

INTRO To DATA SCIENCE

REVIEW

ARUN AHUJA

MOUNT SINAI MEDICAL CENTER

- @ Icahn Institute for Genetics and Genomic Sciences
- Data Scientist @ Integral Ad Science
- Developer @ Morgan Stanley Electronic Trading and Real-Time Systems Group

INTRO To DATA SCIENCE

REVIEW

<i>supervised</i> <i>unsupervised</i>	<i>making predictions</i> <i>discovering patterns</i>
--	--

<i>supervised</i> <i>unsupervised</i>	<i>labeled examples</i> <i>no labeled examples</i>
--	---

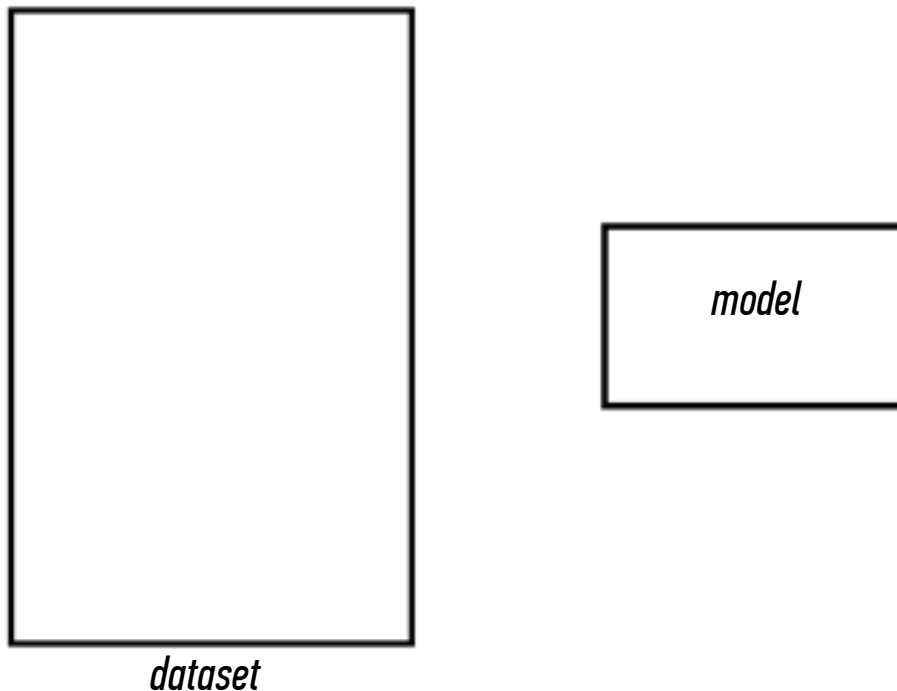
	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

INTRO TO DATA SCIENCE

SUPERVISED LEARNING

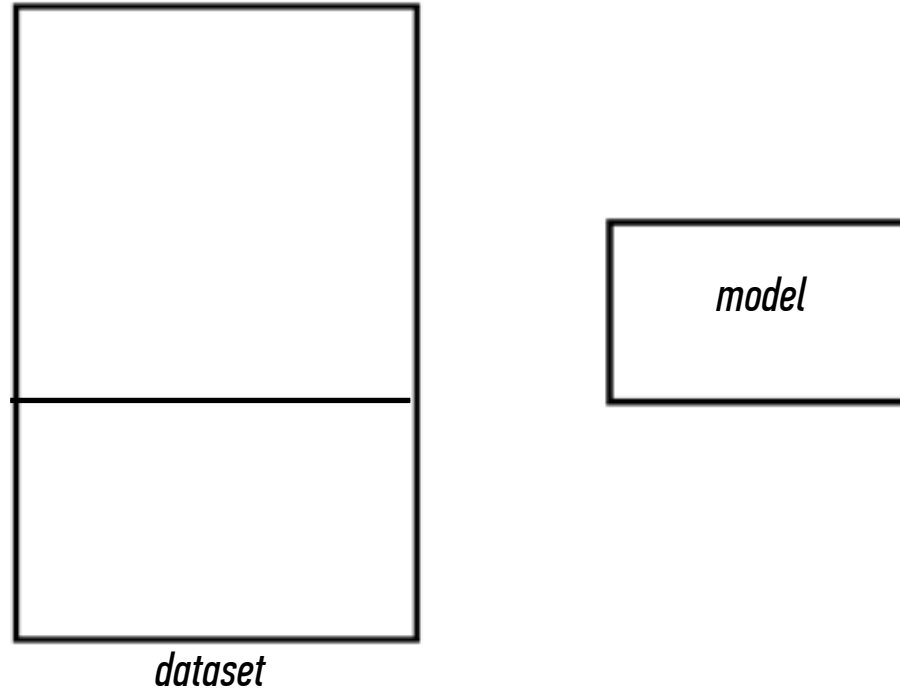
Wednesday, March 19, 14

Q: What steps does a classification problem require?



