# Learning

# – An Introduction to Machine Learning in Python

PyData Chicago 2016
Chicago, The University of Illinois • August 26, 2016

Sebastian Raschka

# Links & Info

## Tutorial Material on GitHub:

https://github.com/rasbt/pydata-chicago2016-ml-tutorial

## Contact:

- o E-mail: mail@sebastianraschka.com

- o Website: http://sebastianraschka.com

- o Twitter: @rasbt

- o GitHub: rasbt

PyData
*Chicago 2016*

# Let's not stress!

This is an introductory tutorial, and we are here to learn!

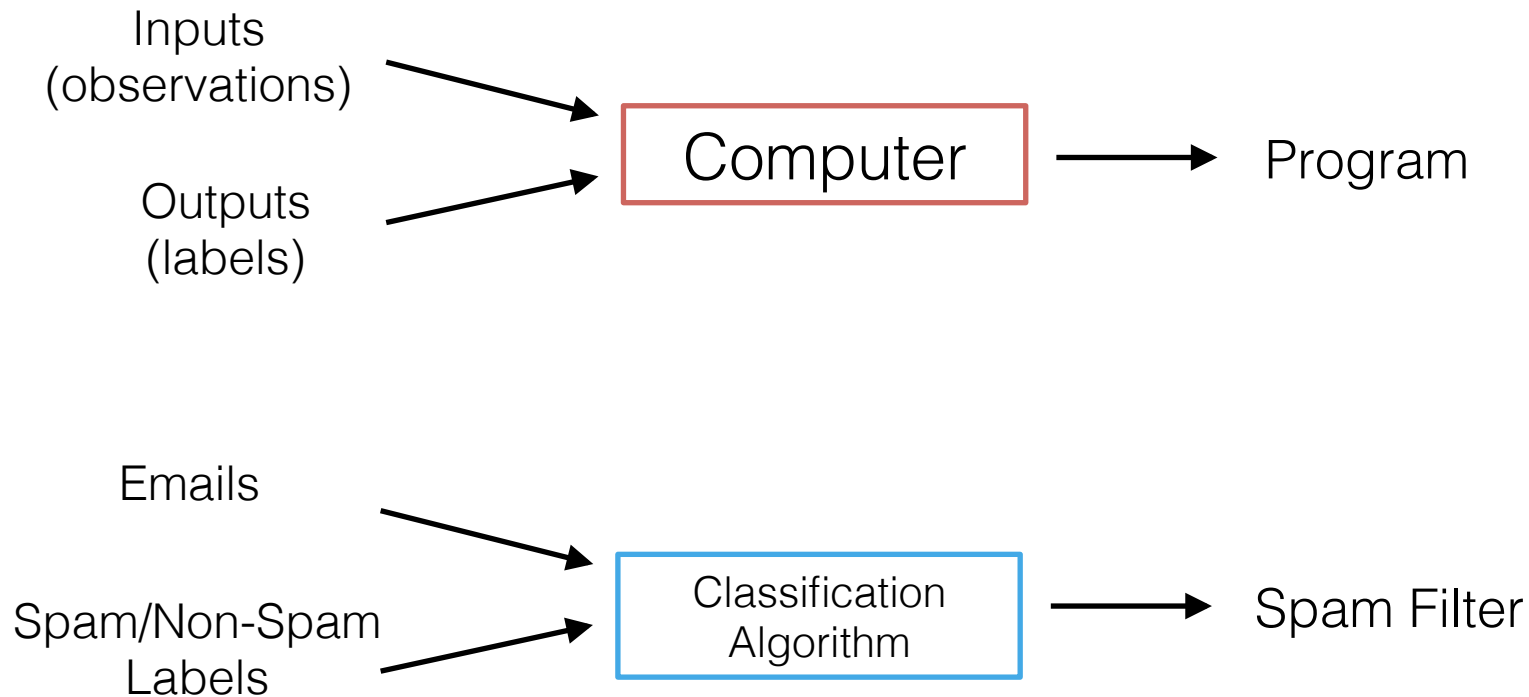## Please ask questions!

# What can machine learning do for us?



https://flic.kr/p/5BLW6G [CC BY 2.0]





https://commons.wikimedia.org/wiki/
File:Google_self_driving_car_at_the_Googleplex.jpg

