

---

# Applied Machine Learning

## ML 101 & Course Introduction



James G. Shanahan <sup>1,2</sup>

<sup>1</sup>Church and Duncan Group,

<sup>2</sup>*School of Informatics, Computing and Engineering, Indiana University*

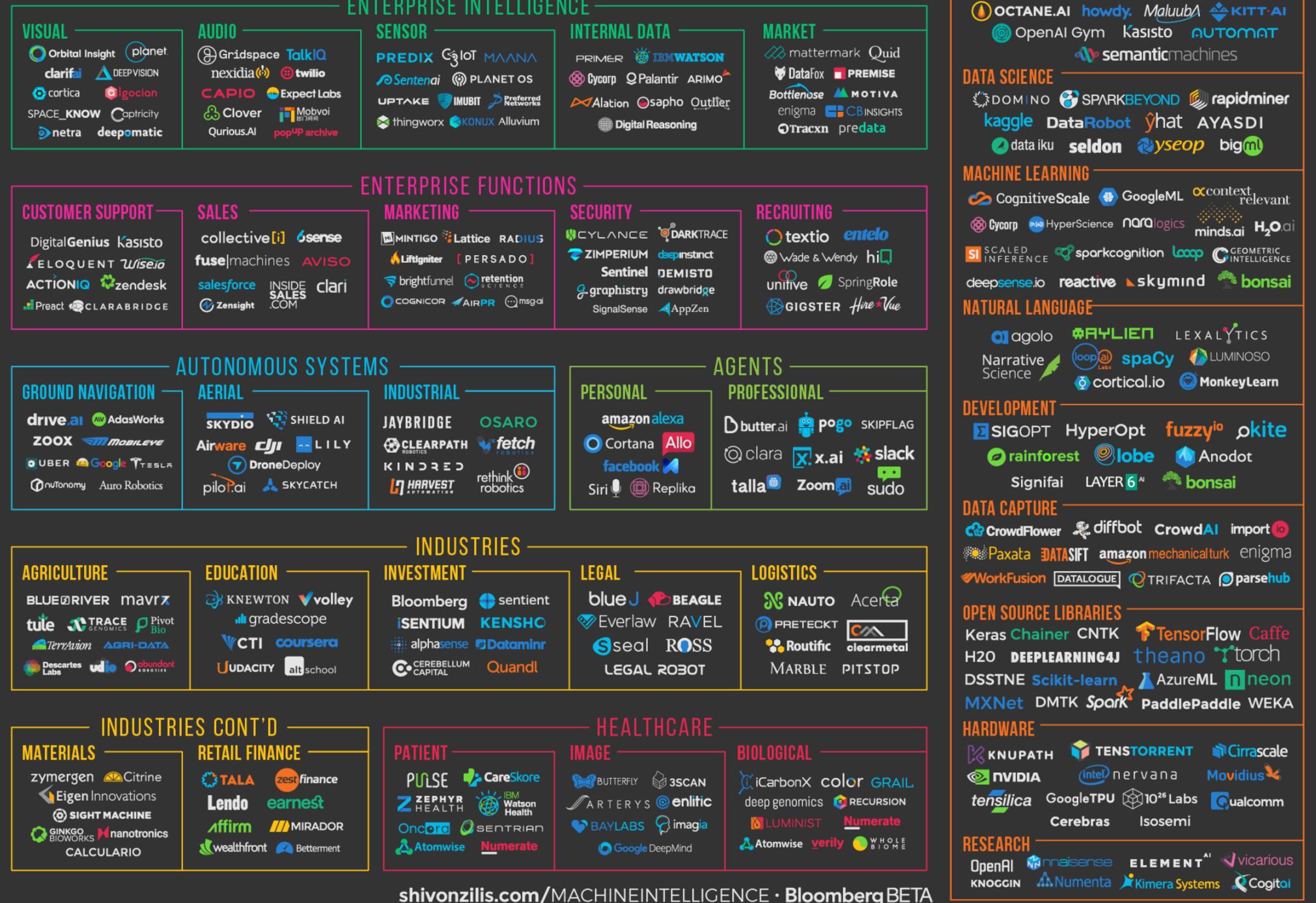
***EMAIL: James\_DOT\_Shanahan\_AT\_gmail\_DOT\_com***

# Outline

---

- **AI/ML 101:**
  - Introduction
  - Linear Regression
  - Beyond linear regression
- **Top AI market trends** to watch in 2017 and beyond
  - Key technical developments
  - Case studies
- **ML investor or entrepreneur**
- **Short Case Study in a Python Notebook**
- **Conclusions**
- **Course Logistics**

# MACHINE INTELLIGENCE 3.0



# AI Revenue and M&As

---

- **AI Software: from \$1.38 billion in 2016 to \$59.75 billion by 2025 [Tractica.com 2017]**
  - Autonomous vehicles, Consumer/social/image, Stock trading, Healthcare, HR, patent, security, digitizing paper
  - 9-years with a compound annual growth rate (CAGR) of 52%.
- **AI Services: \$120 Billion by 2025**
  - training, integration, Apps
- **AI Hardware: \$120 Billion by 2025**
  - GPUs, CPUs, Customer
- **250 AI acquisitions 2012-2017**
  - Over 250 private companies using AI algorithms across different verticals have been acquired 2012-2017(Q1), with 37 acquisitions taking place in Q1 2017 alone. Google acq 12 companies.

