# Azure ML Tutorial

## - Azure Bootcamp in Troy -

2016. 4. 16

Ingyu Lee
(inlee@troy.edu)

# Contents

Microsoft

# Azure Machine Learning

## Microsoft Azure Essentials

Jeff Barnes

Book, Microsoft Azure Essentials: Azure Machine Learning,
http://www.microsoftvirtualacademy.com/ebooks#9780735698178

# Part I What is Azure ML(Machine Learning)?

1. What is machine learning?
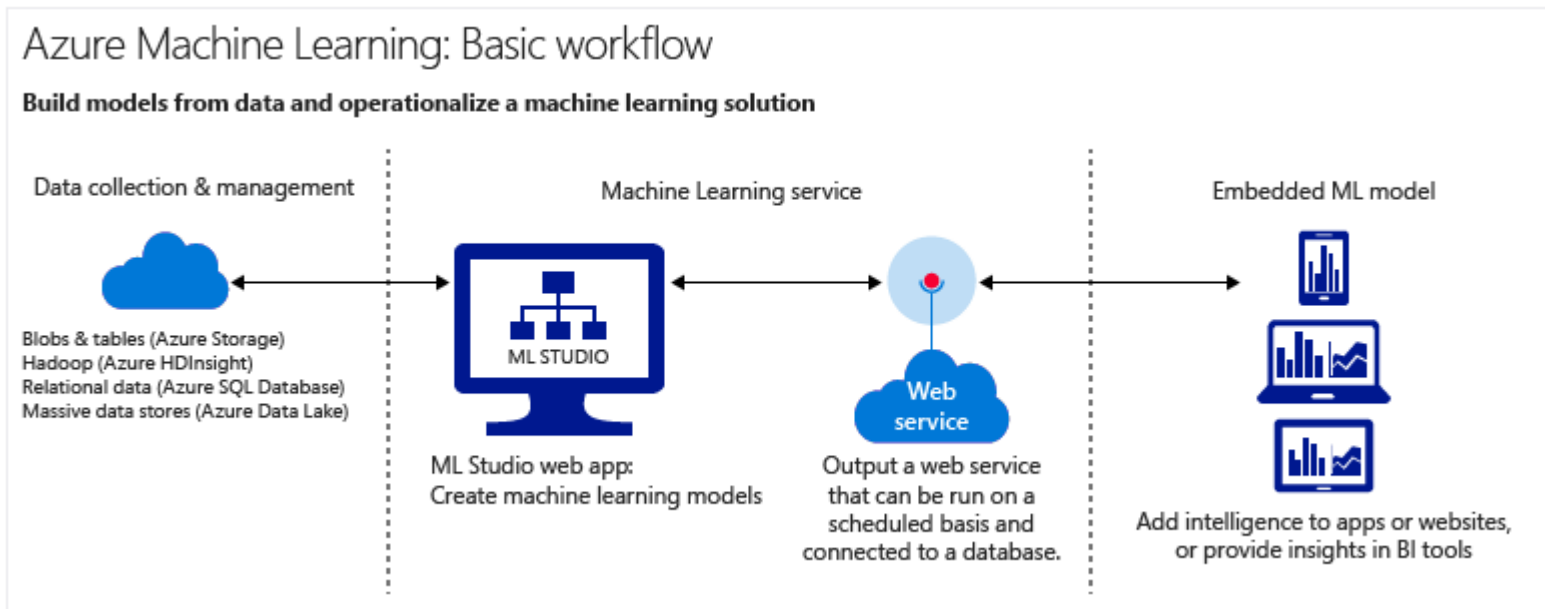    a. Machine learning uses computers to run predictive models that learn from existing data in order to forecast future behaviors, outcomes, and trends.
    b. These forecasts or predictions from machine learning can make apps and devices smarter.
    c. When you shop online, machine learning helps recommend other products you might like based on what you've purchased.
    d. When your credit card is swiped, machine learning compares the transaction to a database of transactions and helps the bank do fraud detection.