Photo by Michasel Shick, CC BY-SA 4.0 lit

# What is machine learning?

PyData
*Chicago 2016*

# 3 types of learning

Supervised

Unsupervised

Reinforcement

PyData
*Chicago 2016*
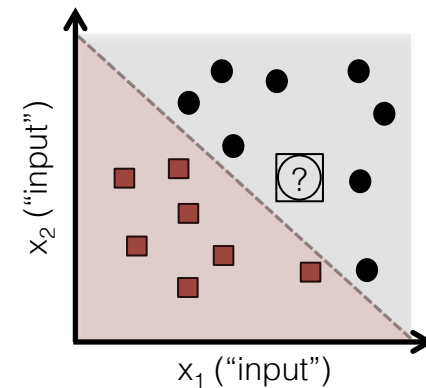
# Working with labeled data
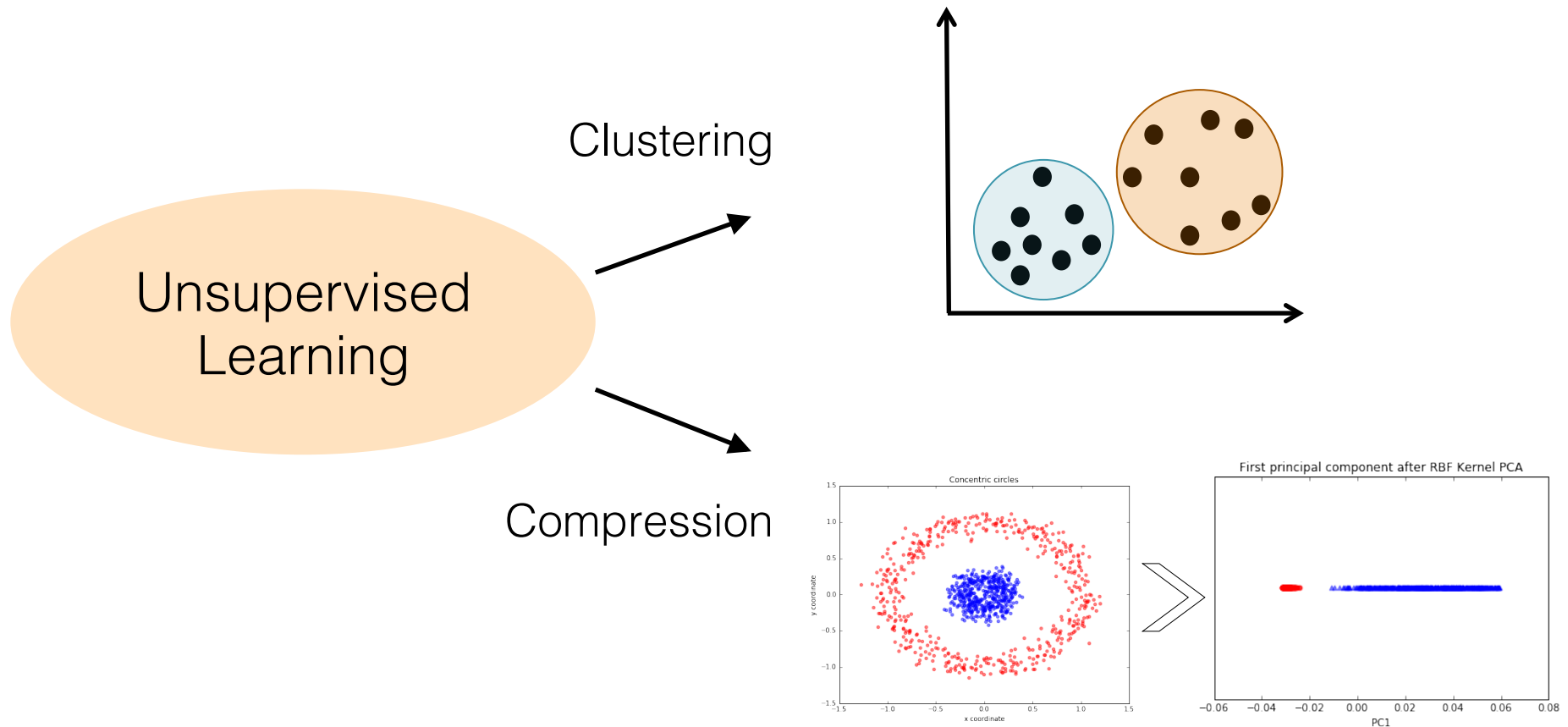


