# MACHINE LEARNING DESIGN, DEMYSTIFIED

SATURN 2018 Tutorial | May 8 | Plano

**Carnegie Mellon University**
Software Engineering Institute

softserve

**Carnegie Mellon University**
Software Engineering Institute

softserve

# INTRODUCTIONS

**Rick Kazman**

Professor, University of Hawaii
Research Scientist, SEI

**Serge Haziyev**

Head of Intelligent
Enterprise, SoftServe

**Iurii Milovanov**

Data Science Practice Leader,
SoftServe

# AGENDA

- **A Bit of Background**
- **Game** ☺
- **Prototyping**
- **Summary & QA**

# MOTIVATION



Architectural drivers

<<precedes>>

**Architectural design**

<<precedes>>

Architectural documentation

<<precedes>>

Architectural evaluation

Clouds

Mobile

IoT

**Machine Learning**

Big Data

Blockchain

# ADD 3.0



ML Design Concepts:
- **Problem**
- **Algorithm Family**
- **Algorithm**

**Design objectives** | **Primary functional requirements** | **Quality attribute scenarios** | **Constraints** | **Concerns**

**Step 1: Review Inputs**

**Step 2: Establish iteration goal and select inputs to be considered in the iteration**

**Step 3: Choose one or more elements of the system to decompose**

**Step 4: Choose one or more design concepts that satisfy the inputs considered in the iteration**

**Step 5: Instantiate architectural elements, allocate responsibilities and define interfaces**

**Step 6: Sketch views and record design decisions**

**Step 7: Perform analysis of current design and review iteration goal and design objectives**

Step 8: Refine as necessary

**Software architecture design**

Input/output artifact

Process step

**Carnegie Mellon University**
Software Engineering Institute

for public release and unlimited distribution.

# SMART DECISIONS GAME

First presented at SATURN 2015

A fun, lightweight way to introduce architectural design and ADD

Available at:

http://smartdecisionsgame.com/

softserve

**Carnegie Mellon University**
Software Engineering Institute

# SHORT QUIZ ☺
## What's the name of this company in AI field?

### 10x increase over the past 2 years!



245.33 USD  Mar 9, 2018

# AI PROGRESS SINCE 1950s



ARTIFICIAL INTELLIGENCE
Early artificial intelligence stirs excitement.

MACHINE LEARNING
Machine learning begins to flourish.

DEEP LEARNING
Deep learning breakthroughs drive AI boom.

1950's   1960's   1970's   1980's   1990's   2000's   2010's

*Source: NVidia*

# MYTHS AND FICTION ABOUT ARTIFICIAL BEINGS



**R.U.R. (Karel Čapek)**
**1921**



**Golem (Bible)**
**~1000 BC**



**Sumerian Anunnaki creating the first man**
**~2300 BC**

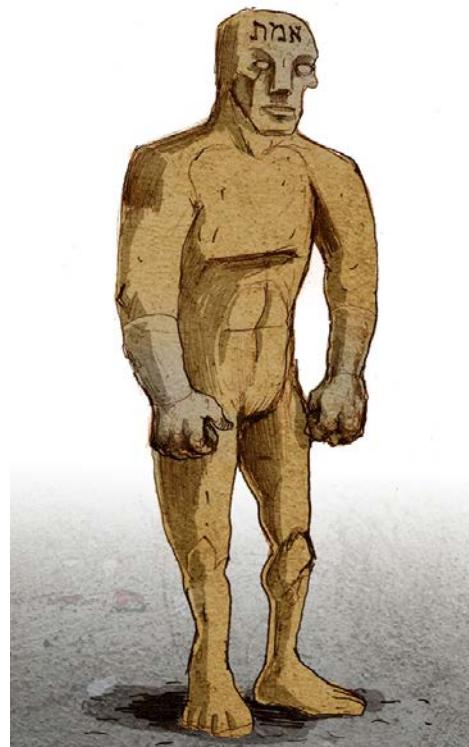# THE CURRENT STATE OF AI



| 5K | 1M | 16M | 71M | 760M | 22B | 86B |

Number of neurons

soft**serve**

# GAME CHALLENGE OVERVIEW
## Business Use Case



Banner A



Banner B

**Which ad will the user choose?**

# GAME CHALLENGE OVERVIEW
## Marketecture Diagram

$$CTR = \frac{\text{Number of click-throughs}}{\text{Number of impressions}} \times 100(\%)$$