# Part I What is Azure ML(Machine Learning)?

2. What is Machine Learning on Microsoft Azure?

    a. Azure Machine Learning is a cloud-based predictive analytics service that makes it possible to quickly build and deploy predictive models as analytic solutions.

    b. Azure Machine Learning not only provides tools to model predictive analytics, but also provides a fully-managed service you can use to deploy your predictive models as ready-to-consume web services.

    c. Azure Machine Learning provides tools for creating complete predictive analytics solutions in the cloud: Quickly create, test, operationalize, and manage predictive models.
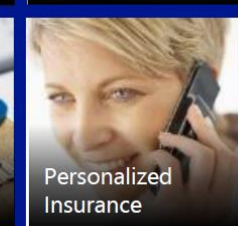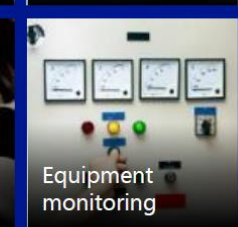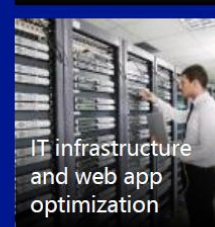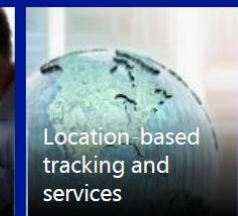
## Azure Machine Learning: Basic workflow

**Build models from data and operationalize a machine learning solution**

| Data collection & management | Machine Learning service | Embedded ML model |
| --- | --- | --- |

Blobs & tables (Azure Storage)
Hadoop (Azure HDInsight)
Relational data (Azure SQL Database)
Massive data stores (Azure Data Lake)

ML STUDIO

Web service

ML Studio web app:
Create machine learning models

Output a web service that can be run on a scheduled basis and connected to a database.

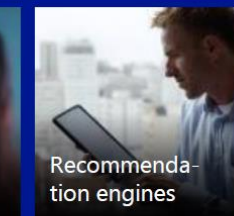Add intelligence to apps or websites, or provide insights in BI tools

# Part I What is Azure ML(Machine Learning)?
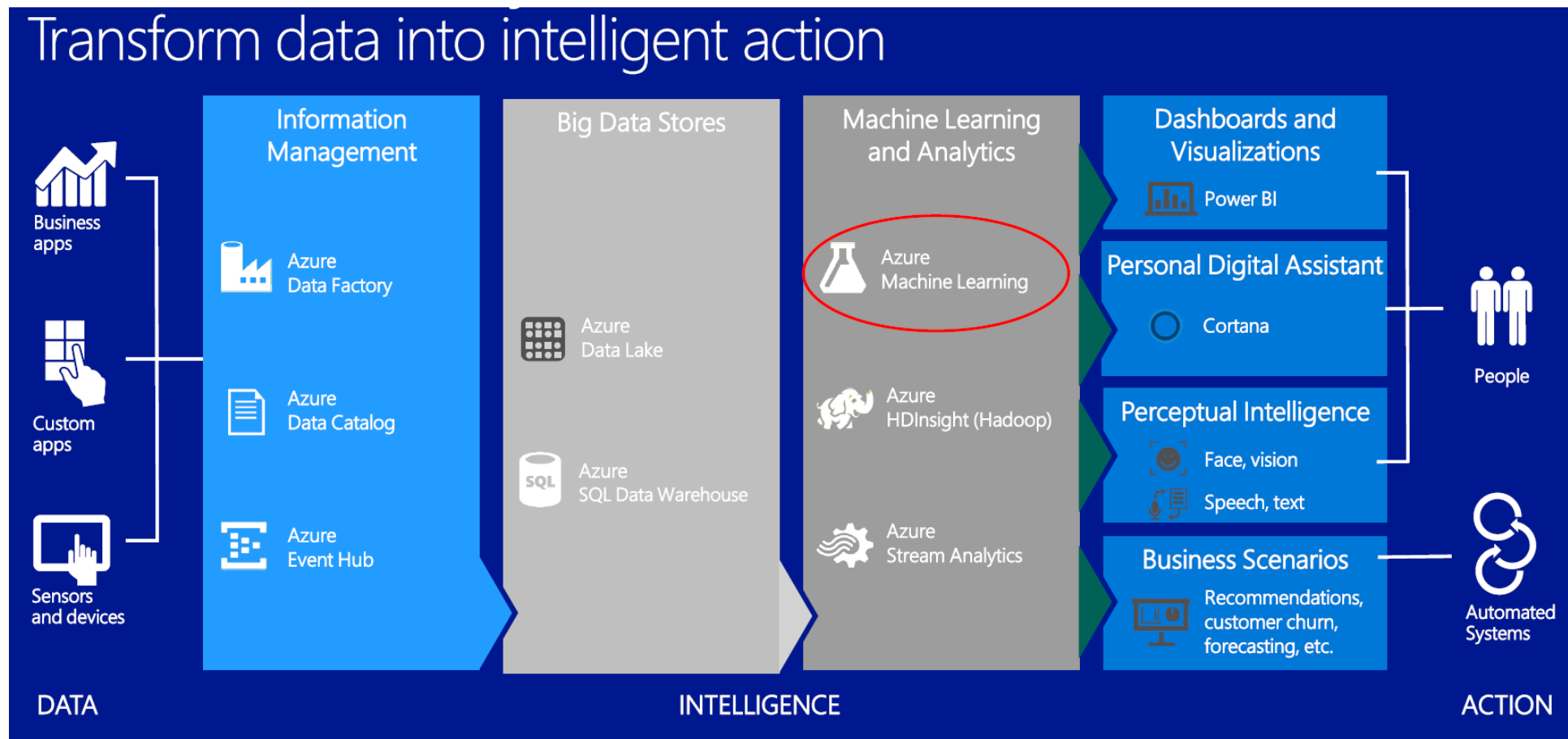
3. What is predictive analytics?
   a. Predictive analytics uses various statistical techniques – in this case, machine learning – to analyze collected or current data for patterns or trends in order to forecast future events.
   b. Azure Machine Learning is a particularly powerful way to do predictive analytics:
      1) You can work from a ready-to-use library of algorithms, create models on an internet-connected PC without purchasing additional equipment or infrastructure, and deploy your predictive solution quickly.
      2) You can also find ready-to-use examples and solutions in the Microsoft Azure Marketplace or Cortana Intelligence Gallery.