Supervised Learning
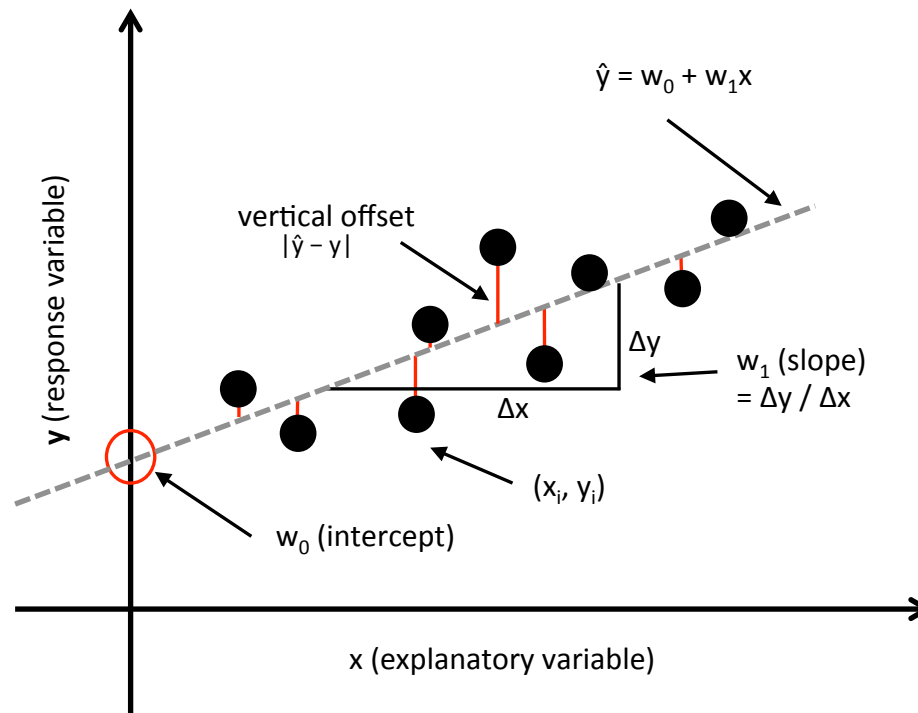
Regression

Classification

# Working with <u>unlabeled</u> data

# Topics

1. Introduction to Machine Learning
2. **Linear Regression**
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. Support Vector Machine Classifiers
7. Model Evaluation & Hyperparameter Tuning
8. Tree-based Methods
9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Simple Linear Regression

# Data representation

features (columns)

$$\mathbf{X}=$$

| $\mathbf{x_0}$ | $\mathbf{x_1}$ | ... | $\mathbf{x_m}$ |
|---|---|---|---|
| $x_{0,0}$ | $x_{0,1}$ | | |
| $x_{1,0}$ | $x_{1,1}$ | | |
| $x_{2,0}$ | $x_{2,1}$ | | |
| $x_{3,0}$ | $x_{3,1}$ | | |
| . | | | |
| . | | | |
| . | | | |
| $x_{n,0}$ | $x_{n,1}$ | ... | $x_{n,m}$ |

samples (rows)