**Ingest**

**Web Servers**

- Hundreds of servers
- Massive logs from multiple sources

**Web Logs**

**Train**

**CTR Prediction System**

**Online user**

Predict CTR and show appropriate ad

**Carnegie Mellon University**
Software Engineering Institute

**softserve**

# WHY DO WE NEED MACHINE LEARNING?

# WHY NOT JUST CODING?!

**Most of the today's AI problems:**

- Deal with an infinite problem space – think about how many words are there in the English language

- Poorly defined – we still do not know how our brain solves problems

**Therefore, traditional rule-based hand-coding for such problems suffers a 'complexity collapse' and is not feasible**

softserve

# MACHINE LEARNING APPROACH

Instead of writing a program by hand, we use a set of examples to train the algorithm



**Developer** writes code



**Algorithm** "writes code"

# ML BUILDING BLOCKS

# TYPES OF LEARNING

# SUPERVISED LEARNING

- **Input examples** and corresponding **ground truth outputs** are provided

- The goal is to learn general rules that map a new example to the predicted output

# SUPERVISED LEARNING

Input Data

Output Data

**Example:** Given a set of **house features** along with corresponding **house prices**, predict a price for a new house based on its features (e.g. size, location, etc.)

# UNSUPERVISED LEARNING

- Only **input examples** are provided

- No explicit information about ground truth

- The algorithm tries to discover the internal structure of the data based on some prior knowledge about desired outcome

# UNSUPERVISED LEARNING

**Input Data**

**Example:** Given a set of **customer transactions** discover what would be the best way to group them into clusters based on **customer similarity**

**Output Preferences**

# SUPERVISED LEARNING

# UNSUPERVISED LEARNING

# ITERATION 1:

What type of learning best fits a given use case?

Select from: supervised or unsupervised

# ITERATION 1:
# Supervised or Unsupervised Learning?

**Historical Dataset**



Banner A                    Banner B

**Train**
**Input:** Logs

**ML Algorithm**

**Predict**
**Input:** User Features
**Output:** Most preferred banner to show

**Logs:**

User X Features | Banner A | Click 0
User Y Features | Banner B | Click 0
User X Features | Banner B | Click 1
...

**Carnegie Mellon University**
Software Engineering Institute

**softserve**

# MACHINE LEARNING CARDS

# MACHINE LEARNING CARDS

PROBLEM — ALGORITHM FAMILY — ALGORITHM

- **ANOMALY DETECTION**
- **REGRESSION**
- **CLASSIFICATION**
- **CLUSTERING**

Algorithm Families:
- SUPPORT VECTOR MACHINES
- DECISION TREE LEARNING
- INSTANCE-BASED LEARNING
- GENERALIZED LINEAR MODELS
- ARTIFICIAL NEURAL NETWORK
- CENTROID-BASED CLUSTERING
- HIERARCHICAL CLUSTERING
- DENSITY-BASED CLUSTERING

Algorithms:
- One-Class SVM
- Linear SVM
- Non-Linear SVM
- Classification/Regression Decision Tree
- Random Forest
- Isolation Forest
- Radius Neighbors
- K-Nearest Neighbors
- Logistic Regression
- Bayesian Naive Classifier
- Linear Regression
- Bayesian Linear Regression
- Feedforward ANN (Multilayer Perceptron)
- K-Means Clustering
- Complete-Linkage Clustering
- Single-Linkage Clustering
- Average-Linkage Clustering
- DBSCAN

**Legend:**

- Problem cards

- Algorithm Family cards

- Algorithm cards

**Carnegie Mellon University**
Software Engineering Institute

**softserve**

# PROBLEM TYPES

# CLASSIFICATION

**Key Highlights:**

- Identifies which category an object belongs to
- Supervised learning problem

**Examples:**

- Detect fraudulent transactions (one-class)
- Categorize emails by spam or not spam (binary)
- Categorize articles based on their topic (multi-class)
- Detect objects on the image (multi-label)

# REGRESSION

**Key Highlights:**

- Predict a continuous value associated with an object
- Supervised learning problem

**Examples:**

- Predict stock prices from market data
- Score a credit application based on historical data
- Estimate demand for a given product

# CLUSTERING

**Key Highlights:**

- Group similar objects into clusters
- Unsupervised learning problem

**Examples:**

- Discover audiences to target on social networks
- Group checking data based on GEO-proximity
- Detect common topics in corporate knowledge base