Predictive analytics should address the likelihood of something happening in the future, even if it is just an instant later...

Churn analysis

Social network analysis

Recommenda-tion engines

Location-based tracking and services

IT infrastructure and web app optimization

Weather forecasting for business planning

Legal discovery and document archiving

Equipment monitoring

Advertising analysis

Pricing analysis

Fraud detection

Personalized Insurance

# Part I Cortana Analytics Suite



Transform data into intelligent action

| DATA | INTELLIGENCE | | | ACTION |
|---|---|---|---|---|

Business apps
Custom apps
Sensors and devices

**Information Management**
- Azure Data Factory
- Azure Data Catalog
- Azure Event Hub

**Big Data Stores**
- Azure Data Lake
- Azure SQL Data Warehouse

**Machine Learning and Analytics**
- Azure Machine Learning
- Azure HDInsight (Hadoop)
- Azure Stream Analytics

**Dashboards and Visualizations**
- Power BI

**Personal Digital Assistant**
- Cortana

**Perceptual Intelligence**
- Face, vision
- Speech, text

**Business Scenarios**
- Recommendations, customer churn, forecasting, etc.

People
Automated Systems

# Part I Azure ML Workflow

**Data**

**Azure Machine Learning**

**Consumers**

**Cloud storage**
Azure Storage
Azure Table
Hive
etc.

**Local storage**
Upload data from PC…

*Data*

**ML Studio**
**(Web IDE)**

*Model*

**ML Web Services**
**(REST API Services)**

*API*

**Excel**

*Manage*

**Workspace:**
**Experiments**
Datasets
Trained models
Notebooks
Access settings

**Azure Marketplace**
**(Applications store)**

*API*

**Azure ML Gallery**
**(community)**

**Business Apps**

| Business problem | Modeling | Deployment | Business value |
|---|---|---|---|

*Reference*: **TechEd 2014**
**Conference**

# Part I Key Machine Learning Terminology and Concepts

1. Data exploration:
   - Process of gathering information about a large and often unstructured data set in order to find characteristics for focused analysis. Data mining refers to automated data exploration.

2. Descriptive analytics:
   - Process of analyzing a data set in order to summarize what happened. The vast majority of business analytics – such as sales reports, web metrics, and social network analysis – are descriptive.

3. Predictive analytics:
   - Process of building models from historical or current data in order to forecast future outcomes.

4. Supervised learning:
   - Algorithms are trained with labeled data – in other words, data comprised of examples of the answers wanted. For instance, a model that identifies fraudulent credit card use would be trained from a data set in which data points indicating known fraudulent and valid charges were labeled.

5. Unsupervised learning:
   - Is used on data with no labels, and the goal is to find relationships in the data. For instance, you might want to find groupings of customer demographics with similar buying habits.

