



Module 3

Alireza Samar

Who Am I?



Alireza Samar @alirezasmr

Founder, CTO at Appoint AI

Graduate Research Assistant at UTM MLDS

Curator at Machine Learning Weekly - mlweekly.com



MLDS

What We've Learned in Previous Modules

- Anaconda Package Manager
- Install and run code on Jupyter
- Version Controlling Concept
- Git
- GitHub and GitHub Desktop



What We've Learned in Previous Modules

- Statistics
- Linear Algebra
- Optimization
- Bayes Rule
- Maximum Likelihood, Gradient Descent, ...



What We've Learned in Previous Modules

- Python Basics (Syntax, Arrays, Loops, Functions and etc)
- Numpy
- Matplotlib
- Exploratory Data Analysis
- SciPy - Scientific Computing Toolkit

Workshop Materials

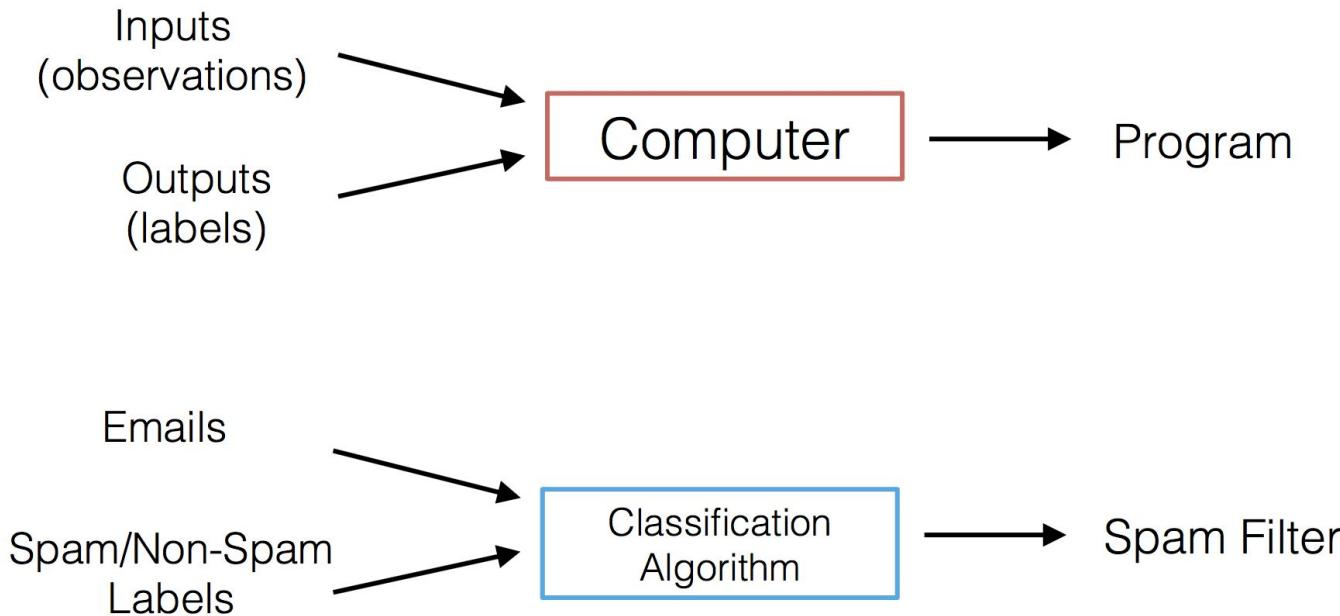
GitHub: **UTMMLDS**

<https://github.com/utmmlds>

The logo consists of the letters "MLDS" in a bold, white, sans-serif font. The letters are partially cut off at the right edge of the slide.

MLDS

What is Machine Learning?