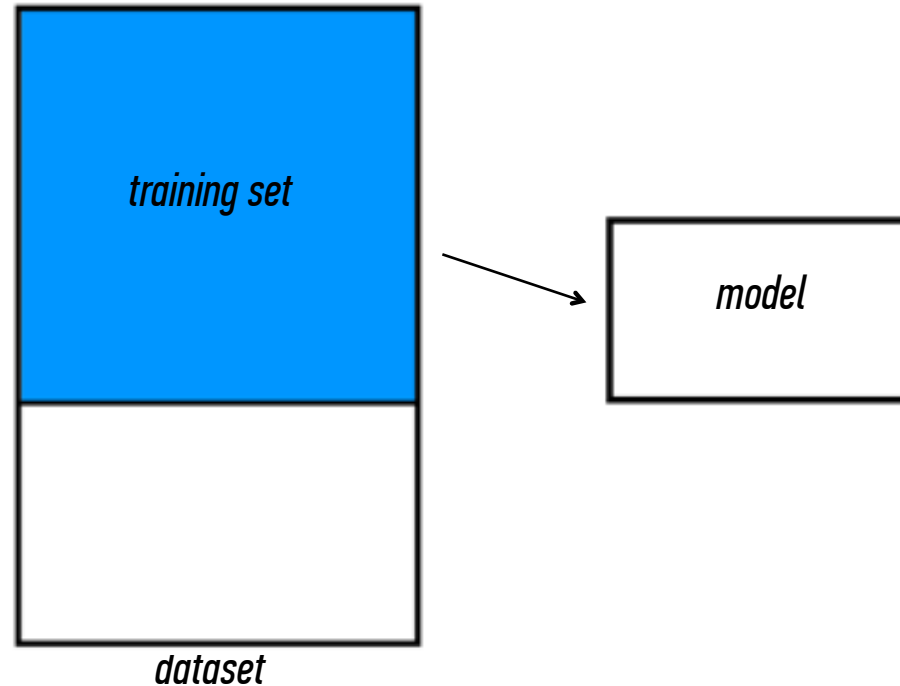
Q: What steps does a classification problem require?

1) split dataset



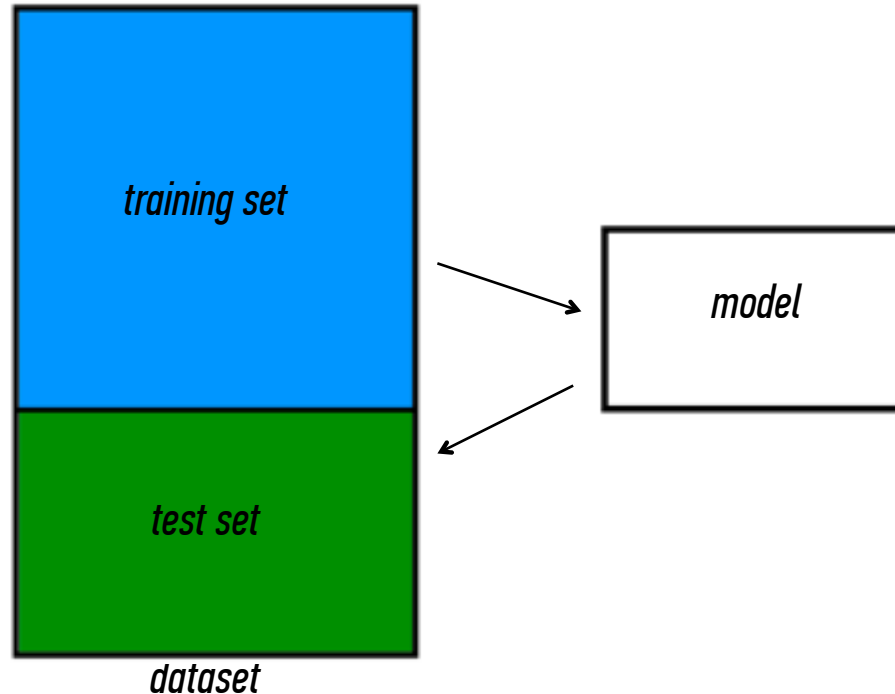
Q: What steps does a classification problem require?