$$\mathbf{y}=$$

| |
|---|
| $y_0$ |
| $y_1$ |
| $y_2$ |
| $y_3$ |
| . |
| . |
| . |
| $y_n$ |

# "Basic" Supervised Learning Workflow

# Coding Example
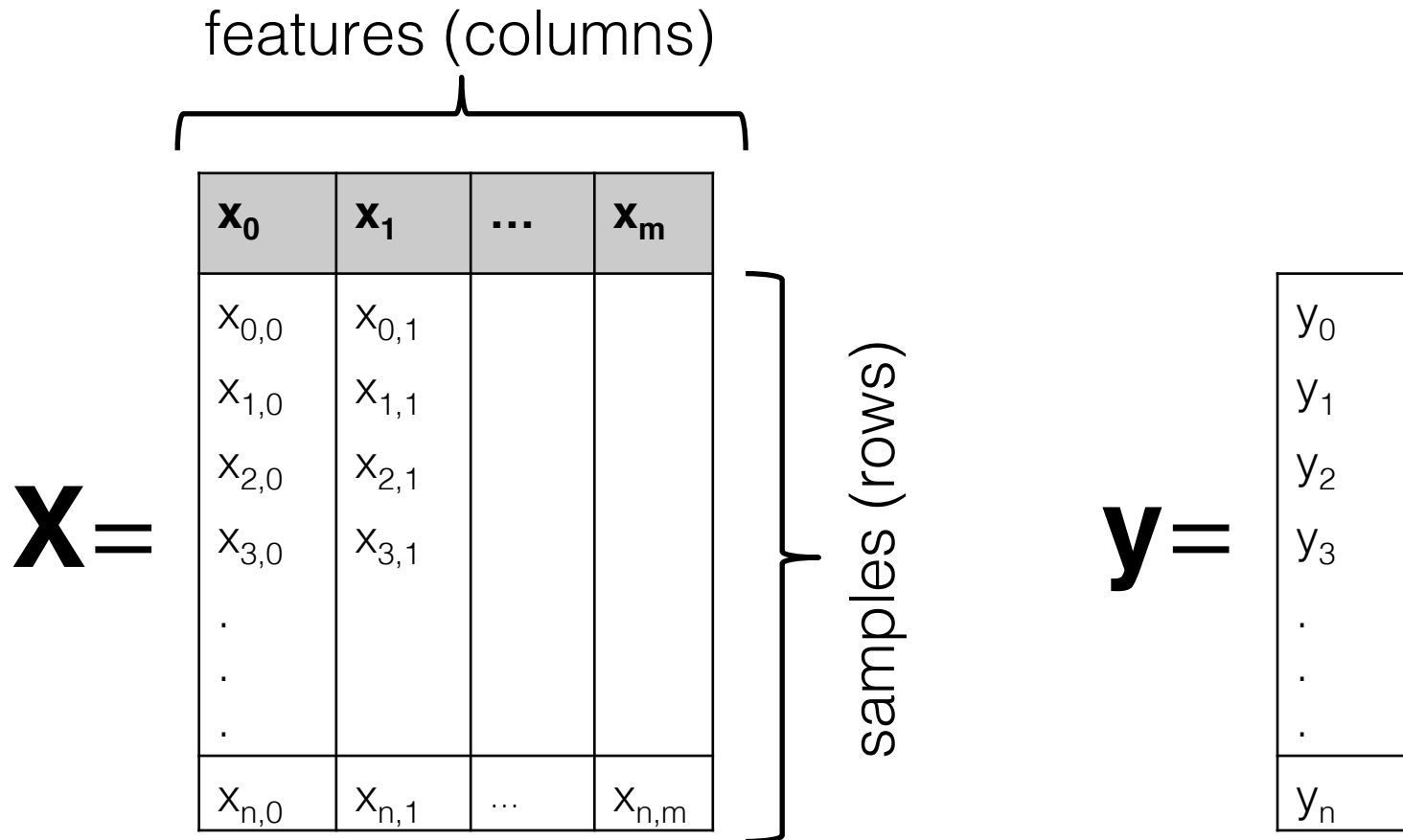
➡️ Jupyter Notebook

PyData
*Chicago 2016*

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. **Introduction to Classification**
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. Support Vector Machine Classifiers
7. Model Evaluation & Hyperparameter Tuning
8. Tree-based Methods
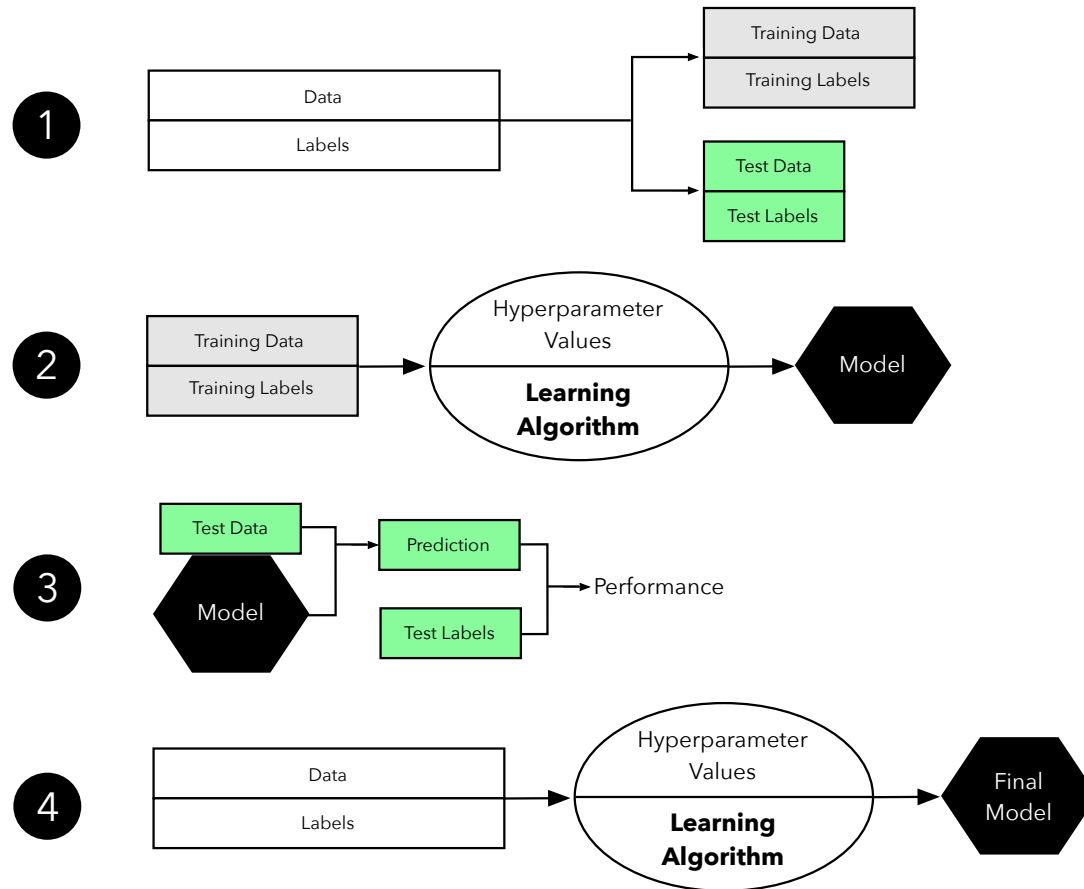9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Scikit-learn API

```python
class SupervisedEstimator(...):
    def __init__(self, hyperparam, ...):
        ...

    def fit(self, X, y):

        ...

        return self

    def predict(self, X):

        ...

        return y_pred

    def score(self, X, y):

        ...

        return score

    ...
```

