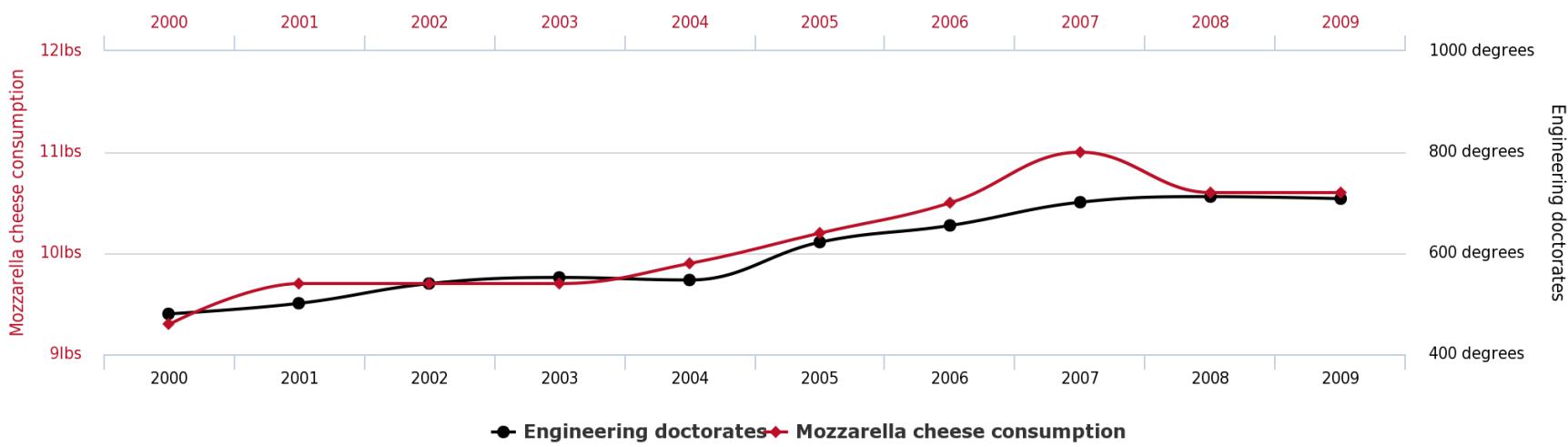
**Divorce rate in Maine**  
correlates with  
**Per capita consumption of margarine**



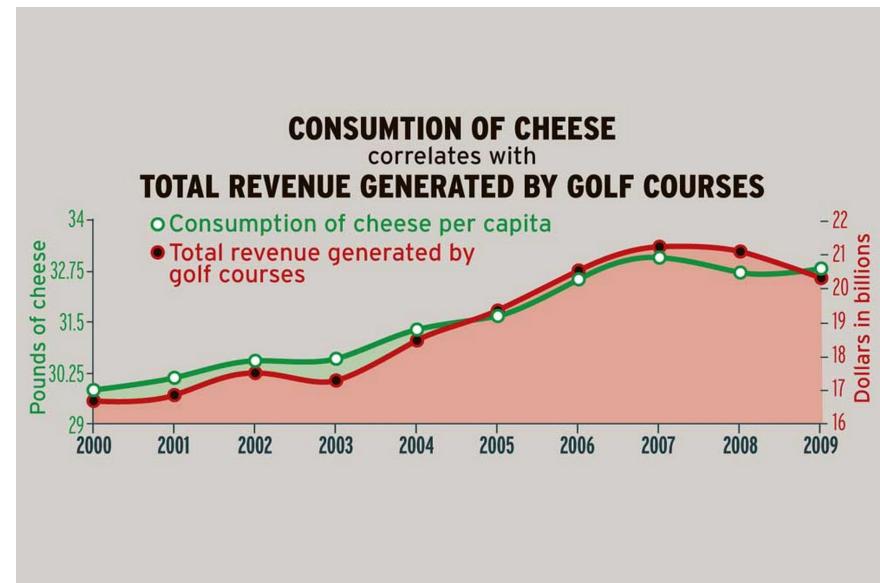
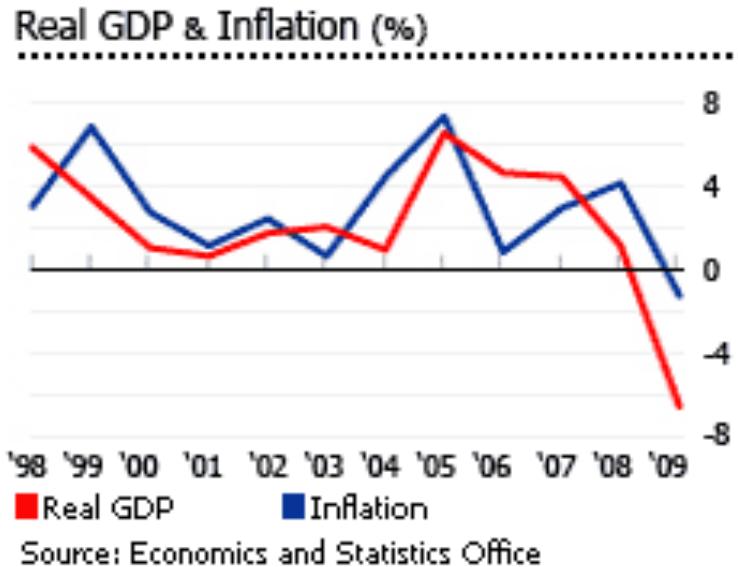
**Total revenue generated by arcades**  
correlates with  
**Computer science doctorates awarded in the US**