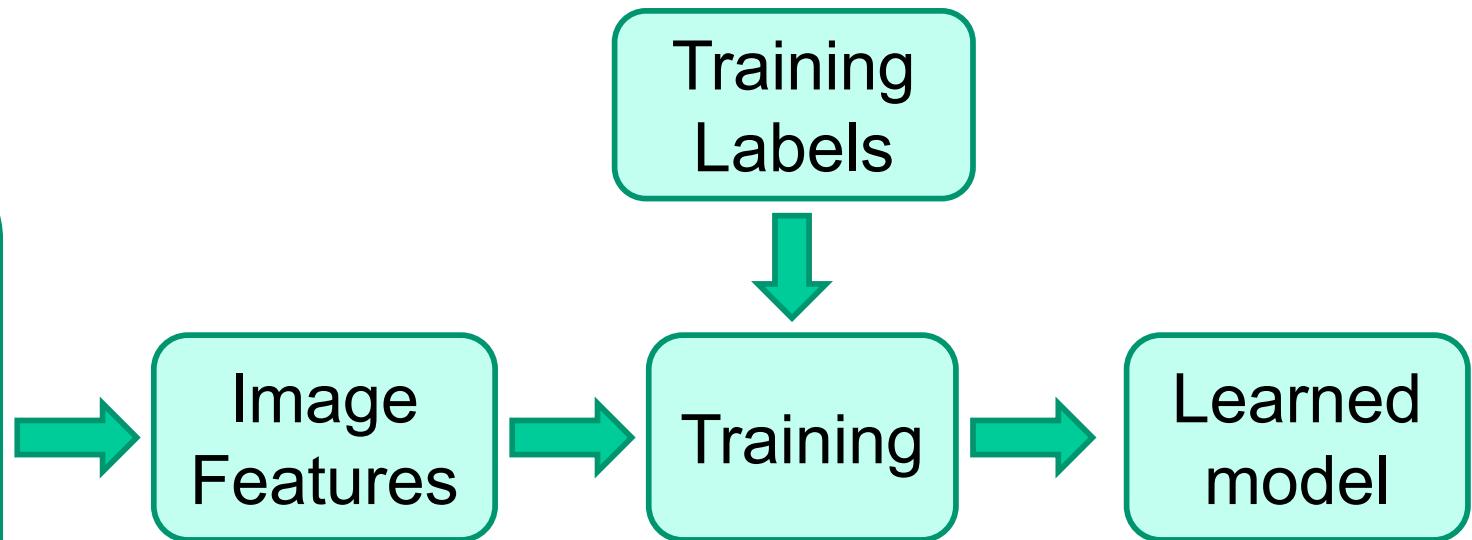
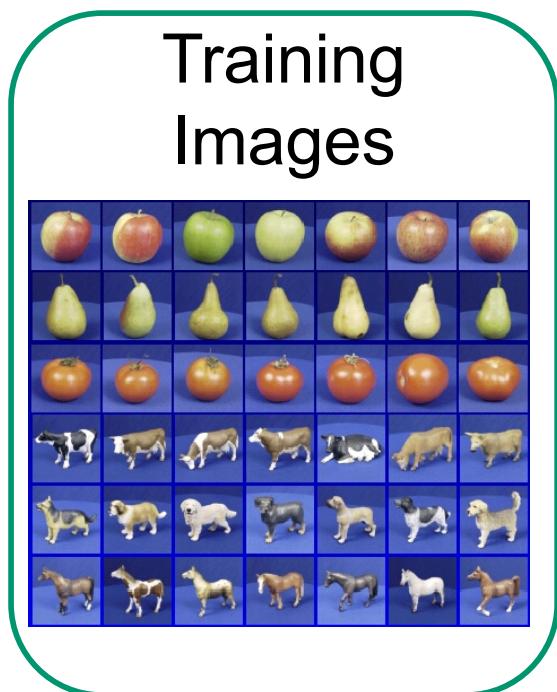
# Outline

---

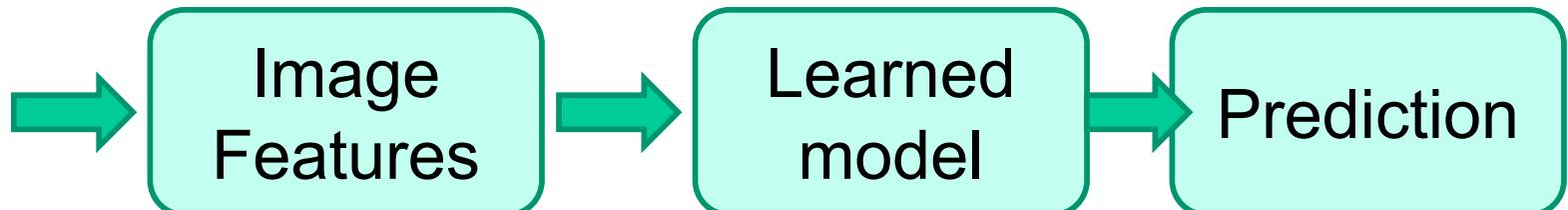
- **AI/ML 101:**
  - Introduction
  - Linear Regression
  - Beyond linear regression
- **Top AI market trends** to watch in 2017 and beyond
  - Key technical developments
  - Case studies
- **ML investor or entrepreneur**
- **Short Case Study in a Python Notebook**
- **Conclusions**
- **Course Logistics**