PyData
*Chicago 2016*

# Iris dataset

Iris-Setosa

Iris-Setosa

Iris-Versicolor

PyData
*Chicago 2016*

# Iris dataset

features (columns)

| | sepal length [cm] | sepal width [cm] | petal length [cm] | petal width [cm] |
|---|---|---|---|---|
| 1 | 5.1 | 3.5 | 1.4 | 0.2 |
| 2 | 4.9 | 3.0 | 1.4 | 0.2 |
| 50 | 6.4 | 3.5 | 4.5 | 1.2 |
| | . | | | |
| | . | | | |
| | . | | | |
| 150 | 5.9 | 3.0 | 5.0 | 1.8 |

$\mathbf{X}=$

samples (rows)

$\mathbf{y}=$

| |
|---|
| setosa |
| setosa |
| versicolor |
| . |
| . |
| . |
| virginica |

petal

sepal

PyData
Chicago 2016
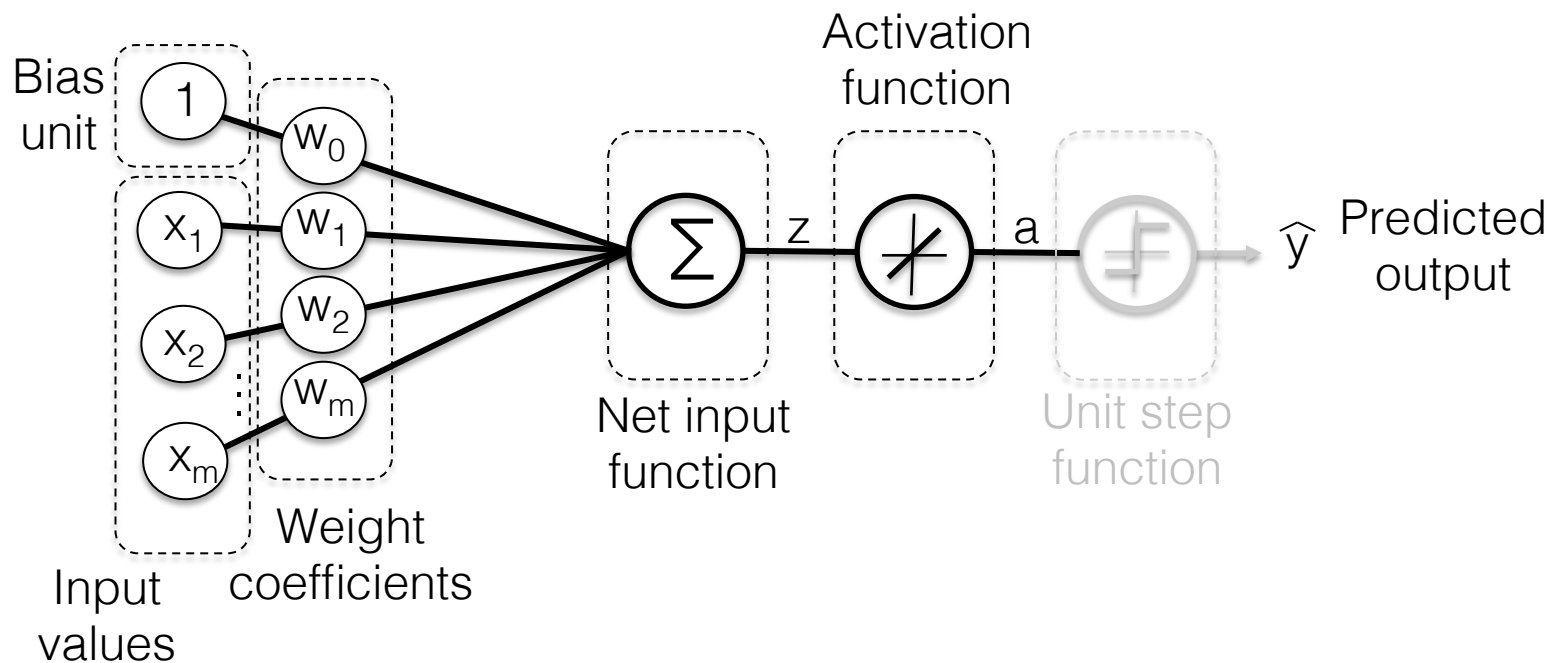
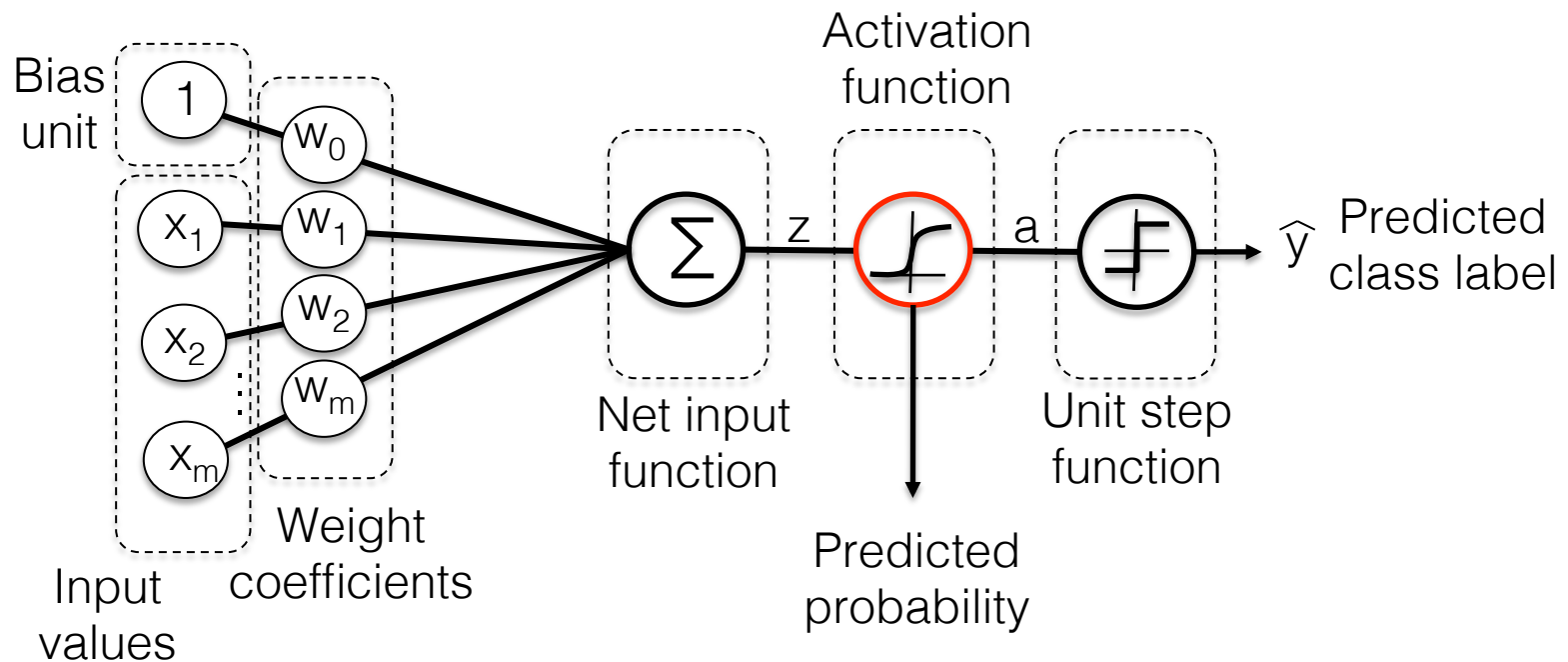# Note about non-stratified splits