# ANOMALY DETECTION

**Key Highlights:**

- Identify observations that do not conform to an expected pattern

- Addresses both supervised and unsupervised learning

**Examples:**

- Identify fraudulent transactions or abnormal customer behavior

- In manufacturing, detect physical parts that are likely to fail in the near future

# ITERATION 2:
What type of problem best fits a given use case?

Select problem card from: classification, regression, clustering or anomaly detection

# ITERATION 2:
## What type of problem?

**Historical Dataset**



Banner A            Banner B

**Logs:**

    User X Features | Banner A | Click 0
    User Y Features | Banner B | Click 0
    User X Features | Banner B | Click 1
    …

**Train**
**Input:** Logs

**ML Algorithm**

**Predict**
**Input:** User Features
**Output:** Most preferred banner to show

# FAMILIES AND ALGORITHMS

# CLASSIFICATION FAMILIES

## Artificial Neural Network
**ALGORITHM FAMILY**

**Description:** ANN is a computational model based on the structure and functions of biological neural networks. [It] can be used for both classification and regression [or] unsupervised learning. The model works well for [big] data and complex non-linear relationships but [is] computationally expensive.



**Characteristics:**
- ★★★ Big Data — performance positively correl[ated with] volume, but at the expense of computation[al...]
- ★ Small Data — the less data available, the [more] other algorithms will outperform neural n[etworks]
- ★★ Imbalanced Data — class imbalance incre[ases the] possibility of overfitting, but it can be mit[igated with] class weights
- ★ Results Interpretation — usually difficult
- ★★★ Online Learning — can be trained sequen[tially]
- ★ Ease of Use — requires substantial under[standing of] ANN field

**Algorithms:** Feedforward ANN (Multilayer Perce[ptron)...]

## Support Vector Machines
**ALGORITHM FAMILY**

**Description:** SVMs are supervised learning algorithms which can be used for both classification or regression problems (except OneClassSVM which is unsupervised). Given labeled training data, the algorithms output an optimal hyperplane which categorizes new examples.



**Characteristics:**
- ★ Big Data — computation and memory-intensive while training, data sampling can be used as a workaround
- ★★ Small Data — good accuracy even for small # of observations
- ★★★ Imbalanced Data — compensates imbalance through class weights, works fine out of the box for moderately imbalanced data
- ★★ Results Interpretation — applicable for regulated fields (i.e. credit scoring)
- ★ Online Learning — most implementations only support batch setting
- ★★ Ease of Use — moderate number of parameters

**Implementations:** Linear SVM, Non-Linear SVM, One-Class SVM.

*MACHINE LEARNING*

## Generalized Linear Models
**ALGORITHM FAMILY**

**Description:** A class of linear models with a common property [where the depen]dent variable is linearly related to t[he...] [vi]a a specified link function. General[ized linear models were] [p]roposed as a way of unifying vario[us...]



[Characteris]tics:
- [Big Da]ta — there are a lot of implementa[tions in various] [packe]ts, that include iterative approach[es and...] [linear] equation systems
- [Small] Data — overfitting resistance due [to...] [statist]ically significant dependencies only
- [Imbal]anced Data — poor, can be handled [by imbal]ance compensations techniques, e.[g...] [samp]ling
- [Result]s Interpretation — results drivers [can be clearly] [interpre]ted
- [Online] Learning — most implementation[s only support] [batch] setting
- [Ease o]f Use — models tuning is user-frien[dly]

[Algorithms:] Logistic Regression, Bayesian Naive [Bayes,] Bayesian Linear Regression.

*MACHINE LEARNING*

## Decision Tree Learning
**ALGORITHM FAMILY**

**Description:** Decision Tree Learning uses a decision tree structure to go from observations about an item to conclusions about the item's target value. It is one of the most interpretable families of machine learning algorithms. This approach can be used for both classification or regression problems.



**Characteristics:**
- ★★ Big Data — interpretability is getting worse on large datasets
- ★★ Small Data — sufficient generalization even for very small dataset, but can lead to overfitting
- ★★ Imbalanced Data — can be handled by stratified bootstrap technique
- ★★★ Results Interpretation — represented by a set of decision rules
- ★★★ Online Learning — can be trained sequentially
- ★★★ Ease of Use — models tuning is user-friendly

**Algorithms:** Classification/Regression Decision Tree, Random Forest, Isolation Forest.

*MACHINE LEARNING*

## Instance-Based Learning
**ALGORITHM FAMILY**

**Description:** Instance-Based Learning (sometimes called [memory-bas]ed learning) is a family of learning algorithms that, [instead of per]forming explicit generalization, compares new [inst]ances with instances seen in training and stored in [...] This approach can be used for both classification or [regression pr]oblems.