# Steps

## Training

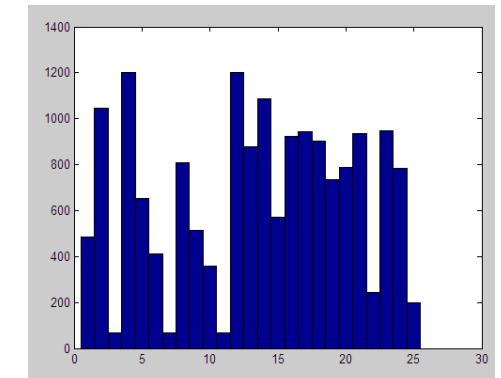


## Testing



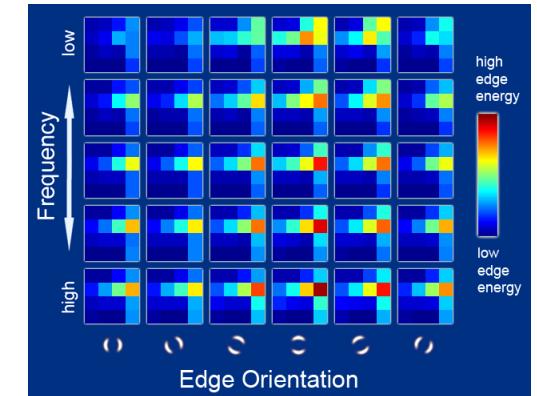
# Features

- Raw pixels



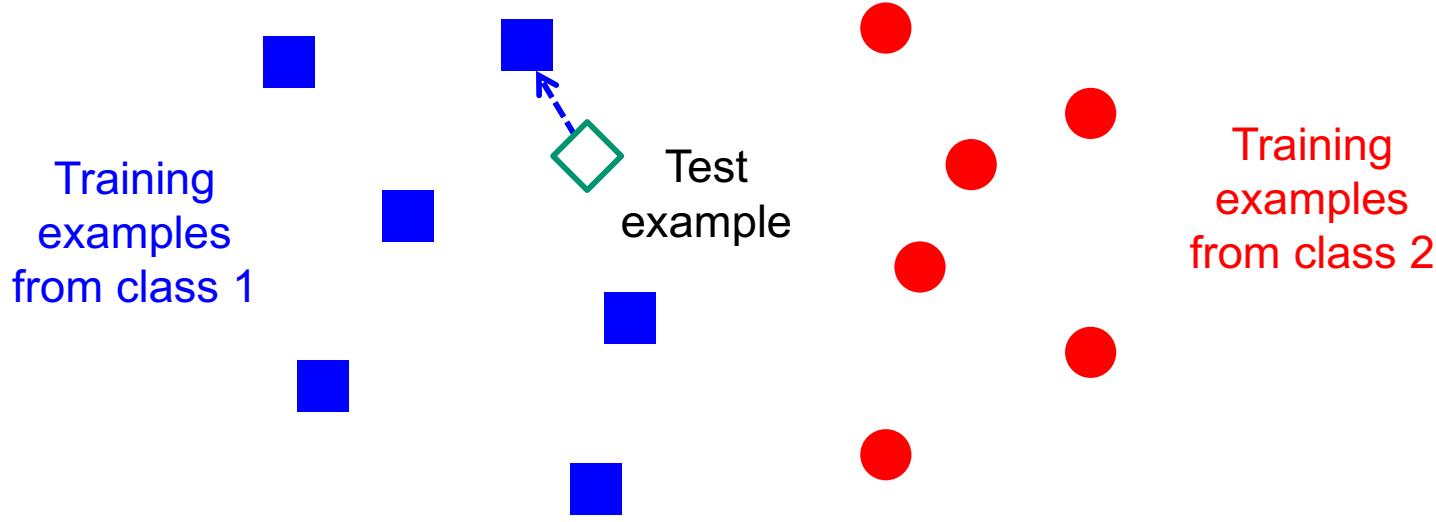
- Histograms

- GIST descriptors



- ...

# Classifiers: Nearest neighbor

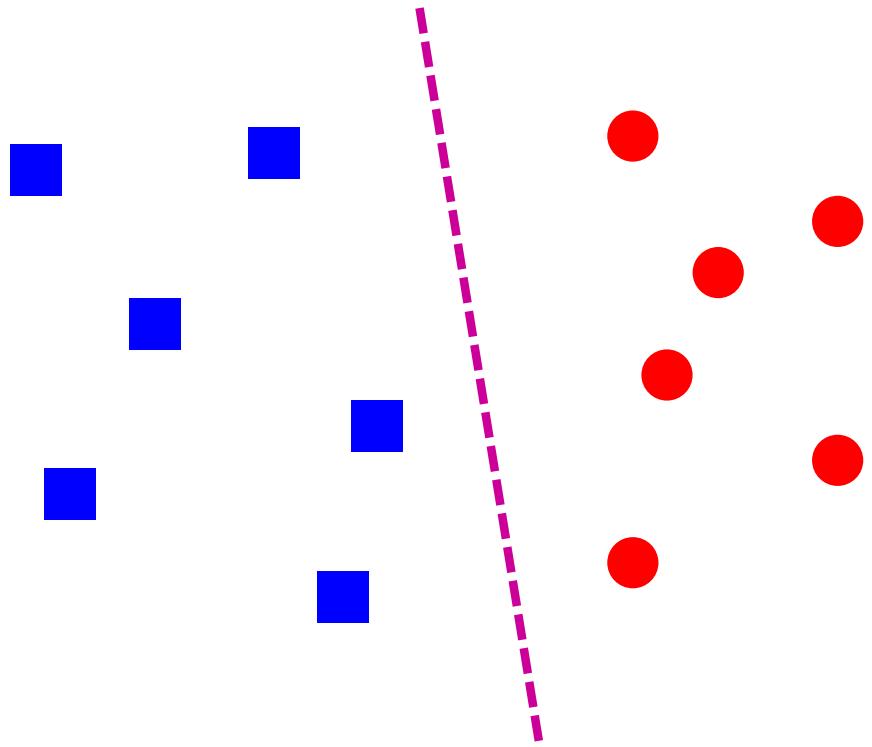


$f(x)$  = label of the training example nearest to  $x$

- All we need is a distance function for our inputs
- No training required!

# Classifiers: Linear

---



- Find a *linear function* to separate the classes:

$$f(x) = \text{sgn}(w \cdot x + b)$$

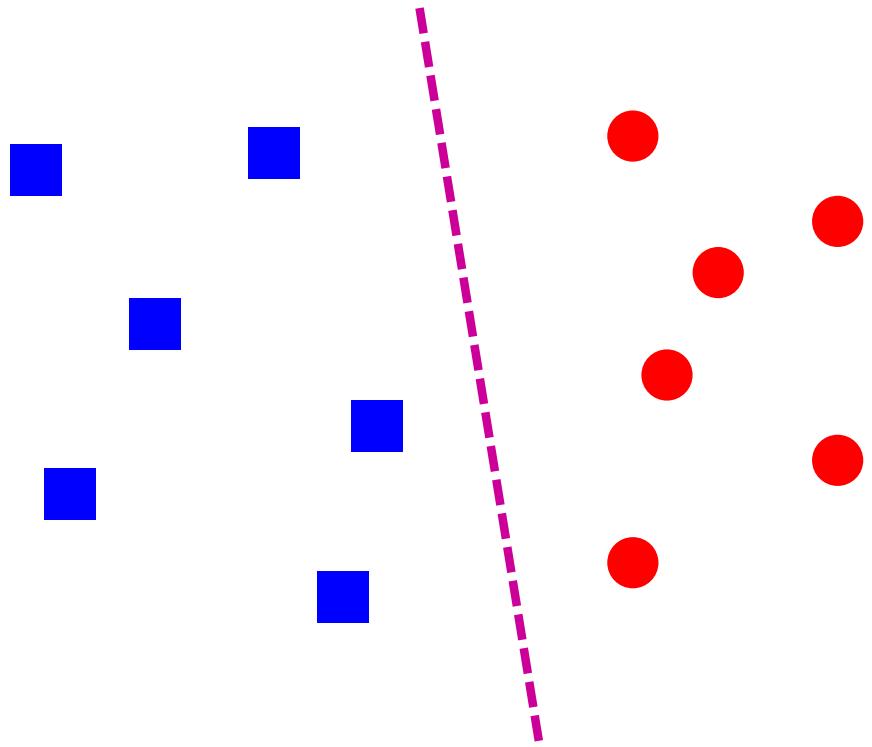
# Many classifiers to choose from

---

- **SVM**
  - **Neural networks**
  - **Naïve Bayes**
  - **Bayesian network**
  - **Logistic regression**
  - **Randomized Forests**
  - **Boosted Decision Trees**
  - **K-nearest neighbor**
  - **RBM**s
  - **Etc.**
- Which is the best one?