# Part I Key Machine Learning Terminology and Concepts

1. Machine learning model
   - **Abstraction of the question you are trying to answer** or the outcome you want to predict. Models are trained and evaluated from existing data.

2. Training from data
   - In Azure Machine Learning, a model is built from an algorithm module that processes training data and functional modules, such as a scoring module.
   - In supervised learning, if you're training a fraud detection model, you'll use a set of transactions that are labeled as either fraudulent or valid. You'll split your data at randomly, and use part to train the model and part to test or evaluate the model.

3. Evaluation data
   - Once you have a trained model, evaluate the model using the remaining test data. You use data you already know the outcomes for, so that you can tell whether your model predicts accurately.

Azure Machine Learning Workflow

Feedback Loop   Create Model   Feedback Loop

Test/Use Model

Data
Analyze
Cleanse

Evaluate Model

Deploy Model

# Part I Key Machine Learning Terminology and Concepts

1. Classification:
   - Organizing data points into categories based on a data set for which category groupings are already known.

2. Regression:
   - Predicting a continuous value based on independent variables, such as predicting the price of a car based on its year and make.

3. Clustering:
   - Partition items into homogeneous groups. Typically used to predict grouping classifications for a given variable.

4. Recommendation:
   - The Netflix contest: Build a better recommender system from Netflix data. Use the crowd's votes to complete the missing entries.

5. Anomaly Detection:
   - Forecast of a value or values from a machine learning model.

- Occam's Razor: The best models are simple models that fit the data well.

- William of Ockham, English friar, philosopher, and theologian (1287-1347) said that among hypotheses that predict equally well, we should choose the one with the fewest assumptions.

# Part II Azure Machine Learning Studio

1. **Projects:** Sets of related Experiments, Trained Models, Datasets, Transforms.

2. **Experiments**: Experiments that have been created, run, and saved as drafts. These include a set of sample experiments that ship with the service to help jumpstart your projects.

3. **Web Services**: A list of experiments that you have published as web services.

4. **Notebooks:** Jupyter notebooks that you have created.

5. **Datasets**: A list of sample datasets that ship with the product, and uploaded data. You can use these datasets to learn about Azure Machine Learning.

6. **Trained Models**: List of any trained models that you saved from your experiments.

7. **Settings**: Configure your account and resources. Invite other users to share your workspace in Azure Machine Learning.

# Part II Azure Machine Learning Studio

With Azure Machine Learning Studio, you can

a.  Create predictive models in Machine Learning Studio, a browser-based tool, by dragging, dropping, and connecting modules.

b.  Use a large library of Machine Learning algorithms and modules in Machine Learning Studio to jump-start your predictive models.

c.  Choose from a library of sample experiments, R and Python packages, and best-in-class algorithms from Microsoft businesses like Xbox and Bing. Extent Studio modules with your own custom R and Python scripts.

d.  In Cortana Intelligence Gallery you can try analytics solutions authored by others or contribute your own using Azure services including Machine Learning, HDInsight (Hadoop), Stream Analytics, and Data Lake Analytics, as well as Azure big data stores and data management services.

# Part II First Experiment with Azure ML

1. **Goal**
   - Create a linear regression model that predicts the price of an automobile based on different variables such as make and technical specifications.
2. **Steps**
   a. **Create a model**
      - **Step 1: Get data**
      - **Step 2: Preprocess data**
      - **Step 3: Define features**
   b. **Train the model**
      - **Step 4: Choose and apply a learning algorithm**
   c. **Score and test the model**
      - **Step 5: Predict new automobile prices**

# Part II First Experiment with Azure ML



Source: http://0xCode.in/azure-ml-for-data-scientist
This work is licensed under a Creative Commons Attribution 4.0 International License

# Part II First Experiment with Azure ML

# Part II First Experiment with Azure ML

**Step 1. Get Data**

a. Start a new experiments.
b. Type automobile in the search box.
c. Drag the dataset to the experiment canvas.
d. Check the data with Visualize.

# Part II First Experiment with Azure ML

**Step 2. Preprocess Data**
a. Type project columns.
b. Click Launch column selector. Select normalized-losses to remove.
c. Drag the Clean Missing Data. Select remove for missing row.

# Part II First Experiment with Azure ML

Step 3. Define features
a.  Drag Project Columns.
b.  Launch column selector.
c.  Select no columns for Begin with and then select include and column names.