## Per capita consumption of mozzarella cheese correlates with Civil engineering doctorates awarded

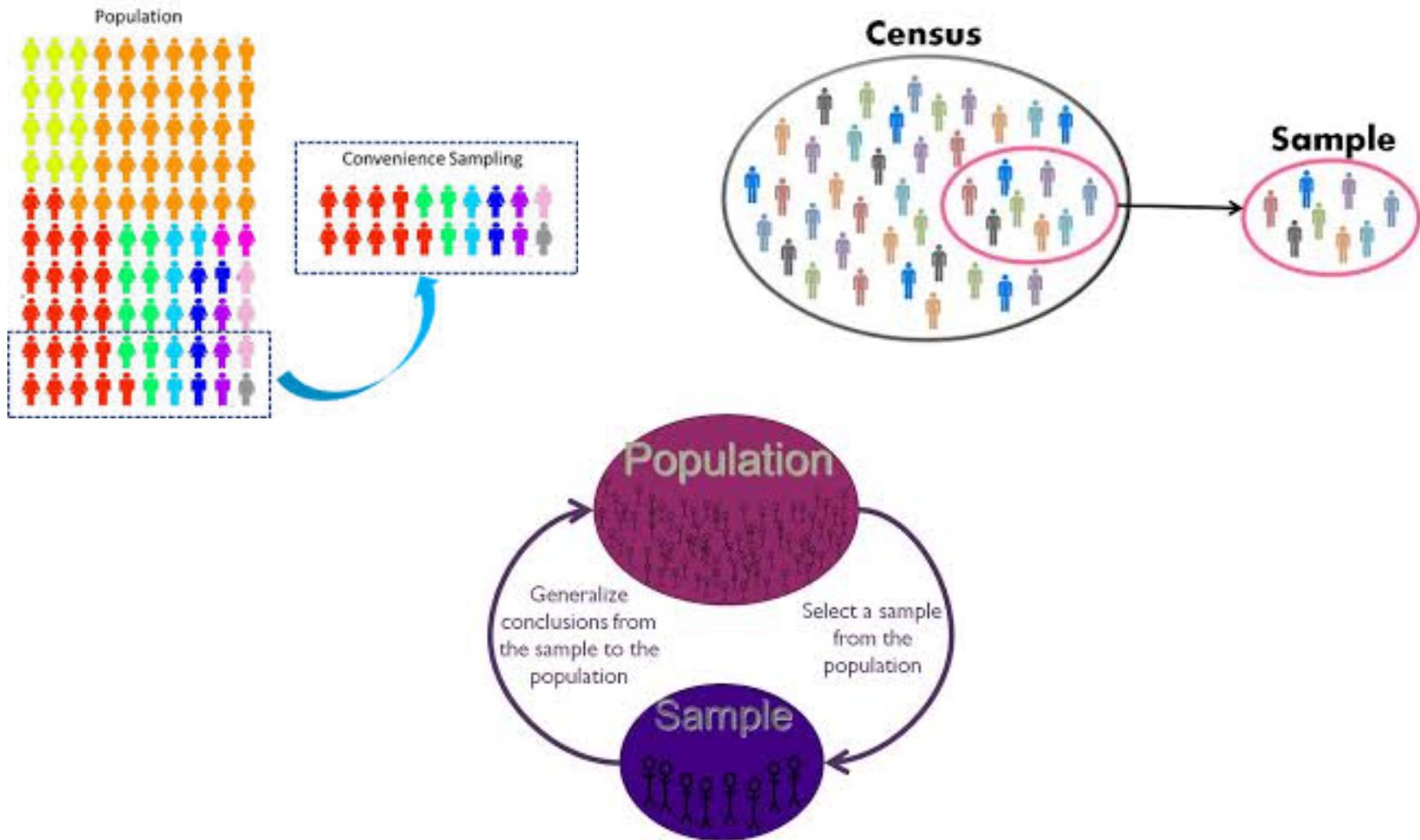


# Which Correlation is more convincing?



How do we distinguish between facts and coincident?

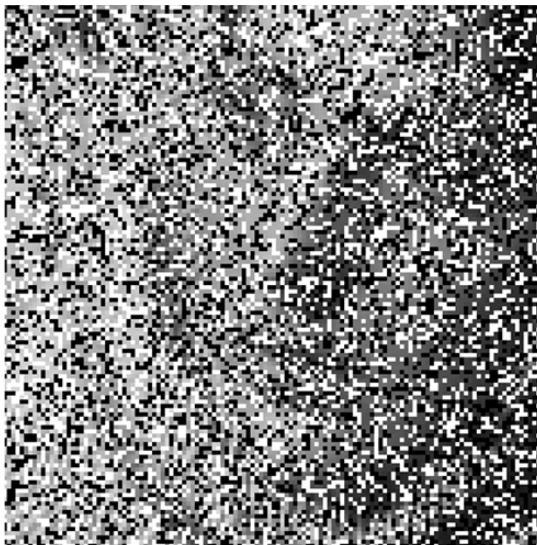
# DATA: LEARNING FROM A TRAINING SET (SAMPLE)



# THE DATA IS THE MODEL



(a)



(b)



(c)



(d)



(e)



(f)

# STATISTICS VS. MACHINE LEARNING

## **First there was statistics:**

Strict criteria for when an hypothesis ("discovery") is statistically significant

Strong assumptions, elaborate computation



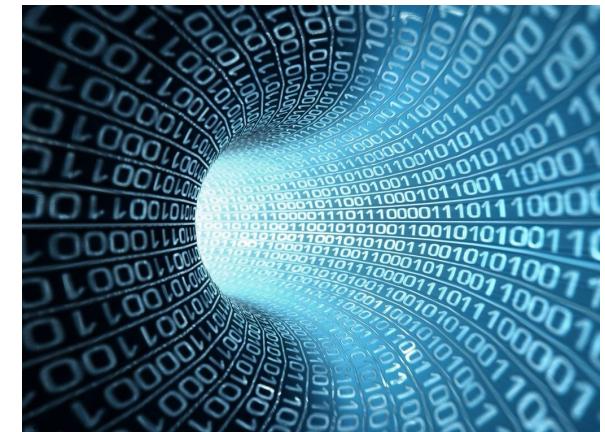
## **Then came Computer Science:**

Emphasize on efficient computation

Output best approximation, even if not certain

## **... And a lot of BIG data**

With lucrative business applications.