- training set → 38 x Setosa, 28 x Versicolor, 34 x Virginica
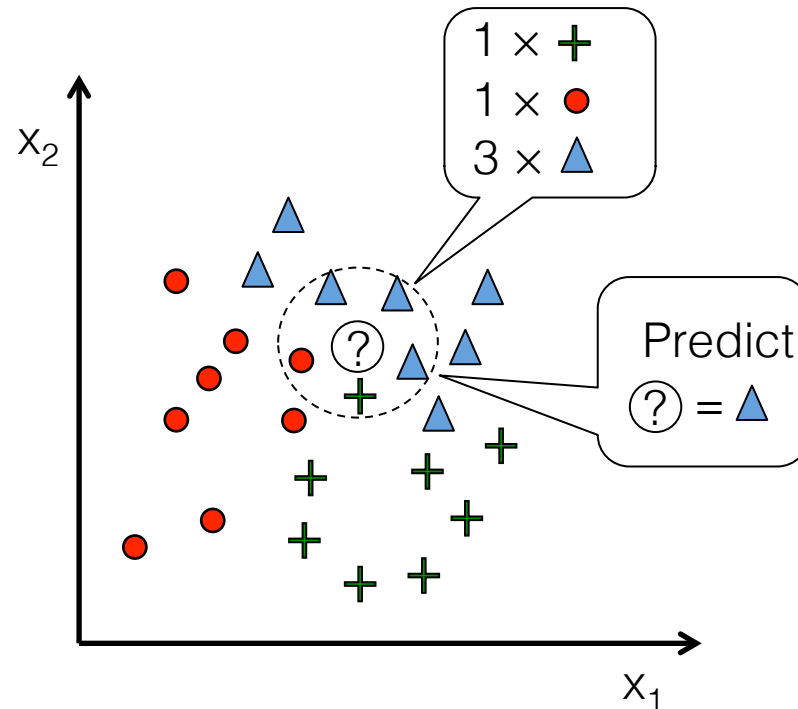- test set → 12 x Setosa, 22 x Versicolor, 16 x Virginica