3 Types of Learning

- Supervised
- Unsupervised
- Reinforcement

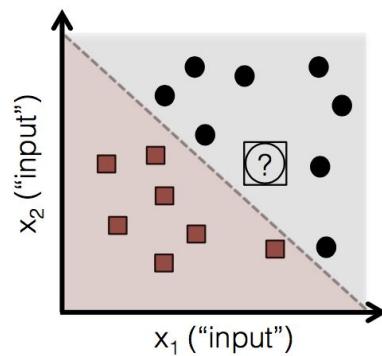
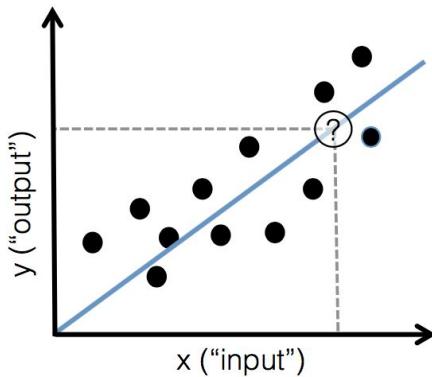


Working with Labeled Data



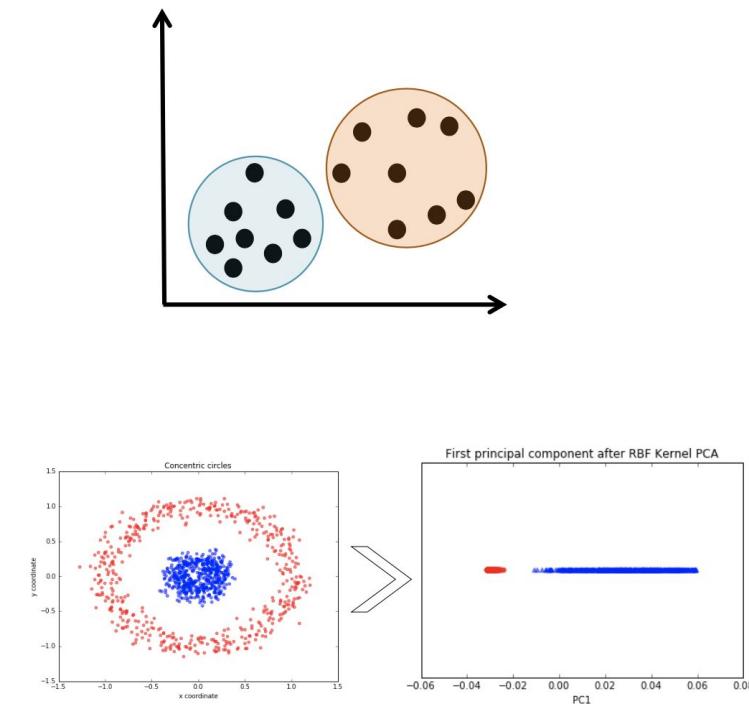
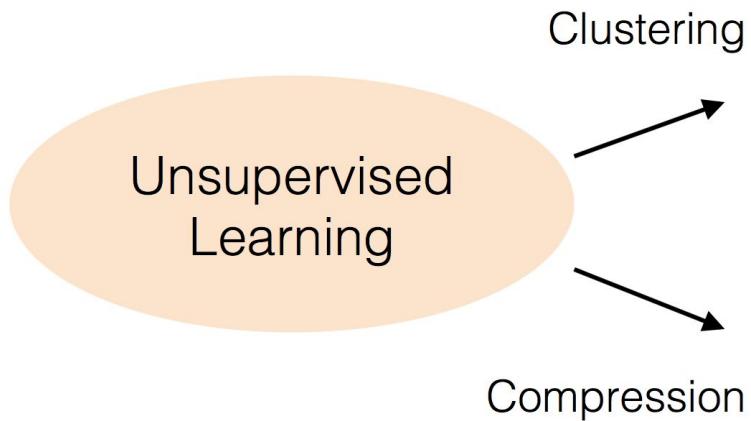
Regression