[Characteris]tics:
- [Big Da]ta — generally, the application of this algorithms [to big] data is not feasible due to memory and prediction [time] restrictions
- [Small] Data — good accuracy, even for a small number of [obser]vations
- [Imbal]anced Data — more robust for data with [imbal]anced classes and is efficient for multiclass [classif]ication with a small number of features
- [Result]s Interpretation — transparent inference; process, [wit]hout the possibility of getting explanation rules
- [Online] Learning — can be trained sequentially
- [Ease o]f Use — models tuning is user-friendly

[Algorithms:] K-Nearest Neighbors, Radius Neighbors.

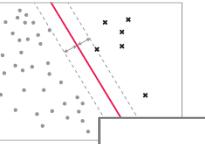*MACHINE LEARNING*

# CLASSIFICATION ALGORITHMS
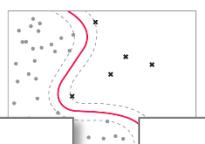
## Linear SVM
**ALGORITHM**

**Description:** An SVM algorithm with a linear kernel.

**Characteristics:**
- ★ Accuracy — depends on da... is required)
- ★ Training Speed — linear de... dataset size
- ★★★ Prediction Speed — depe... features
- ★★ Overfitting Resistance — b... hyperparameters
- ★ Probabilistic Interpretation... expensive cross-validation

**Tips:**
✓ Great performance in high-dim...
✓ Works well in case when # of fe... observations
✓ Margins maximization provides...

**Implementations:** R, Python (sc...

## Non-Linear SVM
**SUPPORT VECTOR MACHINES | ALGORITHM**

**Description:** An SVM algorithm with a non-linear kernel like Gaussian (e.g. RBF).

...curacy for m...
...aining time...
...depends on...

...ance — theore...
...verfitting, b...
...rs
...retation — c...
...idation

...gh-dimensio...
...the problem...

...e learning is...
...tion (scikit-le...

## Bayesian Naive Classifier
**GENERALIZED LINEAR MODELS | ALGORITHM**

**Description:** A simple probabilistic classifier based on Bayes' Theorem with an assumption of strong independence (naive) among features.
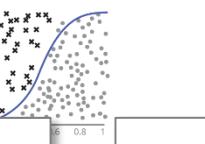
$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B \mid A)P(A)}{P(B)}$$

spam    viagra

...is required...
...bility of eac...
...rges toward...
...logistic reg...
...nds on # of...
...great gene...
...ing set, due to...

...on — classe...
...nary classi...

...minative mo...
...ining data...
...(scikit-learn),...

## Logistic Regression
**GENERALIZED LINEAR MODELS | ALGORITHM**

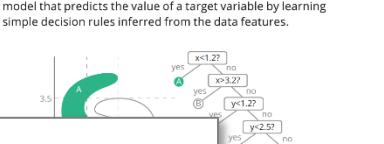**Description:** A powerful tool for two-class and multiclass classification. As a linear regression, it is fast and simple.

...ations into...
...separable, i...

...ary classifi...
...ataset using...

...ns can be e...
...nally easy y...
...sidered rob...
...features. Re...

...— extracted...

...g one-vs-a...
...uence of th...

...learn), Spar...
MA...

## Classification/Regression Trees
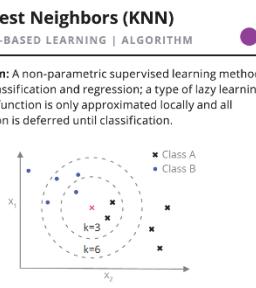**DECISION TREE LEARNING | ALGORITHM**

**Description:** A non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.

...of features; a tree with...
...space is very likely to...

...hms provide the linear...
...he tree size...
...create over-complex...
...om the training data...
...be extracted...

...nd numerical data...