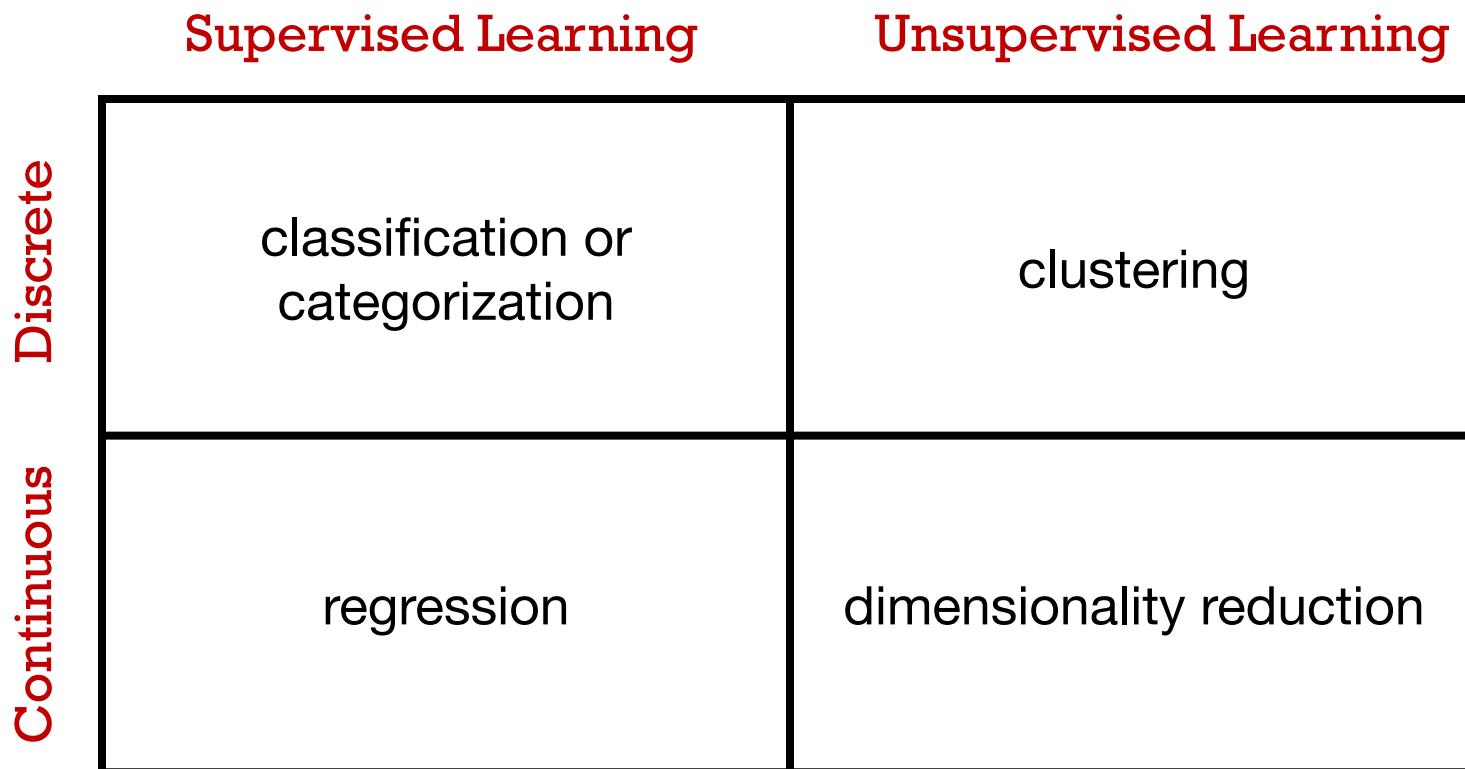
# **GOAL OF THIS PART OF THE COURSE**

**Basic machine learning tools for data analysis**

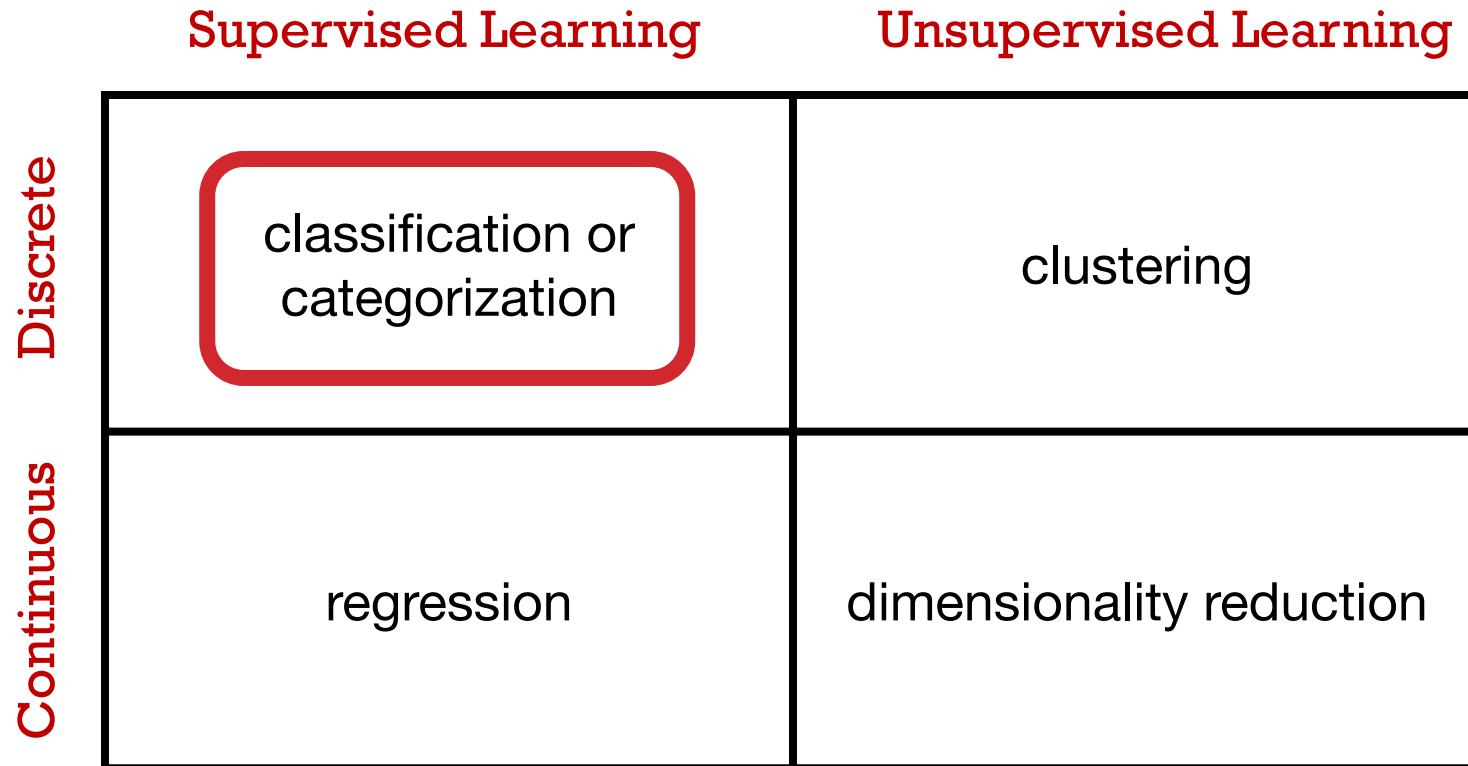
**Very basic concepts in probability and statistics**