- 1) split dataset*
- 2) train model*



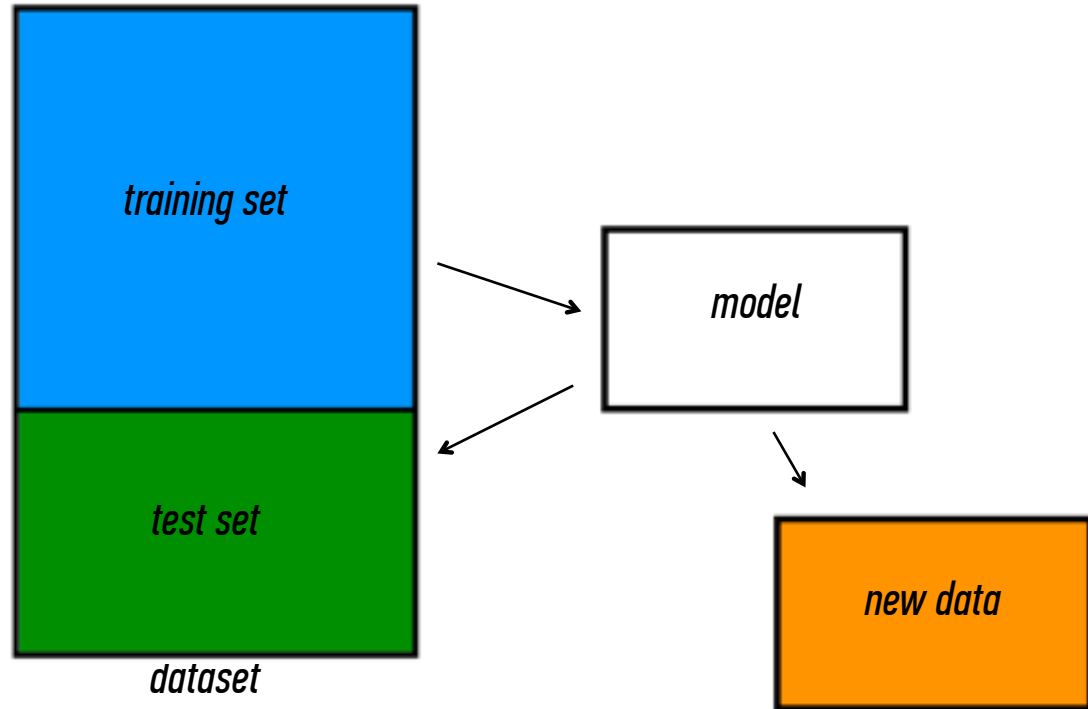
Q: What steps does a classification problem require?

- 1) split dataset*
- 2) train model*
- 3) test model*



Q: What steps does a classification problem require?

- 1) split dataset*
- 2) train model*
- 3) test model*
- 4) make predictions*

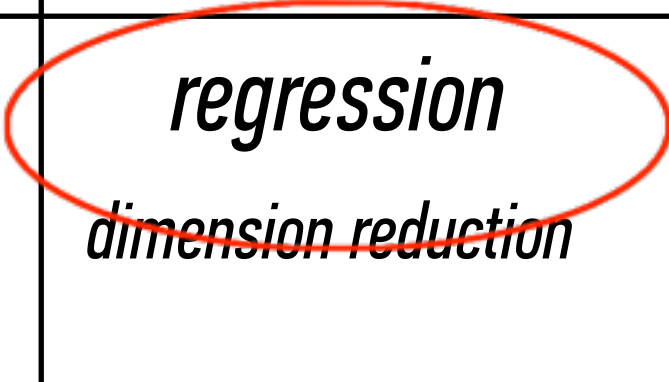


INTRO TO DATA SCIENCE

LINEAR REGRESSION

Wednesday, March 19, 14

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>



*Q: What is a **regression model**?*

A: A functional relationship between input & response variables

*The **simple linear regression model** captures a linear relationship between a single input variable x and a response variable y :*

$$y = \alpha + \beta x + \varepsilon$$

Q: What do the terms in this model mean?