...), Spark (mllib).
**MACHINE LEARNING**

## K-Nearest Neighbors (KNN)
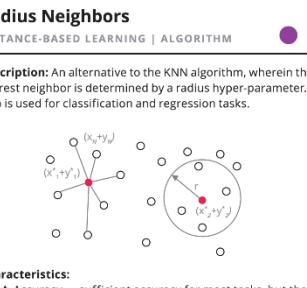**INSTANCE-BASED LEARNING | ALGORITHM**

**Description:** A non-parametric supervised learning method used for classification and regression; a type of lazy learning, where the function is only approximated locally and all computation is deferred until classification.

Class A
Class B

k=3
k=6

**Characteristics:**
- ★★ Accuracy — sufficient accuracy for most tasks, but there is a tradeoff between accuracy vs avoiding overfitting
- ★★ Training Speed — training time is high on large datasets
- ★ Prediction Speed — full training set processing is required
- ★★ Overfitting Resistance — with an increase of k nearest training objects, the probability of overfitting decreases
- ★★★ Probabilistic Interpretation — naturally determined by the inference process

**Tips:**
✓ One of the simplest machine learning algorithms
✓ Good choice for low dimensional space

**Implementations:** R, Python (scikit-learn).

## Radius Neighbors
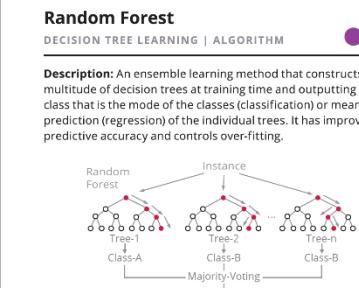**INSTANCE-BASED LEARNING | ALGORITHM**

**Description:** An alternative to the KNN algorithm, wherein the nearest neighbor is determined by a radius hyper-parameter. Also is used for classification and regression tasks.

$(x_w + y_w)$
$(x'_1 + y'_1)$
$(x'_2 + y'_2)$
r

**Characteristics:**
- ★★ Accuracy — sufficient accuracy for most tasks, but there is a tradeoff between accuracy vs avoiding overfitting
- ★★ Training Speed — training time is high on large datasets
- ★ Prediction Speed — full training set processing is required
- ★★ Overfitting Resistance — with an increase of radius the probability of overfitting decreases
- ★★★ Probabilistic Interpretation — naturally determined by the inference process

**Tips:**
✓ One of the simplest machine learning algorithms
✓ Good choice for data that isn't sampled uniformly
✓ For high-dimensional spaces, this method is less effective due to so-called "curse of dimensionality"

**Implementations:** R, Python (scikit-learn).

## Random Forest
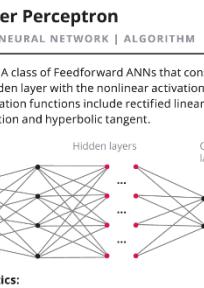**DECISION TREE LEARNING | ALGORITHM**

**Description:** An ensemble learning method that constructs a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. It has improved predictive accuracy and controls over-fitting.

Random Forest    Instance

Tree-1    Tree-2    Tree-n
Class-A   Class-B   Class-B

Majority-Voting
Final-Class

**Characteristics:**
- ★★★ Accuracy — due to bootstrap aggregation of independent trees
- ★★ Training Speed — depends on model complexity, but it can be parallelized
- ★★ Prediction Speed — depends on model complexity but it can be improved due to independent trees
- ★★★ Overfitting Resistance — complexity doesn't lead to overfitting
- ★★ Probabilistic Interpretation — can be extracted

**Tips:**
✓ Naturally works with both categorical and numerical data
✓ High accuracy without extensive tuning

**Implementations:** R, Python (scikit-learn), Spark (mllib), Azure.

## Multilayer Perceptron
**ARTIFICIAL NEURAL NETWORK | ALGORITHM**

**Description:** A class of Feedforward ANNs that consists of at least one hidden layer with the nonlinear activation function. Popular activation functions include rectified linear unit (ReLU), sigmoid function and hyperbolic tangent.