**Understanding the power and pitfalls of data analysis**

# MACHINE LEARNING PROBLEMS



# MACHINE LEARNING PROBLEMS



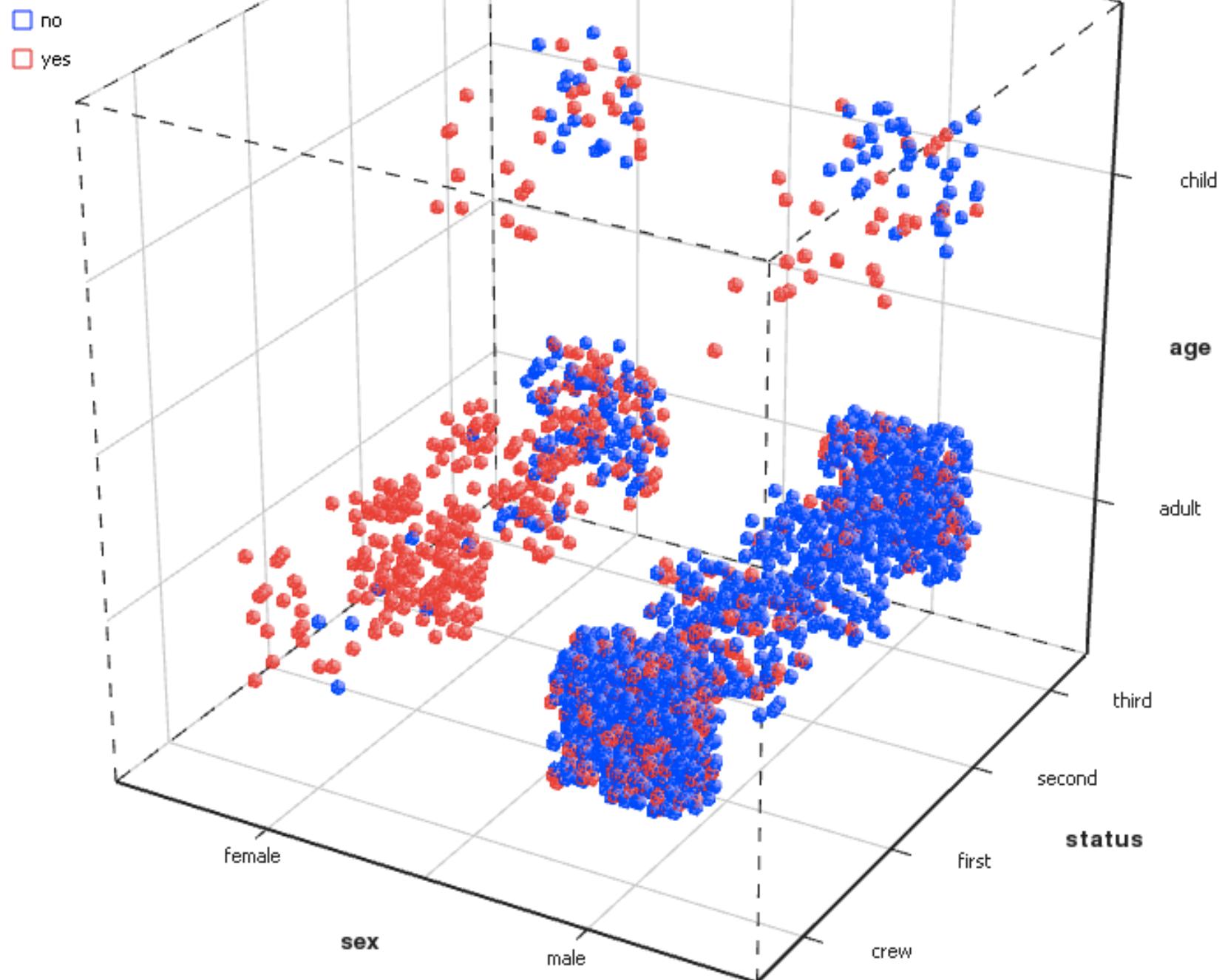
# EXAMPLE: TITANIC DATASET

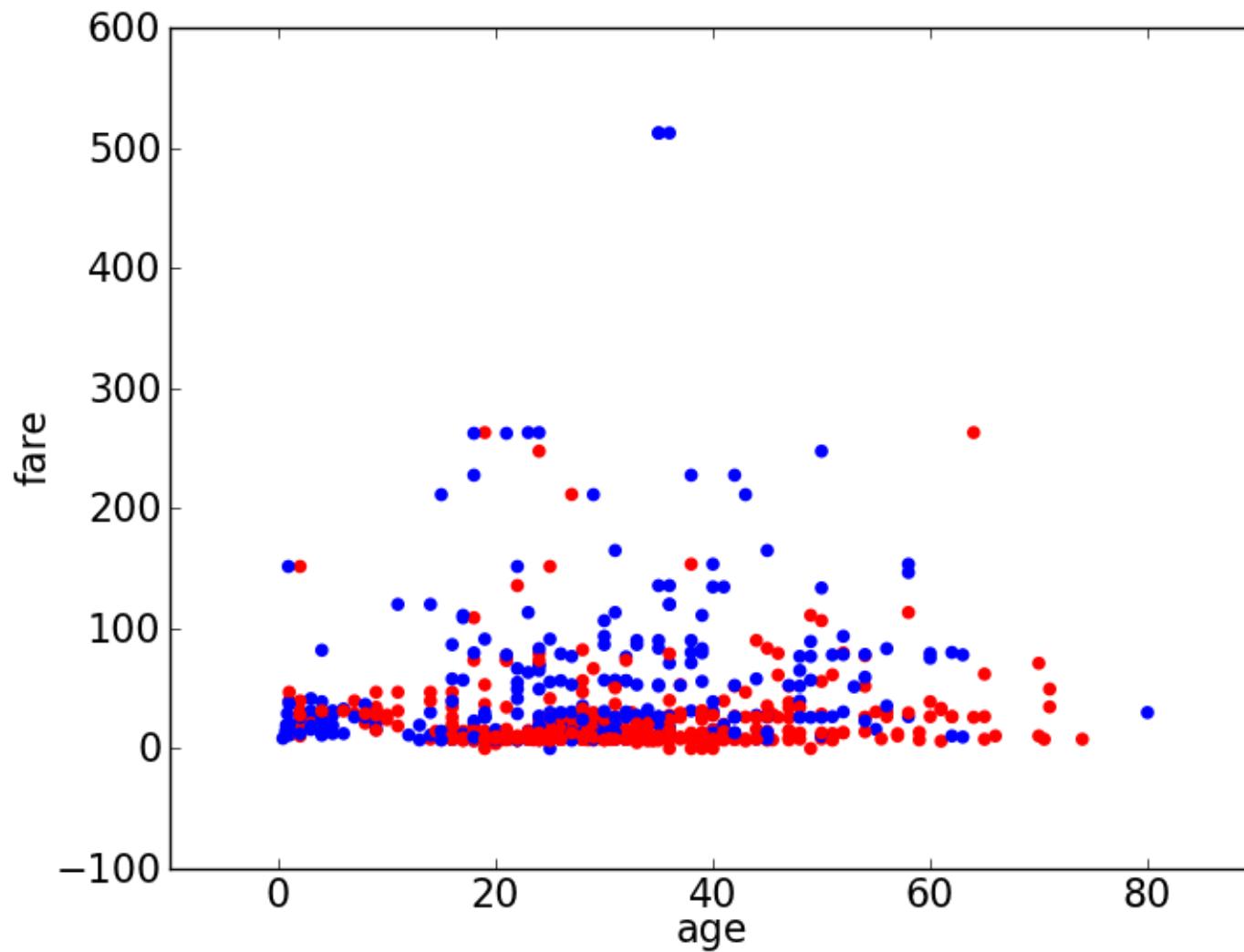
Label    Features

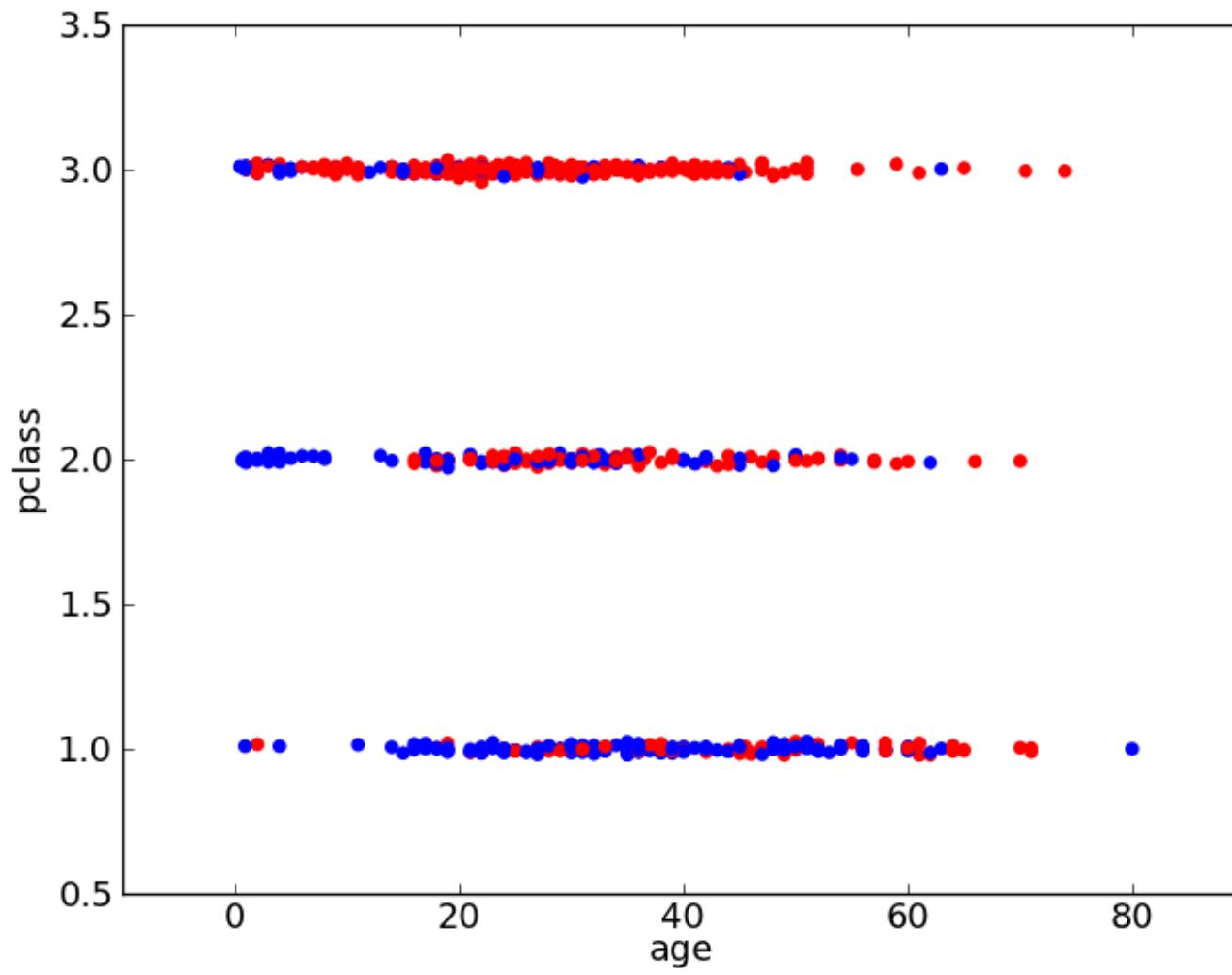
| survived | pclass | sex    | age | sibsp | parch | fare    | cabin | embarked |
|----------|--------|--------|-----|-------|-------|---------|-------|----------|
| 0        | 3      | male   | 22  | 1     | 0     | 7.25    |       | S        |
| 1        | 1      | female | 38  | 1     | 0     | 71.2833 | C85   | C        |
| 1        | 3      | female | 26  | 0     | 0     | 7.925   |       | S        |
| 1        | 1      | female | 35  | 1     | 0     | 53.1    | C123  | S        |