$$y = \alpha + \beta x + \varepsilon$$

*A: y = **response variable** (the one we want to predict)*

*x = **input variable** (the one we use to train the model)*

*α = **intercept** (where the line crosses the y -axis)*

*β = **regression coefficient** (the model “parameter”)*

*ε = **residual** (the prediction error)*

OLS: $\min(\|y - x\beta\|^2)$

L1 regularization: $\min(\|y - x\beta\|^2 + \lambda\|\beta\|)$

L2 regularization: $\min(\|y - x\beta\|^2 + \lambda\|\beta\|^2)$

INTRO TO DATA SCIENCE

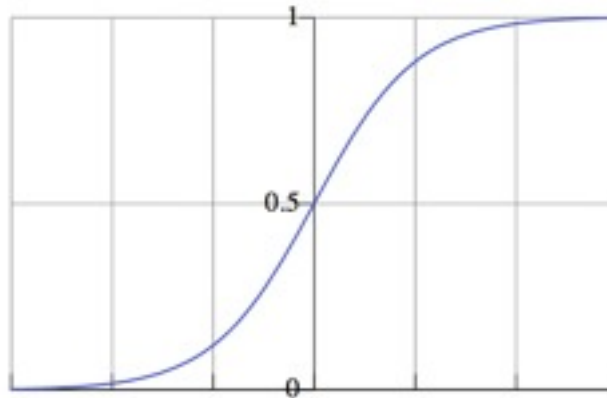
LOGISTIC REGRESSION

Wednesday, March 19, 14

$$E(y|x) = \pi(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

$$E(y|x) = \pi(x) = \frac{e^{\alpha+\beta x}}{1+e^{\alpha+\beta x}}$$

We've already seen what this looks like:



The logit function is an important transformation of the logistic function. Notice that it returns the linear model!

$$g(x) = \ln\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

*The **logit function** is an important transformation of the logistic function. Notice that it returns the linear model!*