# Linear Regression Recap

# Logistic Regression, a generalized linear model

# A "lazy learner:" K-Nearest Neighbors classifier

# Coding Example

➡️ Jupyter Notebook

PyData
Chicago 2016

# Topics

1. Introduction to Machine Learning

2. Linear Regression

3. Introduction to Classification

4. **Feature Preprocessing & scikit-learn Pipelines**

5. Dimensionality Reduction: Feature Selection & Extraction

6. Support Vector Machine Classifiers

7. Model Evaluation & Hyperparameter Tuning

8. Tree-based Methods

9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Scikit-learn API

```python
class UnsupervisedEstimator(...):
    def __init__(self, hyperparam, ...):

        ...

    def fit(self, X, y):

        ...

        return self

    def predict(self, X):

        ...

        return y_pred

    def score(self, X, y):

        ...

        return score

    ...
```

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. **Dimensionality Reduction Feature Selection & Extraction**
6. Support Vector Machine Classifiers
7. Model Evaluation & Hyperparameter Tuning
8. Tree-based Methods
9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. **Support Vector Machine Classifiers**
7. Model Evaluation & Hyperparameter Tuning
8. Tree-based Methods
9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. Support Vector Machine Classifiers
7. **Model Evaluation & Hyperparameter Tuning**
8. Tree-based Methods
9. Unsupervised Learning: Clustering

PyData
*Chicago 2016*

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. Support Vector Machine Classifiers
7. Model Evaluation & Hyperparameter Tuning
8. **Tree-based Methods**
9. Unsupervised Learning: Clustering
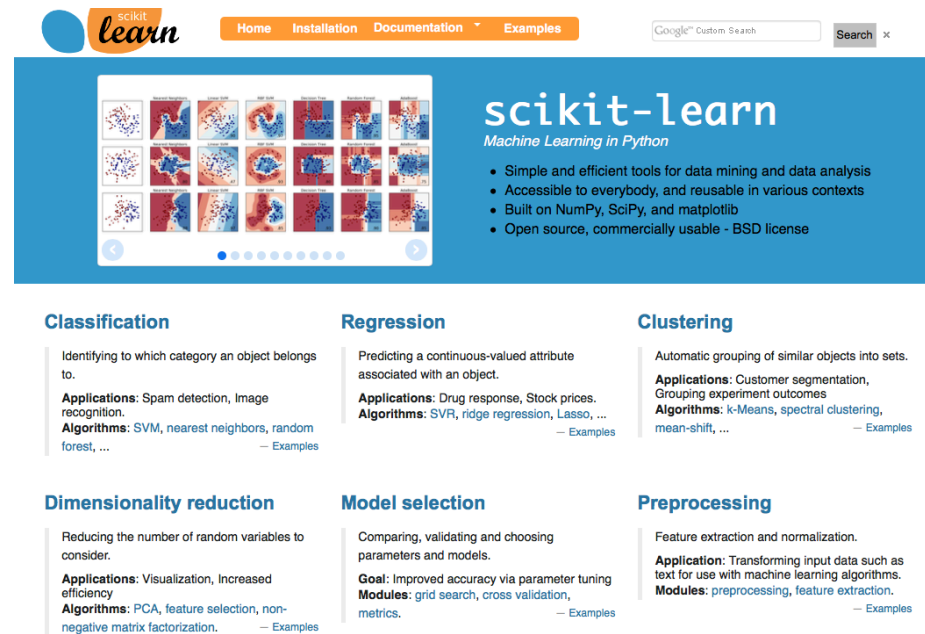
PyData
*Chicago 2016*

# Topics

1. Introduction to Machine Learning
2. Linear Regression
3. Introduction to Classification
4. Feature Preprocessing & scikit-learn Pipelines
5. Dimensionality Reduction: Feature Selection & Extraction
6. Support Vector Machine Classifiers
7. Model Evaluation & Hyperparameter Tuning
8. Tree-based methods
9. **Unsupervised Learning: Clustering**