# Part II First Experiment with Azure ML

**Step 4. Choose and apply a learning algorithm**
a. **Split data (80:20).**
b. **Select learning algorithm (Linear Regression).**
c. **Select Train Model with price column.**
d. **Run the experiment.**

# Part II First Experiment with Azure ML

**Step 5. Predict new automobile prices**

a. **Find and drag Score Model.**

b. **Run the experiment.**

c. **Visualize Score Model.**

d. **Evaluate Model.**

# Part II Azure ML Algorithm Cheat Sheet



## Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.

### ANOMALY DETECTION
- **One-class SVM** — >100 features, aggressive boundary
- **PCA-based anomaly detection** — Fast training

*Finding unusual data points*

### CLUSTERING
- **K-means**

*Discovering structure*

### MULTI-CLASS CLASSIFICATION
- Fast training, linear model — **Multiclass logistic regression**
- Accuracy, long training times — **Multiclass neural network**
- Accuracy, fast training — **Multiclass decision forest**
- Accuracy, small memory footprint — **Multiclass decision jungle**
- Depends on the two-class classifier, see notes below — **One-v-all multiclass**

*Three or more*

*Predicting categories*

**START**

*Two*

### REGRESSION
- **Ordinal regression** — Data in rank ordered categories
- **Poisson regression** — Predicting event counts
- **Fast forest quantile regression** — Predicting a distribution
- **Linear regression** — Fast training, linear model
- **Bayesian linear regression** — Linear model, small data sets
- **Neural network regression** — Accuracy, long training time
- **Decision forest regression** — Accuracy, fast training
- **Boosted decision tree regression** — Accuracy, fast training, large memory footprint

*Predicting values*

### TWO-CLASS CLASSIFICATION
- **Two-class SVM** — >100 features, linear model
- **Two-class averaged perceptron** — Fast training, linear model
- **Two-class logistic regression** — Fast training, linear model
- **Two-class Bayes point machine** — Fast training, linear model

- Accuracy, fast training — **Two-class decision forest**
- Accuracy, fast training, large memory footprint — **Two-class boosted decision tree**
- Accuracy, small memory footprint — **Two-class decision jungle**
- >100 features — **Two-class locally deep SVM**
- Accuracy, long training times — **Two-class neural network**

# Part II Azure ML Algorithm Cheat Sheet

# Part III Azure ML Hands-On: Classification

1. Supervised learning
   a. **Classification**: used for predicting responses that can have just a few known values, such as "married," "single," or "divorced," based on the other columns in the dataset.
   b. **Regression**: predict one or more continuous variables, such as profit or loss, based on other columns in the dataset.
      a. **Simple linear regression**: A single key variable is predicted based on one or more other variables in the dataset.
      b. **Multiple linear regression**: More than one key variable is predicted based on one or more other variables in the dataset.
      c. **Multivariate linear regression**: multiple correlated dependent key variables are predicted, rather than a single key variable.

# Part III Azure ML Hands-On: Classification

1. **Goal:** To predict whether a person's income exceeds $50,000 per year based on his demographics or census data.

2. **Features:**
   - Age, Workclass, Fnlwgt, Education, Education-num, Marital-status, Occupation, Relationship, Race, Sex, Capital-gain, Capital-loss, Hours-per-week, Native-country, Income.

# Part III Azure ML Hands-On: Classification

## Step 1: Get the data

# Part III Azure ML Hands-On: Classification

## Step 2: Visualize the data

# Part III Azure ML Hands-On: Classification

## Step 2: Visualize the data

# Part III Azure ML Hands-On: Classification

## Step 3: Split the data

# Part III Azure ML Hands-On: Classification

Step 4: Train Model. Choose Two-Class Boosted Decision Tree. The income will be dependent.

# Part III Azure ML Hands-On: Classification

**Step 5: Score Model with Scored Labels and Probabilities.**

# Part III Azure ML Hands-On: Classification

**Step 6: Evaluate Model. Check the ROC curve.**

# Part III Azure ML Hands-On: Classification

**Step 7: Publish as a Web service.**

# Part III Azure ML Hands-On: Classification

## Step 8: Test Web service.

# Part III Azure ML Hands-On: Clustering

1. **K-Means Clustering** module: The K-means method finds a specified number of clusters for a set of D-dimensional data points. It starts an initial set of K centroids, and then uses algorithms to iteratively refine the locations of the centroids. The algorithm terminates when the centroids stabilize or when a specified number of iterations are computed.
2. **Train Clustering Model** module: This module takes an untrained clustering model, such as that produced by the K-Means Clustering module, and an unlabeled data set. It returns a trained clustering model that can be passed to the Assign to Clusters module. It also returns labels for the training data.
3. **Assign to Clusters** module: This module takes a trained clustering model, produced by the Train Clustering Model module, and an unlabeled data set. The module then returns the cluster assignments (indexes) for the input data.