| 0        | 3      | male   | 35  | 0     | 0     | 8.05    |       | S        |
| 0        | 3      | male   |     | 0     | 0     | 8.4583  |       | Q        |
| 0        | 1      | male   | 54  | 0     | 0     | 51.8625 | E46   | S        |
| 0        | 3      | male   | 2   | 3     | 1     | 21.075  |       | S        |
| 1        | 3      | female | 27  | 0     | 2     | 11.1333 |       | S        |
| 1        | 2      | female | 14  | 1     | 0     | 30.0708 |       | C        |
| 1        | 3      | female | 4   | 1     | 1     | 16.7    | G6    | S        |
| 1        | 1      | female | 58  | 0     | 0     | 26.55   | C103  | S        |
| 0        | 3      | male   | 20  | 0     | 0     | 8.05    |       | S        |

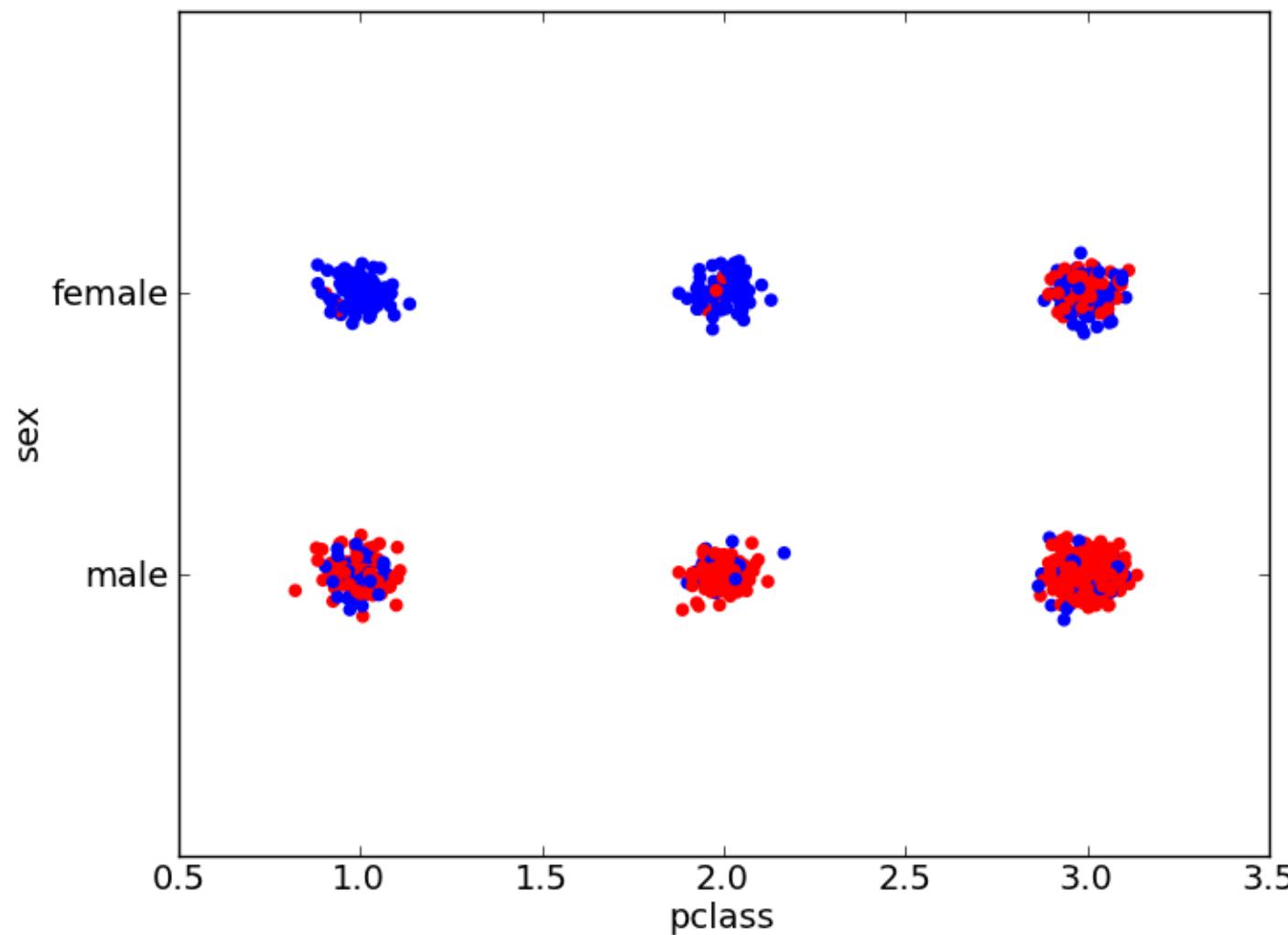
Can we predict survival from these features?