# Classifiers: Linear

---



- Find a *linear function* to separate the classes:

$$f(x) = \text{sgn}(w \cdot x + b)$$

# Comparison

assuming  $x$  in  $\{0, 1\}$

	Learning Objective	Training	Inference
Naïve Bayes	$\text{maximize} \sum_i \left[ \sum_j \log P(x_{ij}   y_i; \theta_j) + \log P(y_i; \theta_0) \right]$	$\theta_{kj} = \frac{\sum_i \delta(x_{ij} = 1 \wedge y_i = k) + r}{\sum_i \delta(y_i = k) + Kr}$	$\theta_1^T \mathbf{x} + \theta_0^T (1 - \mathbf{x}) > 0$ where $\theta_{1j} = \log \frac{P(x_j = 1   y = 1)}{P(x_j = 1   y = 0)}$ , $\theta_{0j} = \log \frac{P(x_j = 0   y = 1)}{P(x_j = 0   y = 0)}$ 
Logistic Regression	$\text{maximize} \sum_i \log(P(y_i   \mathbf{x}, \boldsymbol{\theta})) + \lambda \ \boldsymbol{\theta}\ $ where $P(y_i   \mathbf{x}, \boldsymbol{\theta}) = 1 / (1 + \exp(-y_i \boldsymbol{\theta}^T \mathbf{x}))$	Gradient ascent	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Linear SVM	$\text{minimize} \lambda \sum_i \xi_i + \frac{1}{2} \ \boldsymbol{\theta}\ ^2$ such that $y_i \boldsymbol{\theta}^T \mathbf{x} \geq 1 - \xi_i \quad \forall i$	Linear programming	$\boldsymbol{\theta}^T \mathbf{x} > 0$
Kernelized SVM	complicated to write	Quadratic programming	$\sum_i y_i \alpha_i K(\hat{\mathbf{x}}_i, \mathbf{x}) > 0$
Nearest Neighbor	most similar features $\rightarrow$ same label	Record data	$y_i$ where $i = \operatorname{argmin}_i K(\hat{\mathbf{x}}_i, \mathbf{x})$

# Titanic: predict passenger survival



- **Task: Predict whether a passenger will survive or not**
  - you must predict the fate of the passengers aboard the RMS Titanic, which famously sank in the Atlantic ocean during its maiden voyage from the UK to New York City after colliding with an iceberg.

# Titanic: Python and R

- In python:
  - <http://nbviewer.jupyter.org/urls/dl.dropbox.com/s/zshj6zmo8q6ej85/Titanic-EDA-LogisticRegression.ipynb>
- In R
  - <http://trevorstephens.com/kaggle-titanic-tutorial/getting-started-with-r/>

## Titanic: Getting Started With R

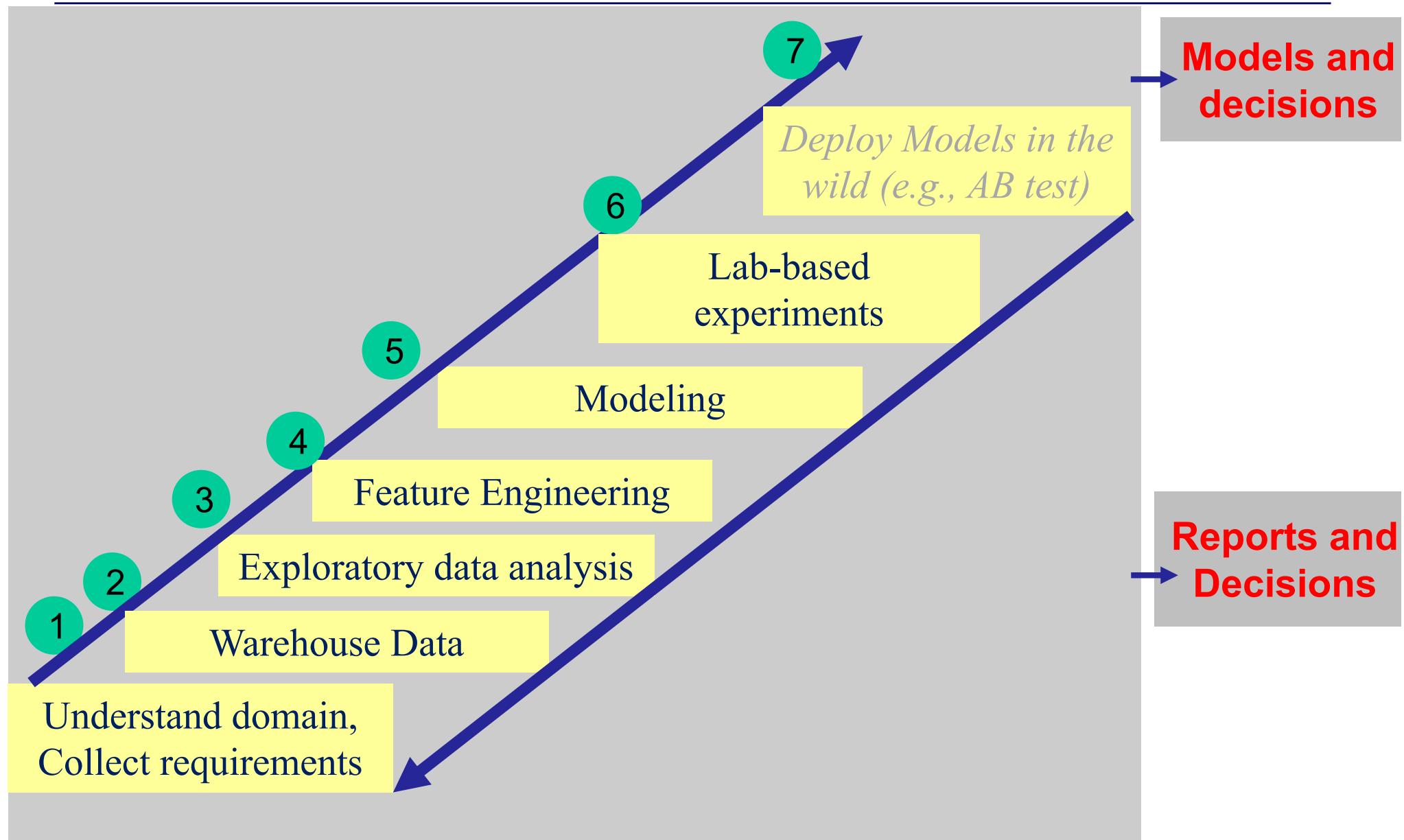
⌚ 3 minutes read

So you're excited to get into prediction and like the look of Kaggle's excellent getting started competition, [Titanic: Machine Learning from Disaster](#)? Great! It's a wonderful entry-point to machine learning with a manageably small but very interesting dataset with easily understood variables.