# Part III Azure ML Hands-On: Grouping wholesale customers

1. **Datasets**
   - **Wholesale customers dataset from the UCI Machine Learning Repository located at http://mlr.cs.umass.edu/ml/datasets/Wholesale+customers**

2. **It includes the annual spending on the following product categories**
   a. **Channel: Retail channel types such as hotel/restaurant/café.**
   b. **Region: Customer region code.**
   c. **Fresh: Annual spending on fresh products.**
   d. **Milk: Annual spending on milk products.**
   e. **Grocery: Annual spending on grocery products.**
   f. **Frozen: Annual spending on frozen products.**
   g. **Detergents-Paper: Annual spending on detergents and paper products**
   h. **Delicatessen: Annual spending on delicatessen products.**

# Part III Azure ML Hands-On: Grouping wholesale customers

**Step 1: Get the data from URL and explore the data.**

# Part III Azure ML Hands-On: Grouping wholesale customers

Step 2: Choose the K-Means Clustering Model and Train the Clustering Model.

# Part III Azure ML Hands-On: Grouping wholesale customers

**Step 3: Setup the K-Means Clustering properties.**



- **Euclidean: distance between two points is the length of line segment connecting them.**
- **Cosine: Measure of similarity between two vectors of an inner space that measure the cosine of the angles between them.**

# Part III Azure ML Hands-On: Grouping wholesale customers

Step 4: Visualize the Training Results.

# Part III Azure ML Hands-On: Grouping wholesale customers

Step 5: Split the data and set up for testing.

# Part III Azure ML Hands-On: Grouping wholesale customers

**Step 6: Publish as a Web Service.**

# Part III Azure ML Hands-On: Grouping wholesale customers

**Step 7: Test the Web Service.**

## wholesale customers

DASHBOARD    CONFIGURATION

General

Parent Experiment

Wholesale Customers 10 - transform to Web Svc 4

Description

No description provided for this web service.

API key

tZgp7V9yaB9QnekxFaZqc3rCd8pqltveERJ/EWCa6wOWywqlLy8jLb9ksQ9YuKuUwDhtvLMlGLZrdWmH060TEA==

Default Endpoint

| URL | TYPE | LAST UPDATED | ↓ | TEST |
|-----|------|--------------|---|------|
| API help page | REQUEST/RESPONSE | 3/18/2015 12:38:05 AM | | Test |
| API help page | BATCH EXECUTION | 3/18/2015 12:38:05 AM | | |

Additional endpoints

Number of additional endpoints created for this web service: 0

Manage endpoints in Azure management portal

# Part III Azure ML Hands-On: Recommendation

1. **Create a model**: A model is a container of your usage data, catalog data, and the recommendation model.

2. **Import catalog data**: This is an optional step. A catalog contains metadata information on the items, if you do not upload catalog data, the recommendation's services will learn about your catalog from the usage data.

3. **Import usage data**
   - By uploading a file that contains the usage data.
   - By sending data acquisition events. Usually you upload a usage file to be able to create an initial recommendation model (bootstrap). You would use this usage until the system gathers enough data by using the data acquisition format.

4. **Building a recommendation model**:
   - This is an asynchronous operation in which the recommendation system takes all the usage data and creates a recommendation model. This operation can take several minutes or several hours depending on the size of the data and the build configuration parameters. When triggering the build, you will get a build ID. Use it to check when the build process has ended before starting to consume recommendations.

# Part III Azure ML: Building the restaurant ratings recommender

1. Data sets
    a. **Restaurant features**: This includes information about each restaurant, such as placeID, location, name, address, state, whether alcohol is served, smoking policy, dress code, accessibility, and price.
    b. **Restaurant customers**: This includes a wide variety of personal attributes include userID, smoker, drink level, dress preference, marital status, birth year, interests, personality traits, religion, favorite color, weight, and budget.
    c. **Restaurant ratings**: This includes key fields like userID, placeID, and rating.