Classification



MLDS

Working with **Unlabeled** Data

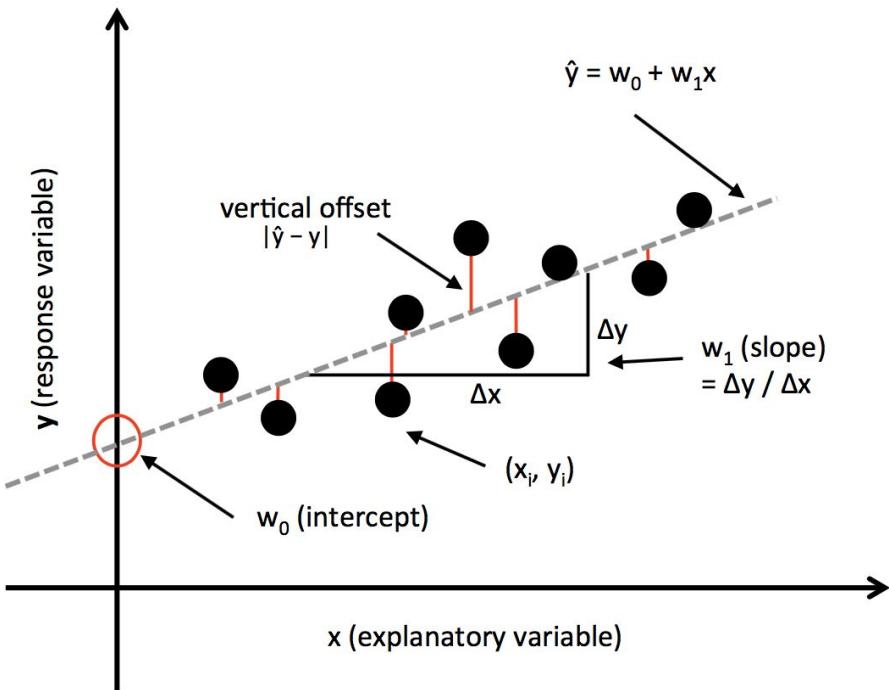


MLDS

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. Model Evaluation & Hyperparameter Tuning

Simple Linear Regression



Data Representation

X =

features (columns)			
x_0	x_1	...	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)