PyData
Chicago 2016

# Topics

1. Introduction to Machine Learning

2. Linear Regression

3. Introduction to Classification

4. Feature Preprocessing & scikit-learn Pipelines

5. Dimensionality Reduction: Feature Selection & Extraction

6. Support Vector Machine Classifiers

7. Model Evaluation & Hyperparameter Tuning

8. Tree-based methods

9. Unsupervised Learning: Clustering

# Further Resources

Documentation:
http://scikit-learn.org

Mailing list:
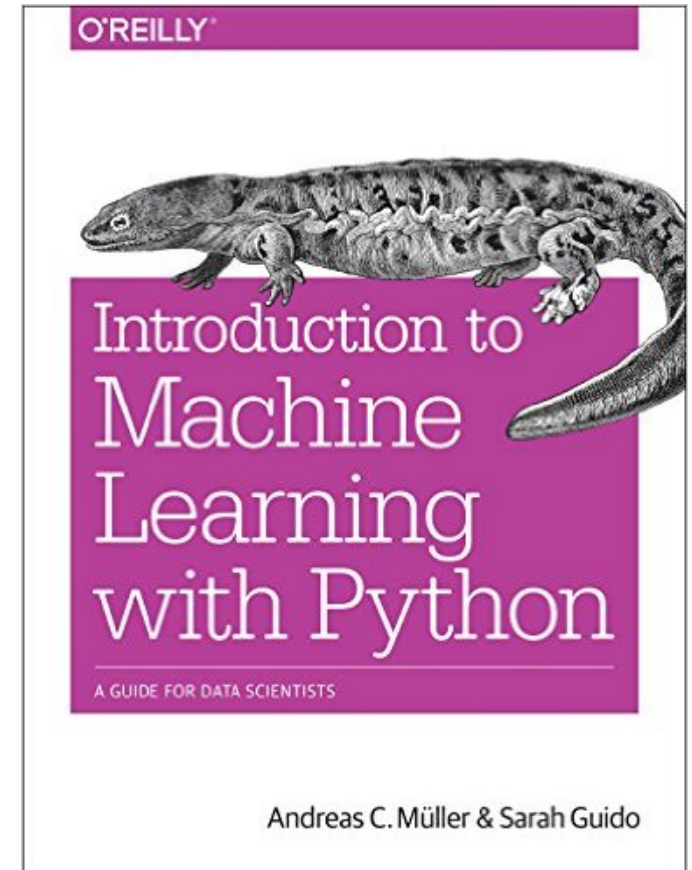https://mail.python.org/mailman/listinfo/scikit-learn

# Further Resources

*Great "math-free," practical guide to machine learning with scikit-learn*

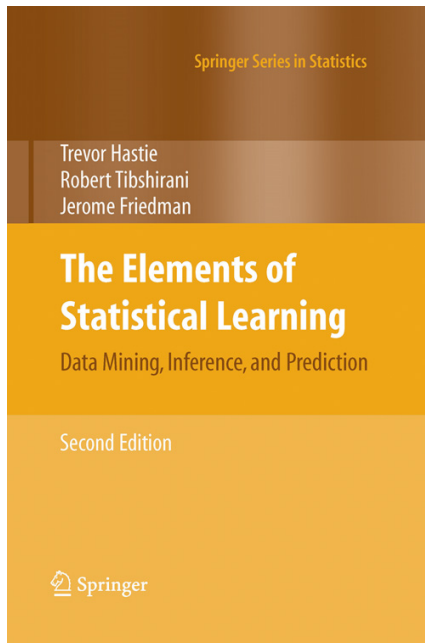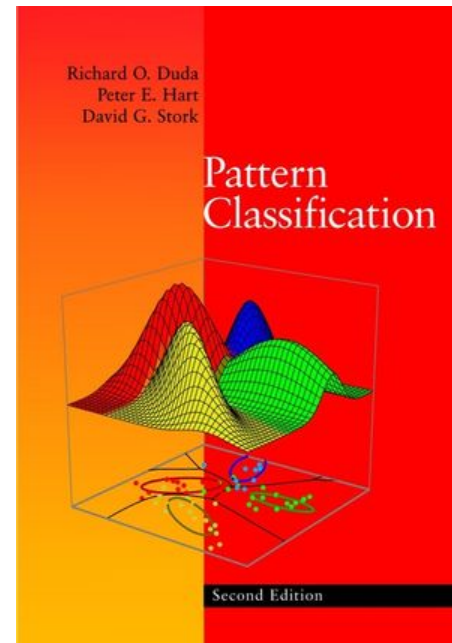By Andreas Mueller (scikit-learn core developer) and Sarah Guido

http://shop.oreilly.com/product/0636920030515.do

Estimated release: October 20, 2016

# Further Resources

My favorite machine learning "math & theory" books





http://statweb.stanford.edu/~tibs/ElemStatLearn/ (free PDF)

http://www.wiley.com/WileyCDA/WileyTitle/productCd-0471056693.html

# Further Resources

My own book, math, from-scratch code, and practical scikit-learn code:

GitHub repository:

https://github.com/rasbt/python-machine-learning-book

Amazon link:

https://www.amazon.com/Python-Machine-Learning-Sebastian-Raschka/dp/1783555130/