Restaurant

▲ Saved Datasets

Restaurant customer data

Restaurant feature data

Restaurant ratings

01 - Restaurant Ratings Experiment

Restaurant feature data    Restaurant customer data    Restaurant ratings

# Part III Azure ML: Building the restaurant ratings recommender

2. **Select Columns**
   a. We can filter out unnecessary columns in the Restaurant Feature Data and Restaurant Customer Data datasets.
   b. **Restaurant** features: placeID, latitude, longitude, price.
   c. **Customer** Data: userID, latitude, longitude, interest, personality.

3. **Split Module**
   - Specify that half of this dataset is to be used for training and the other half of scoring the new recommendation model.

## Select columns

☐ Allow duplicates and preserve column order in selection

**Begin With**  | No columns ▾

| Include ▾ | column names ▾ |

placeID ✖  latitude ✖  longitude ✖  price ✖

## Select columns

☐ Allow duplicates and preserve column order in selection

**Begin With**  | No columns ▾

| Include ▾ | column names ▾ |

userID ✖  latitude ✖  longitude ✖  interest ✖
personality ✖

Properties

◢ **Split**

Splitting mode

Recommender Split ▾

Fraction of training-onl...

0.5

Fraction of test user rati...

0.25

Fraction of cold users

0

Fraction of cold items

0

Fraction of ignored users

0

Fraction of ignored items

0

- 44 -

4. **Train Matchbox Recommender** module
    a. **Training dataset of user-item-rating triples:** Ratings of items by users, expressed as a triple (Use, Item, Rating).
    b. **Training dataset of user features:** Dataset containing features that describe users.
    c. **Training dataset of item features:** Dataset containing features that describe items.

# Part III Azure ML: Building the restaurant ratings recommender

5. **Score Matchbox**
    a. Trained Matchbox recommender.
    b. Dataset to score.
    c. User features: Dataset containing features that describe users.
    d. Item features: Dataset containing features that describe items.

5. **Score Matchbox**
    a. **Rating Prediction: Predict the rating that a customer will give a particular restaurant.**
    b. **Item Recommendation: Predict which restaurants will be most highly rated by the user.**
    c. **Related Users: Predict which customers (users) are most like this customer.**
    d. **Related Items: Predict which restaurants (items) are most like this restaurant.**

◢ Score Matchbox Recommender

Recommender prediction kind

| Item Recommendation | ∨ |

Recommended item selection

| From Rated Items (for model evaluation) | ∨ |

Maximum number of items to recommend to a user

| 5 |

Minimum size of the recommendation pool for a single user

| 2 |

Evaluate Recommender    🔍

◢ ⊞ **Machine Learning**

   ◢ **Evaluate**

     Evaluate Recommender

◢ Score Matchbox Recommender

Recommender prediction kind

| Rating Prediction | ∨ |

# Part III Azure ML: Building the restaurant ratings recommender

## 5. Evaluate Recommender

# Part III Azure ML Exercises: Classification

1. **Two-class classification**
   a. **To classify Iris flowers based on their features.**
   b. **Two flower species (classes 0 and 1).**
   c. **Four features for each flower (sepal length, sepal width, petal length, and petal width)**

2. **Multi-class classification**
   1. **To classify a letter (class), given some attribute values extracted from the hand-written letter images.**
   2. **Twenty-six letters form twenty-six classes.**

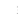# Part III Azure ML Exercises: Clustering

- **Clustering differs from classification in that the training dataset does not have ground-truth labels by itself.**

- **Group the training dataset instances into distinct clusters.**

- **During the training process, the model labels the entries by leaning the differences between their features.**