$$g(x) = \ln\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

*The logit function is also called the **log-odds function**.*

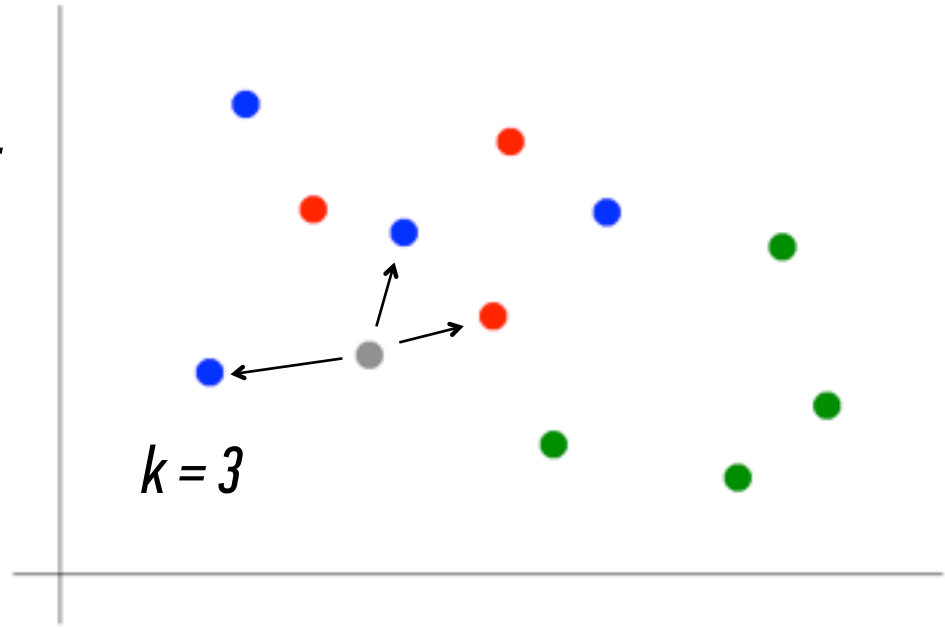
INTRO TO DATA SCIENCE

KNN CLASSIFICATION

Wednesday, March 19, 14

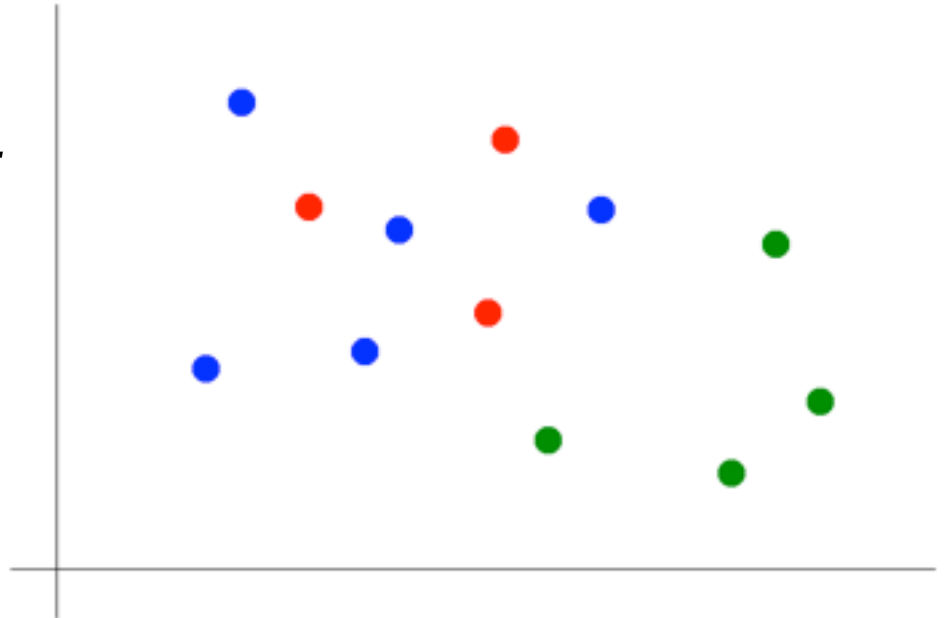
Suppose we want to predict the color of the grey dot.

- 1) Pick a value for k .*
- 2) Find colors of k nearest neighbors.*



Suppose we want to predict the color of the grey dot.

- 1) Pick a value for k .*
- 2) Find colors of k nearest neighbors.*
- 3) Assign the most common color to the grey dot.*



INTRO TO DATA SCIENCE

NAÏVE BAYES

Wednesday, March 19, 14

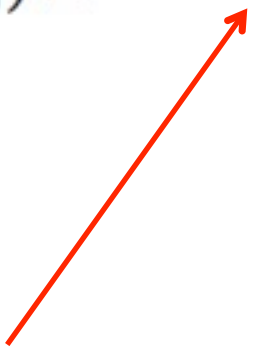
Suppose we have a dataset with features x_1, \dots, x_n and a class label C . What can we say about classification using Bayes' theorem?

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

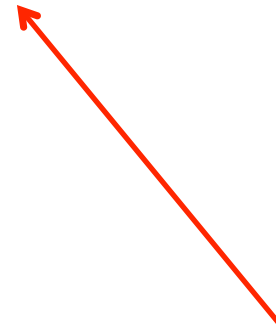
Bayes' theorem can help us to determine the probability of a record belonging to a class, given the data we observe.

source: *Data Analysis with Open Source Tools*, by Philipp K. Janert. O'Reilly Media, 2011.

*This term is the **likelihood function**. It represents the joint probability of observing features $\{x_i\}$ given that that record belongs to class C .*

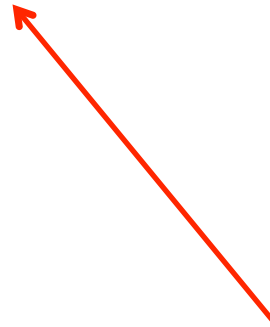
$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


*This term is the **prior probability** of C . It represents the probability of a record belonging to class C before the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


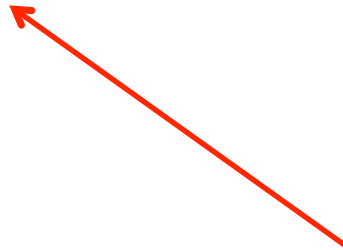
*This term is the **normalization constant**. It doesn't depend on C , and is generally ignored until the end of the computation.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$



*This term is the **posterior probability** of C . It represents the probability of a record belonging to class C after the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$



*This term is the **posterior probability** of C . It represents the probability of a record belonging to class C after the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

The goal of any Bayesian computation is to find (“learn”) the posterior distribution of a particular variable.

*The idea of Bayesian inference, then, is to **update** our beliefs about the distribution of C using the data (“evidence”) at our disposal.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

Then we can use the posterior for prediction.