**survived:**

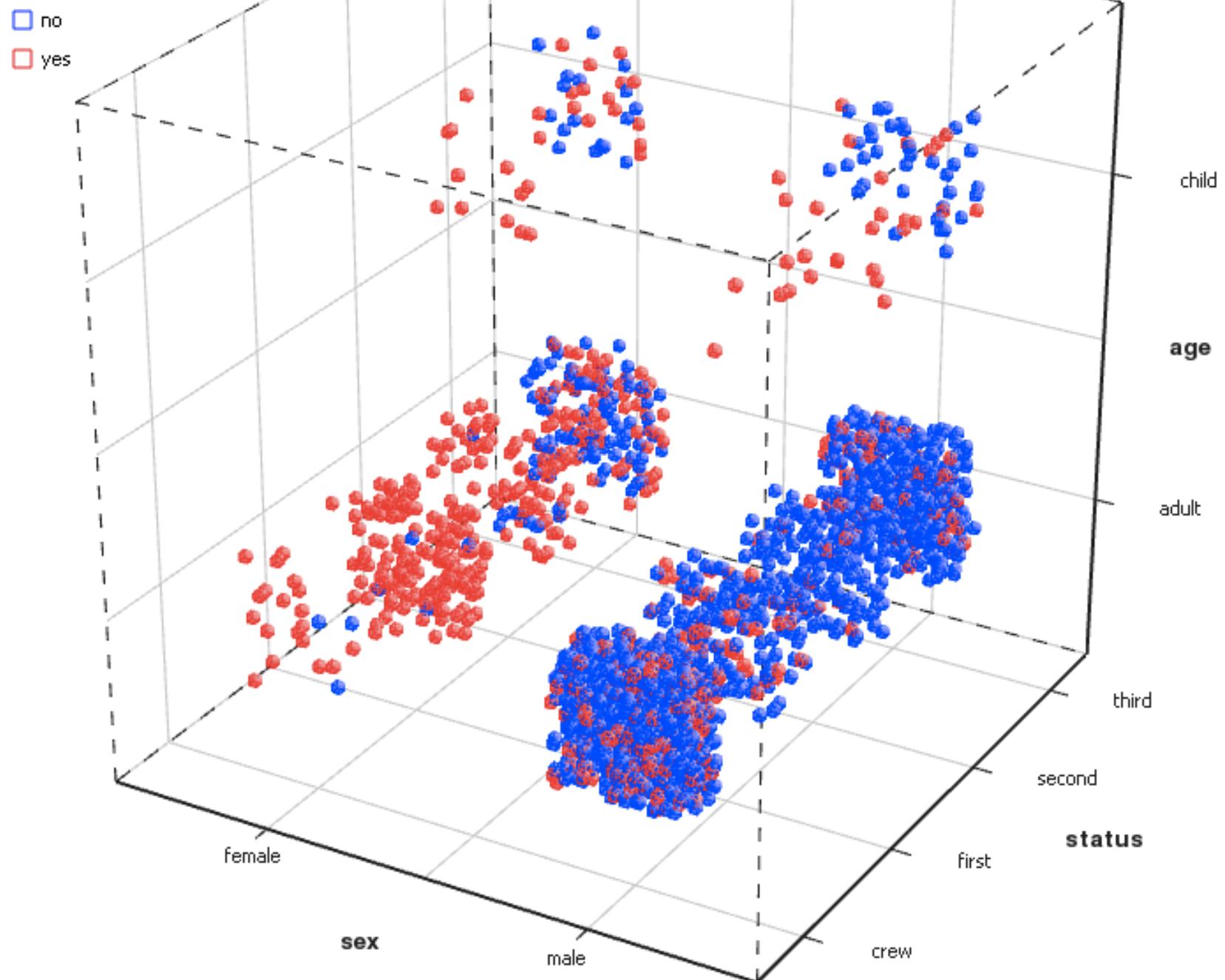




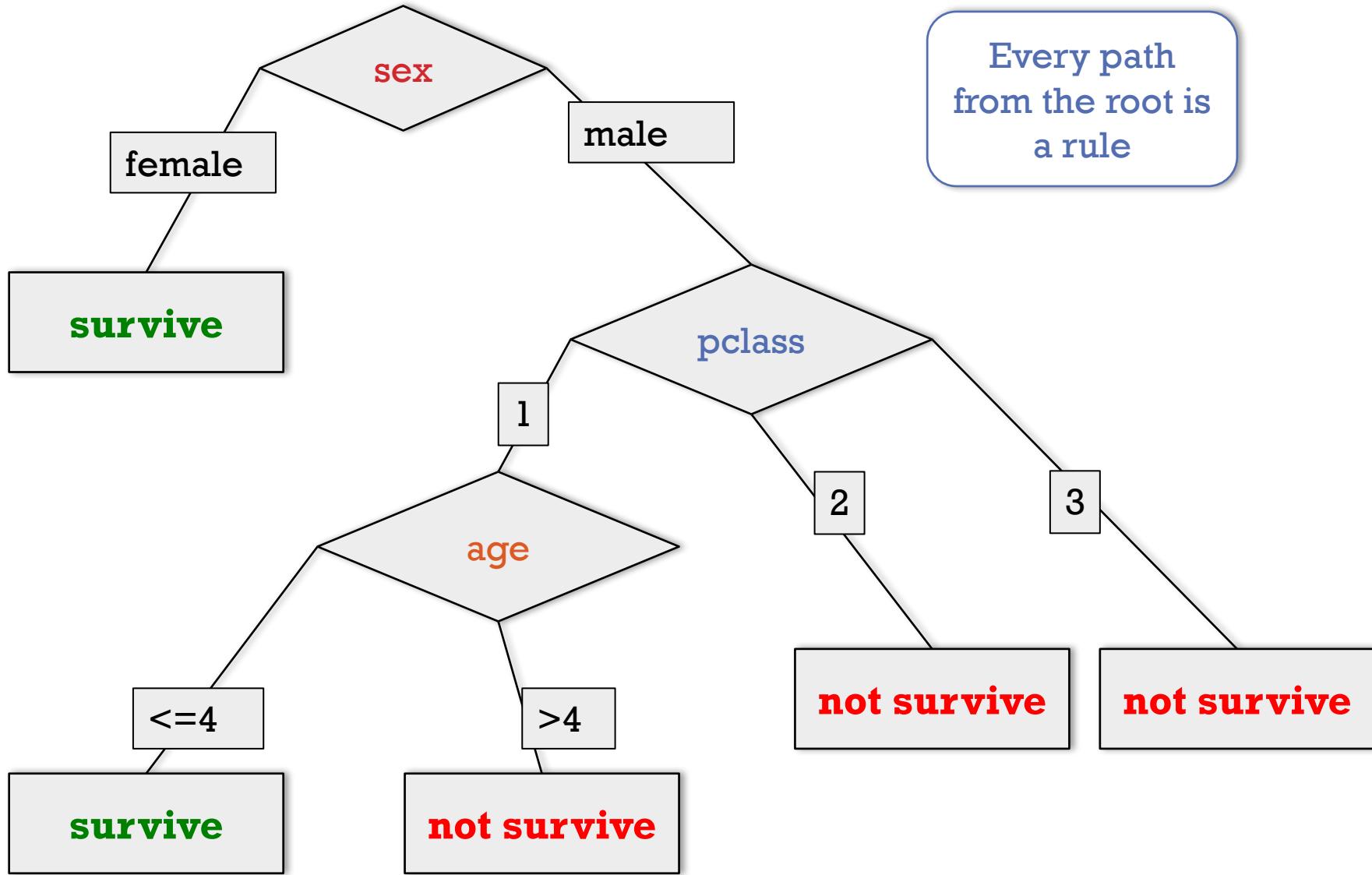




**survived:**



# DECISION TREE



# WHAT IS A CLASSIFIER

Apply a prediction function to a feature representation of an image/data-set to get the desired output:

$$f(\text{apple}) = \text{"apple"}$$

$$f(\text{tomato}) = \text{"tomato"}$$

$$f(\text{cow}) = \text{"cow"}$$

# THE MACHINE LEARNING FRAMEWORK

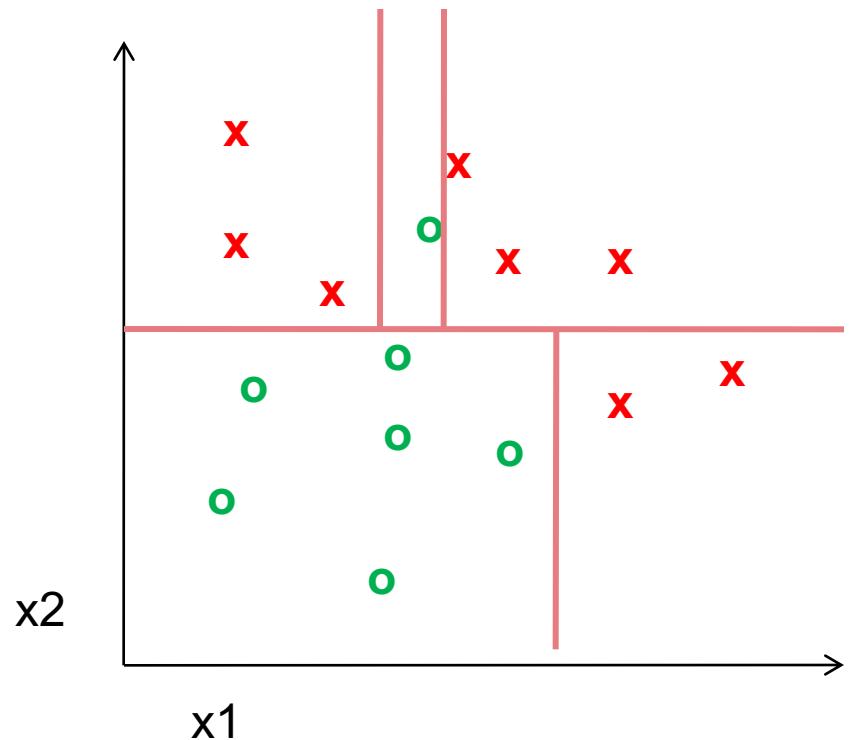
$$y = f(x)$$

output      prediction function      features

Training: given a *training set* of labeled examples  $\{(x_1, y_1), \dots, (x_N, y_N)\}$ , estimate the prediction function  $f$  by minimizing the prediction error on the training set

Testing: apply  $f$  to a never before seen *test example*  $x$  and output the predicted value  $y = f(x)$

# DECISION BOUNDARIES: DECISION TREES



# CLASSIFIER OVERVIEW

| Representation   | Evaluation   | Optimization   |
|--|--|--|
| Instances<br><i>K</i> -nearest neighbor<br>Support vector machines | Accuracy/Error rate<br>Precision and recall<br>Squared error<br>Likelihood | Combinatorial optimization<br>Greedy search<br>Beam search<br>Branch-and-bound |
| Hyperplanes<br>Naive Bayes<br>Logistic regression                  | Posterior probability<br>Information gain                                  | Continuous optimization<br>Unconstrained                                       |
| Decision trees   | K-L divergence   | Gradient descent   |
| Sets of rules<br>Propositional rules<br>Logic programs             | Cost/Utility<br>Margin   | Conjugate gradient<br>Quasi-Newton methods                                     |
| Neural networks  |  | Constrained  |
| Graphical models<br>Bayesian networks<br>Conditional random fields |  | Linear programming<br>Quadratic programming                                    |