In this competition, you must predict the fate of the passengers aboard the RMS Titanic, which famously sank in the Atlantic ocean during its maiden voyage from the UK to New York City after colliding with an iceberg.



# Typical Abstract Data Analytics Pipeline



# Python Notebook

Jupyter Titanic-EDA-LogisticRegression (autosaved) ?

File Edit View Insert Cell Kernel Navigate Widgets Help

Markdown

Contents [-] ↻ ↺

- 0.1 Kaggle Competition | Titanic Machine Learning from Disaster
- 0.2 Goal for this Notebook:
  - 0.2.1 This Notebook will show basic examples of:
  - 0.2.2 Data Handling
  - 0.2.3 Data Analysis
  - 0.2.4 Valuation of the Analysis
  - 0.2.5 Required Libraries:
- 0.3 Data Handling
  - 0.3.1 Let's read our data in using pandas:
  - 0.4 Let's take a look:
  - 0.5 Take care of missing values:
  - 0.6 Let's take a Look at our data graphically:
  - 0.7 Exploratory Visualization:
  - 0.8 Now let's tease more structure out of the data,
- 0.9 Let's break the previous graph down by gender
  - 0.9.1 Great! But let's go down even further:
- 1 Supervised Machine Learning
  - 1.0.1 Logistic Regression:
  - 1.0.2 The skinny, as explained by yours truly:
- 2 So how well did this work?
  - 2.1 Now lets use our model to predict the test set values and then save t
  - 2.2 Read the test data
  - 2.3 Examine our dataframe
  - 2.4 Add our independent variable to our test data. (It is usually left blank
  - 2.5 Results as scored by Kaggle: RMSE = 0.77033 That result is pretty g
  - 2.6 Support Vector Machine (SVM)
- 3 From me
  - 3.1 Random Forest
    - 3.1.1 Follow me on github, and twitter for more books to come soon!

## 0.1 Kaggle Competition | Titanic Machine Learning from Disaster

The sinking of the RMS Titanic is one of the most infamous shipwrecks in history after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. The sinking led to better safety regulations for ships.

One of the reasons that the shipwreck led to such loss of life was that in those days there was some element of luck involved in surviving the sinking, so that the women, children, and the upper-class.

In this contest, we ask you to complete the analysis of what sorts of passengers had the best chances of survival using machine learning to predict which passengers survived the tragedy.

This Kaggle Getting Started Competition provides an ideal starting point for learning machine learning."

From the competition [homepage](#).

## 0.2 Goal for this Notebook:

Show a simple example of an analysis of the Titanic disaster in Python using Jupyter Notebook. This notebook is intended for those who are new to the field or those who are already in the field and looking to see an example of a machine learning project.

### 0.2.1 This Notebook will show basic examples of:

- Importing Data with Pandas

### 0.2.2 Data Handling

- Importing Data with Pandas

# Python Notebook

---

## Goal for this Notebook:

Show a simple example of an analysis of the Titanic disaster in Python using a full complement of PyData utilities. This is aimed for those looking to get into the field or those who are already in the field and looking to see an example of an analysis done with Python.

## This Notebook will show basic examples of:

### Data Handling

- Importing Data with Pandas
- Cleaning Data
- Exploring Data through Visualizations with Matplotlib

### Data Analysis

- Supervised Machine learning Techniques:
  - Logit Regression Model
  - Plotting results
  - Support Vector Machine (SVM) using 3 kernels
  - Basic Random Forest
  - Plotting results

### Valuation of the Analysis

- K-folds cross validation to validate results locally
- Output the results from the IPython Notebook to Kaggle

# Read in the data as a dataframe

## 0.3 Data Handling

### 0.3.1 Let's read our data in using pandas:

```
[1]: df = pd.read_csv("data/train.csv")
```

Show an overview of our data:

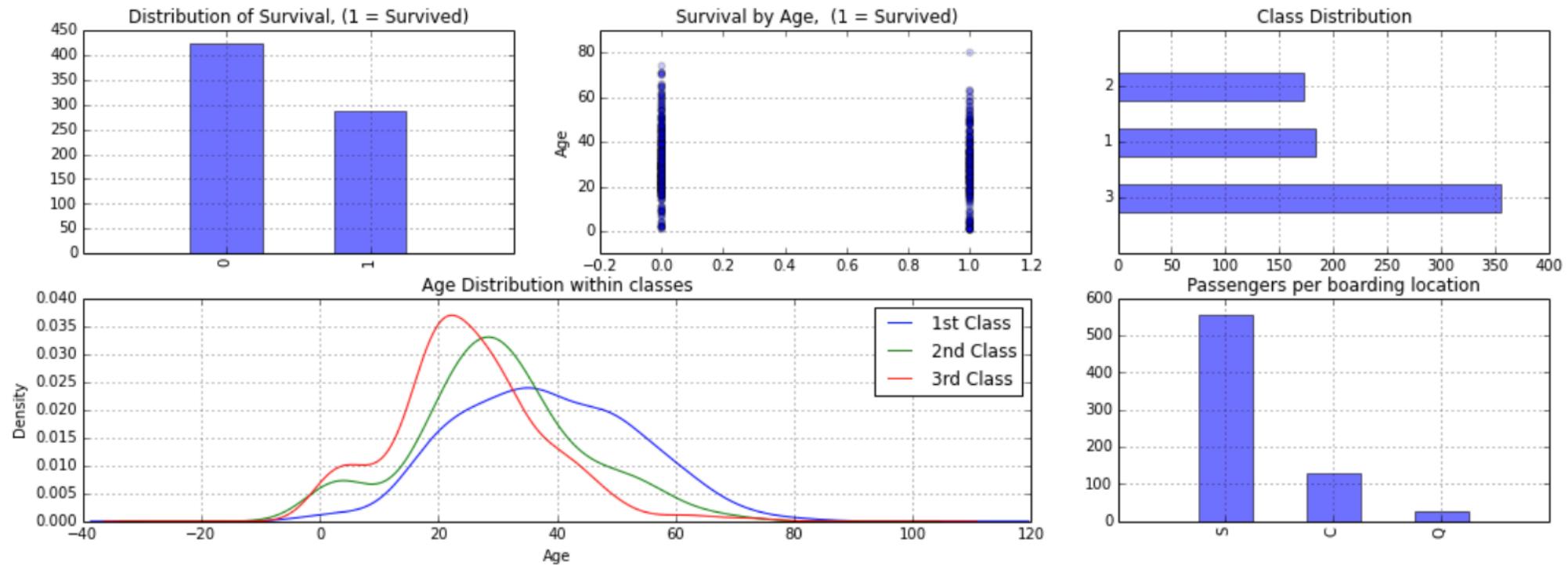
```
[2]: df
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Th...	female	38	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35	0	0	373450	8.0500	NaN	S
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	NaN	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54	0	0	17463	51.8625	E46	S
7	8	0	3	Palsson, Master. Gosta Leonard	male	2	3	1	349909	21.0750	NaN	S

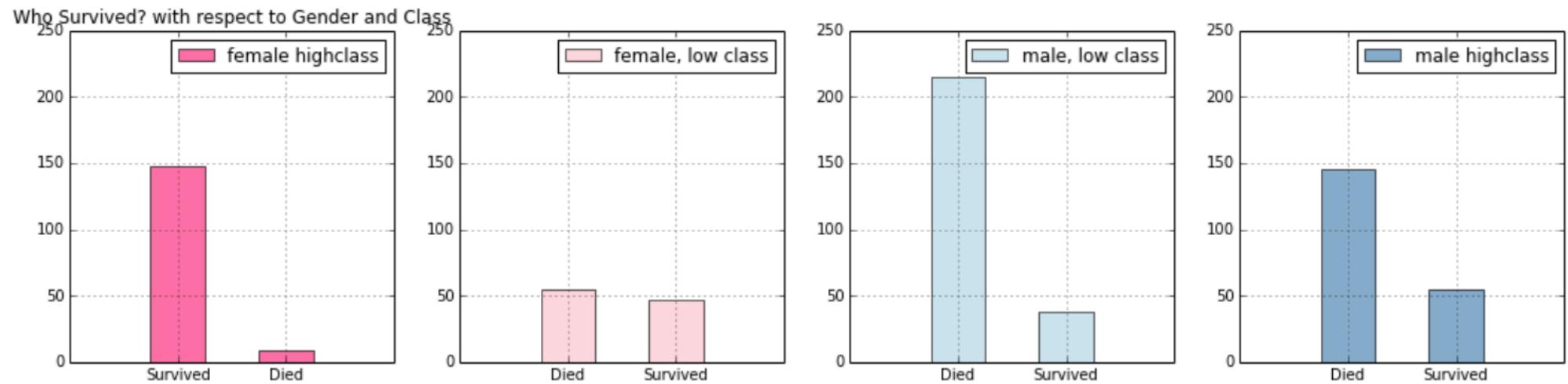
# Exploratory Data Analysis

```
# specifies the parameters of our graphs  
plt.title("Passengers per boarding location")
```

```
<matplotlib.text.Text at 0x119c1dd90>
```



# Exploratory Data Analysis



# Model passenger survival Classification task using logistic regression

```
import statsmodels.api as sm
```

```
1 # model formula
2 # here the ~ sign is an = sign, and the features of our dataset
3 # are written as a formula to predict survived. The C() lets our
4 # regression know that those variables are categorical.
5 # Ref: http://patsy.readthedocs.org/en/latest/formulas.html
6 formula = 'Survived ~ C(Pclass) + C(Sex) + Age + SibSp + C(Embarked)'
7 # create a results dictionary to hold our regression results for easy analysis later
8 results = {}
```

```
1 # create a regression friendly dataframe using patsy's dmatrices function
2 y,x = dmatrices(formula, data=df, return_type='dataframe')
3
4 # instantiate our model
5 model = sm.Logit(y,x)
6
7 # fit our model to the training data
8 res = model.fit()
9
10 # save the result for outputting predictions later
11 results['Logit'] = [res, formula]
12 res.summary()
```

# Ensemble of trees

---

```
|: # import the machine learning library that holds the randomforest
|: import sklearn.ensemble as ske
|
|: # Create the random forest model and fit the model to our training data
|: y, x = dmatrices(formula_ml, data=df, return_type='dataframe')
|: # RandomForestClassifier expects a 1 demensional NumPy array, so we convert
|: y = np.asarray(y).ravel()
|: #instantiate and fit our model
|: results_rf = ske.RandomForestClassifier(n_estimators=100).fit(x, y)
|
|: # Score the results
|: score = results_rf.score(x, y)
|: print "Mean accuracy of Random Forest Predictions on the data was: {0}".format(score)
```

# R Code for Titanic

---

<http://trevorstephens.com/kaggle-titanic-tutorial/getting-started-with-r/>

- All code is available on my [Github repository](#).
- A series of tutorials in five parts is available here:
  - [Part 1: Booting Up R](#)
  - [Part 2: The Gender-Class Model](#)
  - [Part 3: Decision Trees](#)
  - [Part 4: Feature Engineering](#)
  - [Part 5: Random Forests](#)

# New features: Title

- **Title feature (extracted from Name)**

```
> combi$Name <- as.character(combi$Name)
> combi$Name[1]
[1] "Braund, Mr. Owen Harris"
```

```
> table(combi$title)
```

Capt	Col	Don	Dona	Dr	Jonkheer	Lady
1	4	1	1	8	1	1
Major	Master	Miss	Mlle	Mme	Mr	Mrs
2	61	260	2	1	757	197
Ms	Rev	Sir	the Countess			
2	8	1	1			

# Reduce titles to 2 groups

For the ladies, we have Dona, Lady, Jonkheer (\*see comments below), and of course our Countess. All of these are again the rich folks, and may have acted somewhat similarly due to their noble birth. Let's combine these two groups and reduce the number of factor levels to something that a decision tree might make sense of:

```
> combi$Title[combi$Title %in% c('Capt', 'Don', 'Major', 'Sir')] <- 'Sir'  
< ncombi$Title[combi$Title %in% c('Dona', 'Lady', 'the Countess', 'Jonkheer')] <- 'Lady'
```

Our final step is to change the variable type back to a factor, as these are essentially categories that we have created:

```
> combi$Title <- factor(combi$Title)
```

# Family variable

```
> table(combi$FamilyID)
```

11Sage	3Abbott	3Appleton	3Beckwith	3Boulos
11	3	1	2	3
3Bourke	3Brown	3Caldwell	3Christy	3Collyer
3	4	3	2	3
3Compton	3Cornell	3Coutts	3Crosby	3Danbom
3	1	3	3	3 . . .

Hmm, a few seemed to have slipped through the cracks here. There's plenty of FamilyIDs with only one or two members, even though we wanted only family sizes of 3 or more. Perhaps some families had different last names, but whatever the case, all these one or two people groups is what we sought to avoid with the three person cut-off. Let's begin to clean this up:

```
> famIDs <- data.frame(table(combi$FamilyID))
```

Now we have stored the table above to a dataframe. Yep, you can store most tables to a dataframe if you want to, so let's take a look at it by clicking on it in the explorer:

The screenshot shows the RStudio interface with the 'Tutorial4.R' script open. In the environment browser, there is a data frame named 'famIDs'. A tooltip indicates it contains '97 observations of 2 variables'. The data frame has two columns: 'Var1' and 'Freq'. The first few rows are:

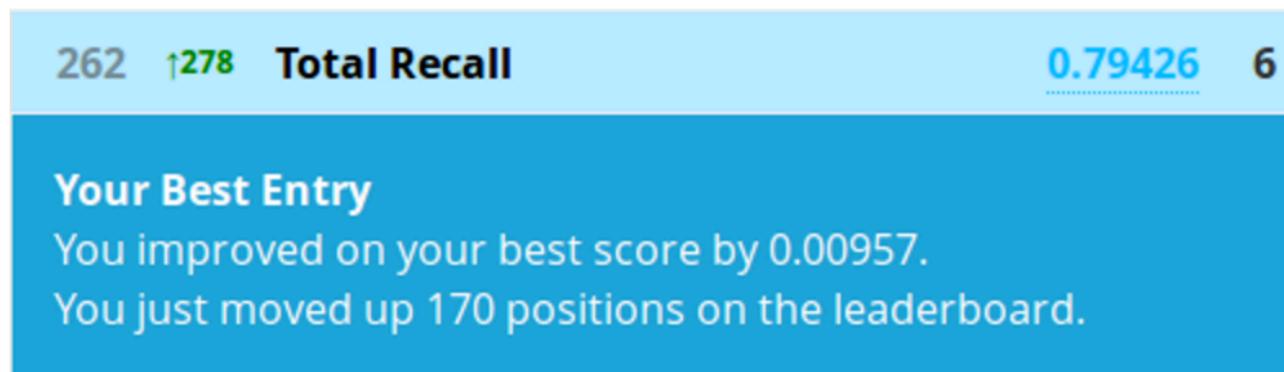
	Var1	Freq
1	11Sage	11
2	3Abbott	3
3	3Bourke	3
4	3Brown	4
5	3Caldwell	3
6	3Christy	2
7	3Coutts	3
8	3Crosby	3
9	3Danbom	3
10	3 . . .	3 . . .

# Improved score with new features + DT

---

- ..

But all that aside, you know should know how to create a submission from a decision tree, so let's see how it performed!



# Exercise

- Download Titanic sample code
- Do a submission and see where you come on the Kaggle dashboard

<https://www.kaggle.com/c/titanic>



Knowledge • 4,775 teams

## Titanic: Machine Learning from Disaster

Fri 28 Sep 2012

Sat 31 Dec 2016 (3 months to go)

Dashboard

- Home
- Data
- Make a submission

Information

- Description
- Evaluation
- Rules
- Prizes
- Frequently Asked Question...
- Getting Started With Excel
- Getting Started With Python...
- Getting Started With R...
- New: Getting Started with R
- Submission Instructions

Forum

Kernels

- New Script
- New Notebook

Leaderboard

Visualization

My Submissions

Competition Details » Get the Data » Make a submission

## Predict survival on the Titanic using Excel, Python, R & Random Forests

If you're new to data science and machine learning, or looking for a simple intro to the Kaggle competitions platform, this is the best place to start. Continue reading below the competition description to discover a number of tutorials, benchmark models, and more.

### Competition Description

The sinking of the RMS Titanic is one of the most infamous shipwrecks in history. On April 15, 1912, during her maiden voyage, the Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. This sensational tragedy shocked the international community and led to better safety regulations for ships.

One of the reasons that the shipwreck led to such loss of life was that there were not enough lifeboats for the passengers and crew. Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others, such as women, children, and the upper-class.

# Outline

---

- **AI/ML 101:**
  - Introduction
  - Linear Regression
  - Beyond linear regression
- **Top AI market trends** to watch in 2017 and beyond
  - Key technical developments
  - Case studies
- **ML investor or entrepreneur**
- **Short Case Study in a Python Notebook**
- **Conclusions**
- **Course Logistics**

# On The BART train in San Francisco



# Conclusions



## Machine learning will change how we work, rest and play

- 20<sup>th</sup> century life transformers: Electricity, automobile, Internet, mobile phones **Analytical → Predictive → Decisive**
- Executives are realizing that this new technology, AI, could change everything, but nobody knows exactly how or when
- Companies have at their disposal the full set of building blocks to begin embedding machine intelligence in their businesses.

## Education:

- Tools; teach machine learning using automatic differentiation

## Social impact:

- Creating a workforce that combines lower headcount with higher productivity
- Ethical concerns

## Huge opportunities

- AI: USD 6.5 billion (WorldWide) in 2016 → \$300B by 2025
- Corporate Arms Race for AI startups (250 acqs in last 5 years)

# Outline

---

- **AI/ML 101:**
  - Introduction
  - Linear Regression
  - Beyond linear regression
- **Top AI market trends** to watch in 2017 and beyond
  - Key technical developments
  - Case studies
- **ML investor or entrepreneur**
- **Short Case Study in a Python Notebook**
- **Conclusions**
- **Course Logistics**

# Course Schedule

---

## PART 1: Introduction

1. Machine Learning introduction and course overview
2. KNN: classification, regression + EDA + ML pipelines
3. Optimization theory: gradient descent

## PART 2: Supervised machine learning

4. Linear Regression
5. Probabilistic approaches to ML
6. Classification: logistic/softmax regression
7. **MID-TERM exam**
8. Perceptron and Support Vector Machines
9. Non-gradient-based approaches: decision trees, random forest

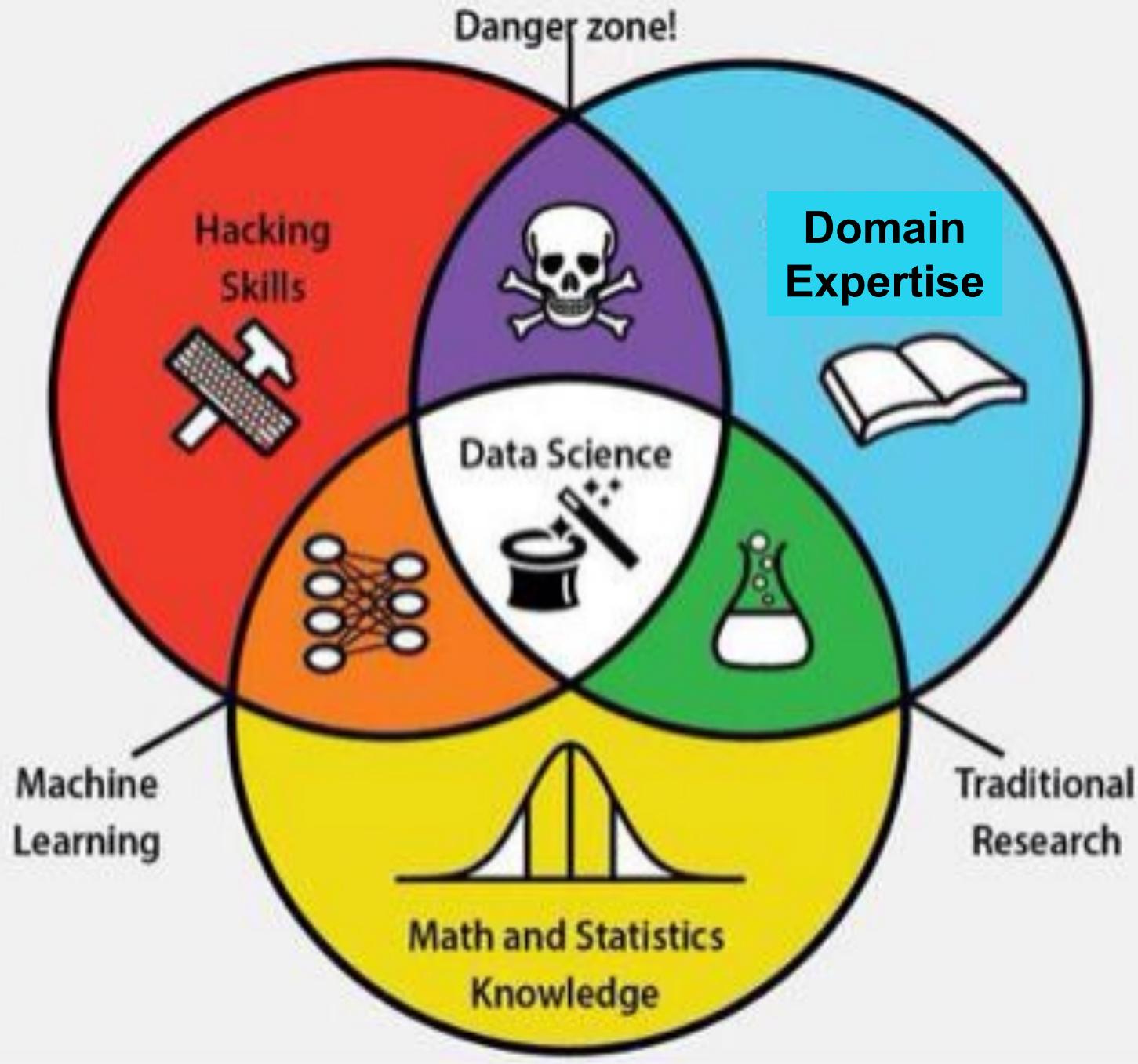
## PART 3: Unsupervised learning

10. Clustering: K-means, Hierarchical clustering, DBSCAN
11. Dimensionality reduction: PCA, tSNE, Word2Vec

## PART 4: Advanced/Recent topics

12. Neural networks: backprop, core concepts, deep learning
13. Recommender systems: popularity-based, item-to-item, matrix factorizations
14. Reinforcement learning: multi-armed bandit, Bellman equation, Q-learning
15. **FINAL Exam**

# DS Skillset



A venn diagram with a Danger Bearing

[Drew Conway]

# Implementing algorithms from scratch

---

- **Focus on theory, implementation and practice**
- **Avoid the “Danger Zone” (lots of great libraries but....)**
- **There are several different reasons why implementing algorithms from scratch can be useful:**
  - it can help us to understand the inner workings of an algorithm
  - Let's us figure how to parallelize
  - we can add new features to an algorithm or experiment with different variations of the core idea
  - we want to invent new algorithms or implement algorithms no one has implemented/shared yet
  - we are not satisfied with the API and/or we want to integrate it more "naturally" into an existing software library
  - we could try to implement an algorithm more efficiently
  - we circumvent licensing issues (e.g., Linux vs. Unix) or platform restrictions

---

<http://cs231n.github.io/python-numpy-tutorial/>

## Python Numpy Tutorial