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

INTRO TO DATA SCIENCE

COMPARISON

<p><i>linear</i></p> <p><i>scalability</i></p> <p><i>interpretation</i></p> <p><i>configuration</i></p> <p><i>feature-select</i></p> <p><i>overfitting</i></p>	
--	--

KNN

	<i>KNN</i>
<i>linear</i>	<i>N</i>

	<i>KNN</i>
<i>linear</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>

	<i>KNN</i>
<i>linear</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>
<i>interpretation</i>	<i>-</i>

	<i>KNN</i>
<i>linear</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>
<i>interpretation</i>	<i>-</i>
<i>configuration</i>	<i>+</i>

	<i>KNN</i>
<i>linear</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>
<i>interpretation</i>	<i>-</i>
<i>configuration</i>	<i>+</i>
<i>feature-select</i>	<i>-</i>

	<i>KNN</i>
<i>linear</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>
<i>interpretation</i>	<i>-</i>
<i>configuration</i>	<i>+</i>
<i>feature-select</i>	<i>-</i>
<i>overfitting</i>	<i>> K</i>

	<i>KNN</i>	<i>Logistic</i>
<i>linear</i>	<i>N</i>	
<i>scalability</i>	<i>$+/-$</i>	
<i>interpretation</i>	<i>$-$</i>	
<i>configuration</i>	<i>$+$</i>	
<i>feature-select</i>	<i>$-$</i>	
<i>overfitting</i>	<i>$> K$</i>	

	<i>KNN</i>	<i>Logistic</i>
<i>linear</i>	<i>N</i>	<i>Y</i>
<i>scalability</i>	<i>+/-</i>	<i>+</i>
<i>interpretation</i>	<i>-</i>	<i>+</i>
<i>configuration</i>	<i>+</i>	<i>+</i>
<i>feature-select</i>	<i>-</i>	<i>+</i>
<i>overfitting</i>	<i>> K</i>	<i>L1 / L2</i>

	<i>KNN</i>	<i>Logistic</i>	<i>NB</i>
<i>linear</i>	<i>N</i>	<i>Y</i>	<i>Y</i>
<i>scalability</i>	<i>+/-</i>	<i>+</i>	<i>+</i>
<i>interpretation</i>	<i>-</i>	<i>+</i>	<i>+</i>
<i>configuration</i>	<i>+</i>	<i>+</i>	<i>+</i>
<i>feature-select</i>	<i>-</i>	<i>+</i>	<i>+</i>
<i>overfitting</i>	<i>> K</i>	<i>L1 / L2</i>	<i>Prior</i>

	<i>KNN</i>	<i>Logistic</i>	<i>NB</i>	<i>RF</i>
<i>linear</i>	<i>N</i>	<i>Y</i>	<i>Y</i>	<i>N</i>
<i>scalability</i>	<i>+/-</i>	<i>+</i>	<i>+</i>	<i>-</i>
<i>interpretation</i>	<i>-</i>	<i>+</i>	<i>+</i>	<i>-</i>
<i>configuration</i>	<i>+</i>	<i>+</i>	<i>+</i>	<i>+</i>
<i>feature-select</i>	<i>-</i>	<i>+</i>	<i>+</i>	<i>+</i>
<i>overfitting</i>	<i>> K</i>	<i>L1 / L2</i>	<i>Prior</i>	<i>n tree</i>

	<i>KNN</i>	<i>Logistic</i>	<i>NB</i>	<i>RF</i>	<i>SVM</i>
<i>linear</i>	<i>N</i>	<i>Y</i>	<i>Y</i>	<i>N</i>	<i>Y/N</i>
<i>scalability</i>	<i>+/-</i>	<i>+</i>	<i>+</i>	<i>-</i>	<i>-</i>
<i>interpretation</i>	<i>-</i>	<i>+</i>	<i>+</i>	<i>-</i>	<i>-</i>
<i>configuration</i>	<i>+</i>	<i>+</i>	<i>+</i>	<i>+</i>	<i>-</i>
<i>feature-select</i>	<i>-</i>	<i>+</i>	<i>+</i>	<i>+</i>	<i>-</i>
<i>overfitting</i>	<i>> K</i>	<i>L1 / L2</i>	<i>Prior</i>	<i>n tree</i>	<i>C-cost</i>

QUESTION

***HOW
DO YOU
REPRESENT
YOUR
DATA?***

continuous

categorical

quantitative

qualitative

	<i>continuous</i>	<i>categorical</i>
<i>color</i>	<i>RGB-values</i>	<i>{red, blue}</i>
<i>ratings</i>	<i>1 – 10 rating</i>	<i>Good / Bad</i>

QUESTION

***HOW
DO YOU
MEASURE
OF
QUALITY?***

<i>supervised</i> <i>unsupervised</i>	<i>test out your predictions</i> <i>...</i>
--	--

<i>supervised</i>	<i>Accuracy, MAE, AUC</i>
<i>unsupervised</i>	<i>...</i>