# RECAP: DECISION TREES

## Representation

- A set of rules: IF...THEN conditions

## Evaluation

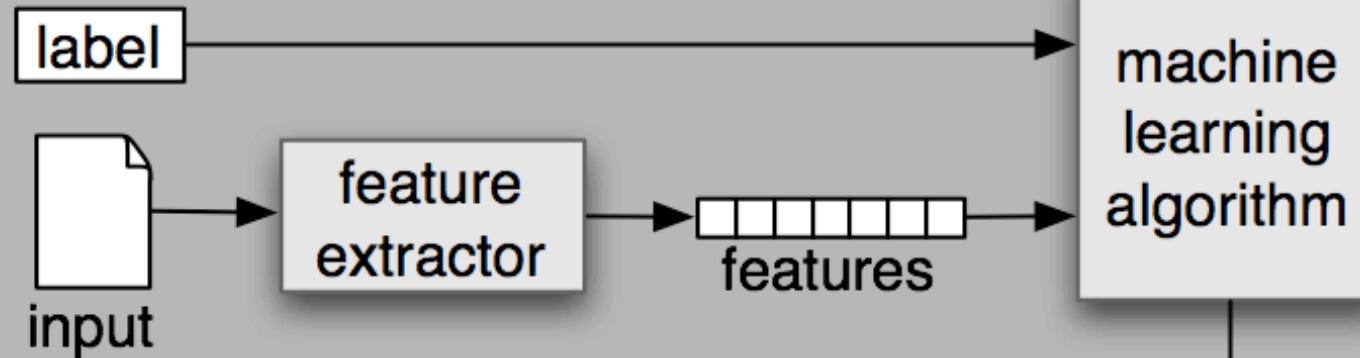
- coverage: # of data points that satisfy conditions
- accuracy = # of correct predictions / coverage

## Optimization

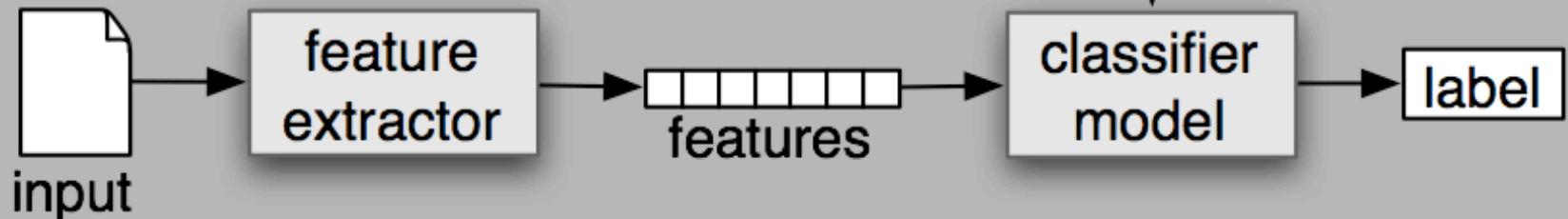
- Build decision tree that maximize accuracy

# ML PIPELINE (SUPERVISED)

## (a) Training



## (b) Prediction



# FEATURES

| Fact Table |                    |
|------------|--------------------|
| -          | <u>Shop ID</u>     |
| -          | <u>Customer ID</u> |
| -          | <u>Date ID</u>     |
| -          | <u>Product ID</u>  |
| -          | Amount             |
| -          | Volume             |
| -          | Profit             |
| -          | ...                |

| Fact Table |                    |
|------------|--------------------|
| -          | <u>Shop ID</u>     |
| -          | <u>Customer ID</u> |
| -          | <u>Date ID</u>     |
| -          | <u>Product ID</u>  |
| -          | Amount             |
| -          | Volume             |
| -          | Profit             |
| -          | Delivery Time      |
| -          | ...                |

| Product |                   |
|---------|-------------------|
| -       | <u>Product ID</u> |
| -       | Type_ID           |
| -       | Brand_ID          |
| -       | Length            |
| -       | Height            |
| -       | Depth             |
| -       | Weight            |
| -       | ...               |

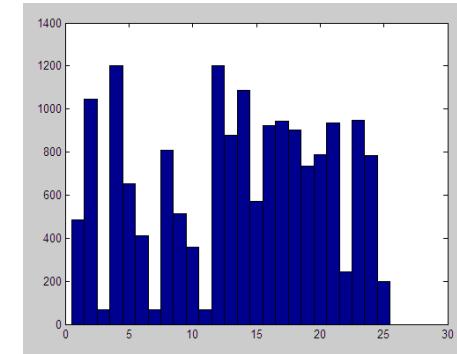
| Product_Type |                |
|--------------|----------------|
| -            | <u>Type ID</u> |
| -            | Name           |
| -            | Description    |
| -            | ...            |

| Brand |                 |
|-------|-----------------|
| -     | <u>Brand ID</u> |
| -     | Name            |
| -     | ...             |

| Customer State | Product Type | Product Weight | Volume (L*H*D) | Month | Delivery Time |
|----------------|--------------|----------------|----------------|-------|---------------|
|                |              |                |                |       |               |
|                |              |                |                |       |               |

# IMAGE FEATURES

## Raw pixels

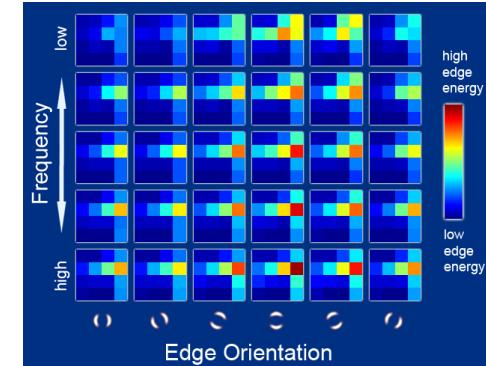


## Histograms

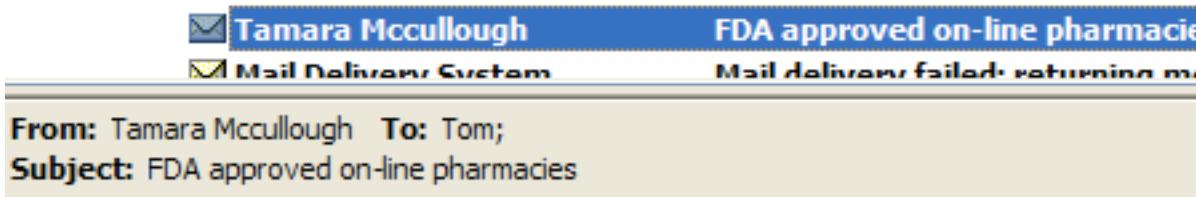
## Geometric features – edge

Corner,...

...



# TEXT FEATURES



**FDA approved on-line pharmacies.**  
Chose your product and site below:

[\*\*Canadian pharmacy\*\*](#) - Cialis Soft Tabs - \$5.78, **Viagra Professional** - \$1.38, Human Growth Hormone - \$43.37, Meridia - \$3.32, Tramadol - \$1.38

[\*\*HerbalKing\*\*](#) - Herbal pills for **Hair enlargement**. Techniques, products, dangerous pumps, exercises and surgeries.

[\*\*Anatrim\*\*](#) - Are you ready for Summer? Use **Anatrim**, the most powerful steroid ever developed.

Spam

Not Spam

## Bag of Words

*Viagra*  
*Soft*  
*Herbel*  
*Pills*  
*Are*  
...

## N-Grams

*herbel pills*  
*pills for*  
*for Hair*  
*Hair enlargement*  
*enlargement Techniques*  
...