# Part III Azure ML Exercises: Regression

- To predict the price of a car based on its features including make, fuel type, body type, drive wheel, etc.
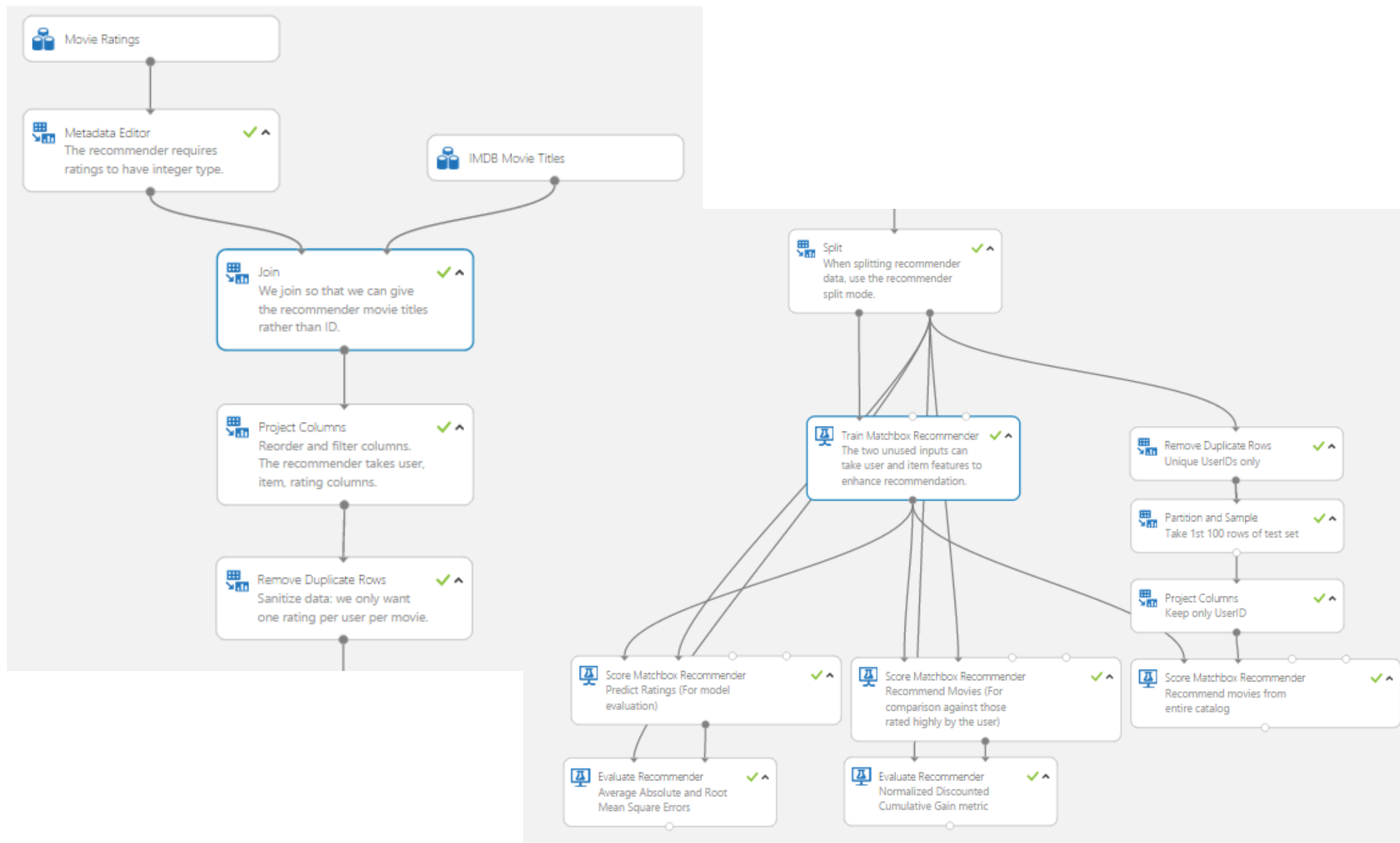
# Part III Azure ML Exercises: Recommendation

# Part IV Lessons and Resources

1. Data wrangling is important
   a. More time is spent on data wrangling than model building.
   b. Different sources, formats, schemas, missing values and noisy data.
   c. Data manipulation modules are very popular: Execute R script, SQL Transform, etc.

2. Azure ML Modeling
   a. Modeling depends on the Business Application Domain and Data.
   b. Feature Engineering is essential.
   c. Parameter Tuning is needed.
   d. Learn R or Python Script.

3. Resources
   a. Getting Started: https://studio.azureml.net
   b. Gallery: http://gallery.azureml.net/
   c. Site/ML Studio/Docs: http://azure.microsoft.com/en-us/services/machine-learning/
   d. Blog: http://blogs.technet.com/b/machinelearning/
   e. edX: Microsoft DAT203X Data Science and Machine Learning Essentials https://courses.edx.org/courses/course-v1:Microsoft+DAT203x+3T2015/info
   f. Microsoft Virtual Academy: https://mva.microsoft.com/

# Q and A