Input layer    Hidden layers    Output layer

**Characteristics:**
- ★★★ Accuracy — works well for both linear and nonlinear dependencies
- ★★ Training Speed — depends heavily on model complexity and on training dataset size
- ★★★ Prediction Speed — depends on # of model features, scales well
- ★ Overfitting Resistance — requires big training set and regularization
- ★★★ Probabilistic Interpretation — thanks to the softmax activation

**Tips:**
✓ Good performance for high dimensional space
✓ Works well with numerical and categorical data

**Implementations:** R, Python (scikit-learn), Spark (mllib, classificatory only), Azure.

**MACHINE LEARNING**
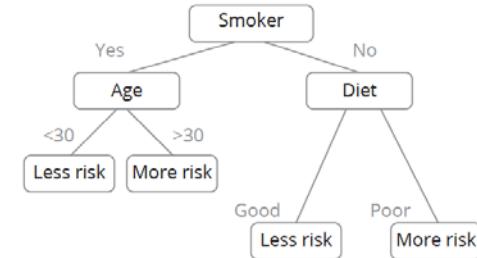
softserve

# DECISION DRIVERS

# FAMILY DRIVERS

- **Big Data** – scalability and ability to leverage from new data

- Small Data – ability to learn from a few examples

- **Imbalanced Data** – ability to distinguish rare events

- Results Interpretation – human-friendly results

- Online Learning – ability to continuously train from new data

- **Ease of Use** – number of parameters to manually tune



**Decision Tree Learning**
ALGORITHM FAMILY

**Description:** Decision Tree Learning uses a decision tree structure to go from observations about an item to conclusions about the item's target value. It is one of the most interpretable families of machine learning algorithms. This approach can be used for both classification or regression problems.

**Characteristics:**
- ★★ Big Data — interpretability is getting worse on large datasets
- ★★ Small Data — sufficient generalization even for very small dataset, but can lead to overfitting
- ★★ Imbalanced Data — can be handled by stratified bootstrap technique
- ★★★ Results Interpretation — represented by a set of decision rules
- ★★★ Online Learning — can be trained sequentially
- ★★★ Ease of Use — models tuning is user-friendly

**Algorithms:** Classification/Regression Decision Tree, Random Forest, Isolation Forest.

MACHINE LEARNING

**Carnegie Mellon University**
Software Engineering Institute

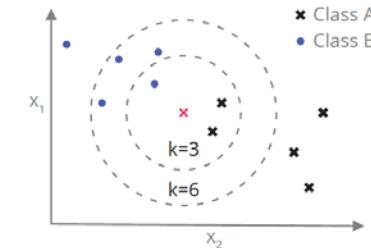softserve

# ALGORITHM DRIVERS

- **Accuracy** – ability to solve complex problems

- **Training Speed** – training runtime performance

- **Prediction Speed** – production runtime performance

- Overfitting Resistance – ability to generalize to new data

- Probabilistic Interpretation – return results as probabilities

## K-Nearest Neighbors (KNN)

INSTANCE-BASED LEARNING | ALGORITHM

**Description:** A non-parametric supervised learning method used for classification and regression; a type of lazy learning, where the function is only approximated locally and all computation is deferred until classification.



**Characteristics:**
- ★★ Accuracy — sufficient accuracy for most tasks, but there is a tradeoff between accuracy vs avoiding overfitting
- ★★ Training Speed — training time is high on large datasets
- ★ Prediction Speed — full training set processing is required
- ★★ Overfitting Resistance — with an increase of k nearest training objects, the probability of overfitting decreases
- ★★★ Probabilistic Interpretation — naturally determined by the inference process