```
distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
min_index = np.argmin(distances) # get the index with smallest distance
Ypred[i] = self.ytr[min_index] # predict the label of the nearest example
```

# Grading Scale

## Grading Scale

Activity	Score
Assignments	20%
Midterm exam	30%
Class participation	20%
End of semester exam	30%

Final Grade	Total Score
A "Excellent"	>= 85%
B "Good"	from 75% to 84%
C "Satisfactory"	from 60% to 74%
D "Poor"	from 50% to 59%
F "Unacceptable"	< 50%

# Weekly Schedule

---

- One Practical session: to be announced
- Office hours on demand: to be announced
- Homework assignments: one every 2 weeks
- Use the STAR methodology to present you Questions/Answers
  - Knowledge base questions
  - Problem solving

# Course Resources

---

- **LMS: Canvas**
- **Github**
- **Docker Container**

[Code](#)[Issues 0](#)[Pull requests 0](#)[Projects 0](#)[Wiki](#)[Insights](#)[Settings](#)

course materials. Autumn 2017

[Edit](#)[Add topics](#)[16 commits](#)[1 branch](#)[0 releases](#)[1 contributor](#)Branch: [master](#)[New pull request](#)[Create new file](#)[Upload files](#)[Find file](#)[Clone or download](#)

 vladimir-chernykh	master solution for KNN	Latest commit 06e614a 2 hours ago
 Assignments	master solution for KNN	2 hours ago
 Dockers	2nd update for new image	a day ago
 Exam	end of term exam placeholder	4 days ago
 Labs	2nd update for new image	a day ago
 Tutorials	minor syntax updates for newer modules vesrions	5 days ago
 .gitignore	initial commit	5 days ago
 README.md	initial commit	5 days ago
 dlcourse.sh	initial commit	5 days ago
 start.sh	initial commit	5 days ago
 stop.sh	initial commit	5 days ago

# This class will be demanding

---

But it is a high **R**ol class

- Plan on spending 10 hours +/-10 hours per week on this class

---



# End of lecture