**Tips:**
- ✔ One of the simplest machine learning algorithms
- ✔ Good choice for low dimensional space

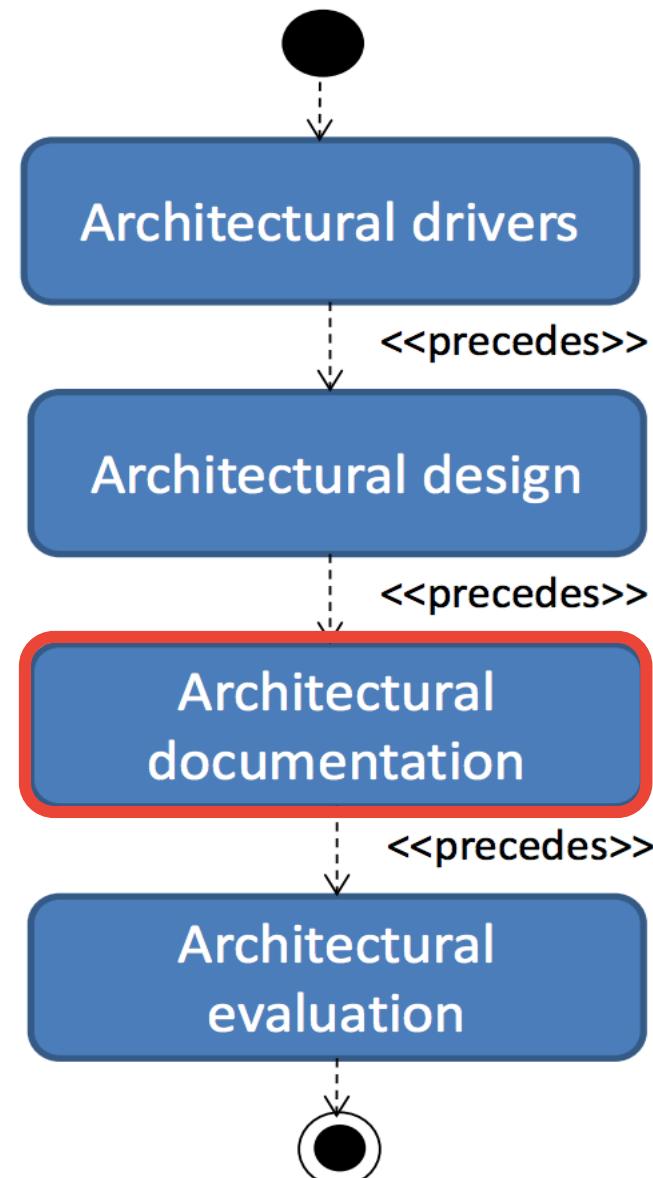**Implementations:** R, Python (scikit-learn).

MACHINE LEARNING

**Carnegie Mellon University**
Software Engineering Institute

soft**serve**

# ITERATION 3:

Select a family and an algorithm card that would best fit a given use case
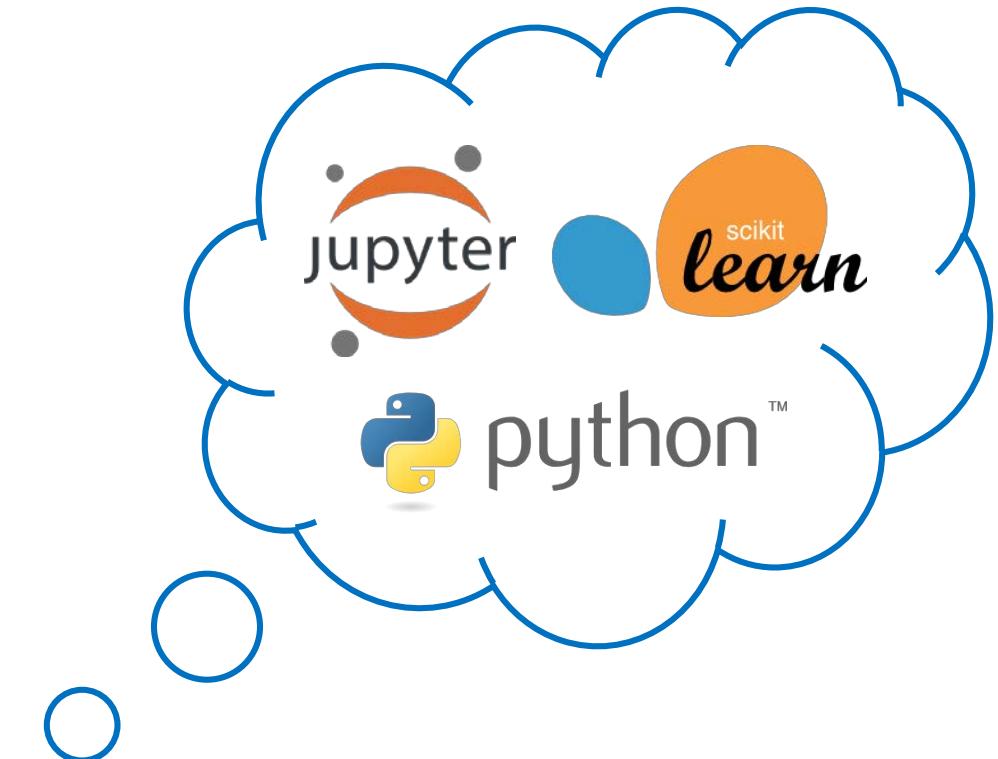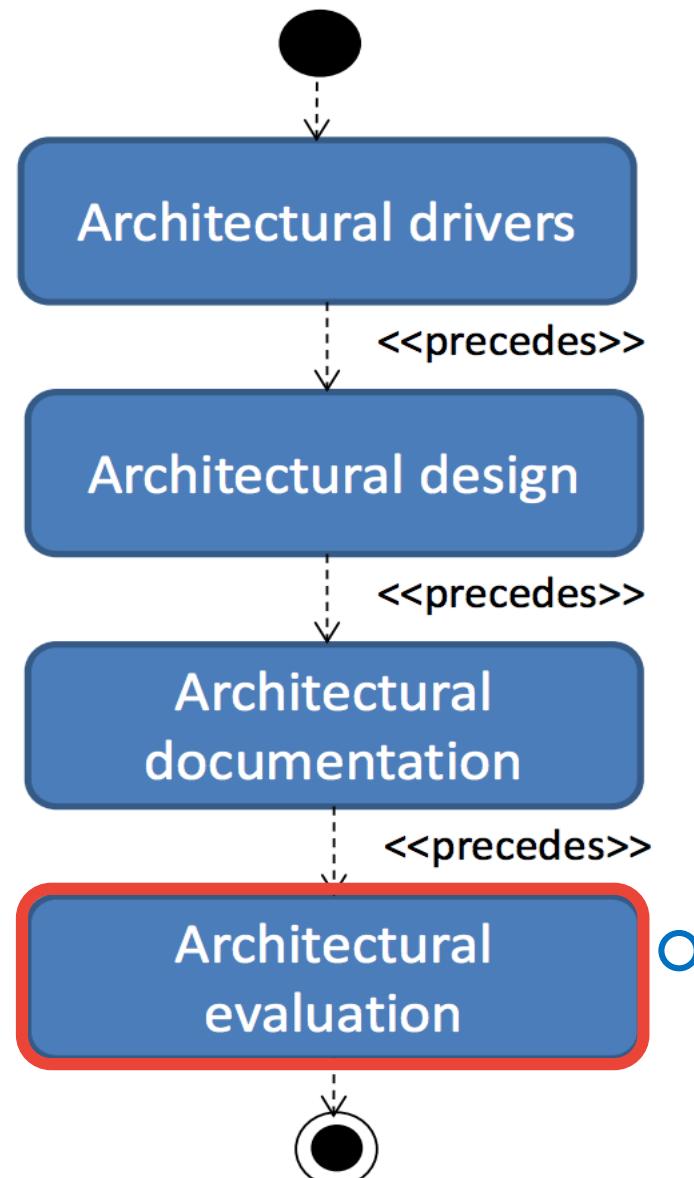
Family Key Drivers: Big Data, Imbalanced Data, Ease of Use
Algorithm Key Drivers: Accuracy, Training and Prediction Speed

# DESIGN PROCESS

# DESIGN PROCESS

Architectural drivers

<<precedes>>

Architectural design

<<precedes>>

Architectural documentation

<<precedes>>

Architectural evaluation

jupyter  scikit learn

python™

# PROTOTYPING AND **EVALUATION SESSION**

# PROTOTYPING FOR EVALUATION

softserve

# RESULTS SUMMARY

| Algorithm name | Training Time | Prediction Time | Tuning Time | Initial Accuracy | Final Accuracy |
|---|---|---|---|---|---|
| Random Forest | 2.61 | 0.47 | 94.44 | 81.61% | **83.05%** |
| KNeighbors | 0.41 | **44.29** | 84.27 | 80.57% | **83.05%** |
| Logistic Regression | 0.12 | 0.05 | **45.94** | **82.93%** | 82.93% |
| MLP | 0.80 | 0.08 | 164.04 | **66.25%** | 82.90% |
| SVM | **177.78** | 54.87 | **973.73** | 82.83% | 82.83% |
| Linear SVM | 5.93 | 0.04 | 82.91 | 82.69% | 82.69% |
| Decision Trees | 0.03 | **0.005** | 52.97 | 73.16% | 82.36% |
| Naive Bayes | **0.02** | 0.01 | 0 | 78.46% | **78.46%** |

**Carnegie Mellon University**
Software Engineering Institute

softserve

# KEY TAKEAWAYS

- Machine Learning solution design is an iterative process

- ADD principles help make ML design decisions in a systematic way

- ML Cards aim to select candidate algorithms from a wide variety of alternatives

- Prototyping is necessary to validate design decisions

**Carnegie Mellon University**
Software Engineering Institute

soft**serve**

QUESTIONS?
WE'VE GOT THE
ANSWERS.