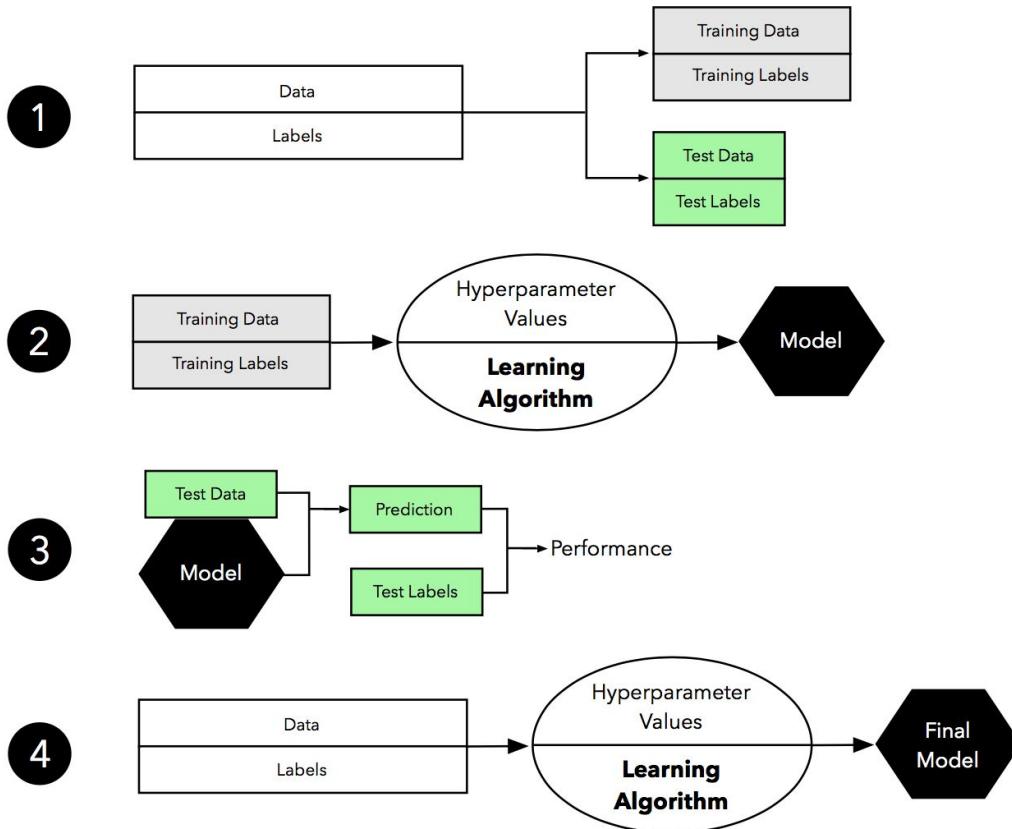
y =

y_0
y_1
y_2
y_3
.
.
y_n



MLDS

“Basic” Supervised Learning Workflow



Coding Time!



MLDS

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. Model Evaluation & Hyperparameter Tuning

Scikit-learn API

```
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
    ...

```



Iris Dataset

Iris-Setosa



Iris-Setosa



Iris-Versicolor



MLDS

Iris Dataset

X=

	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	

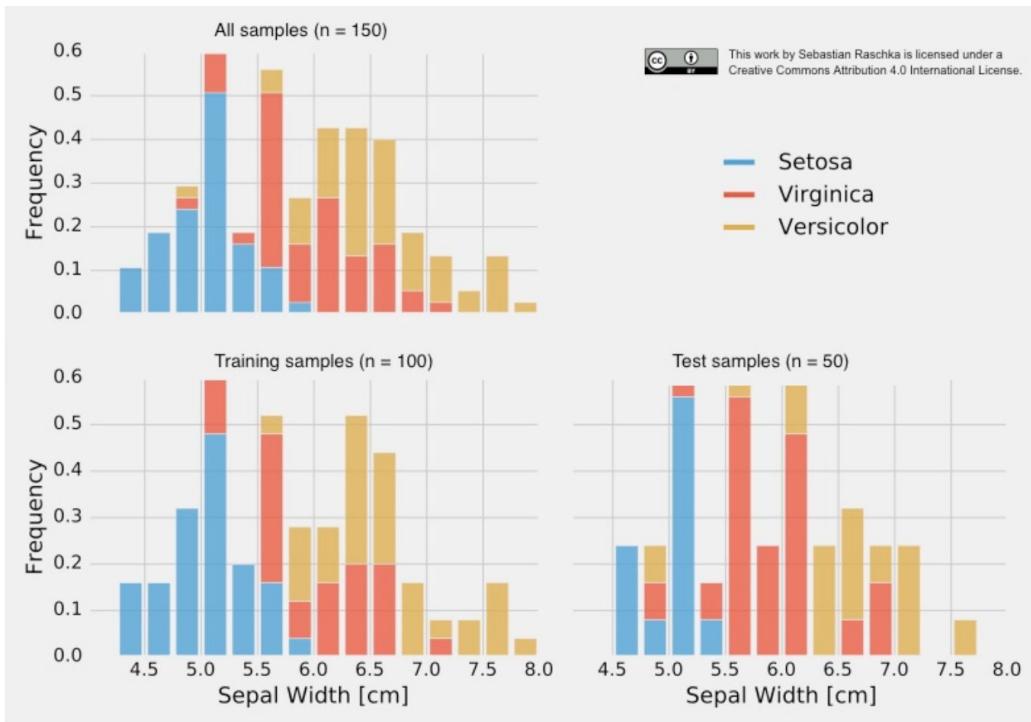
samples (rows)



y=

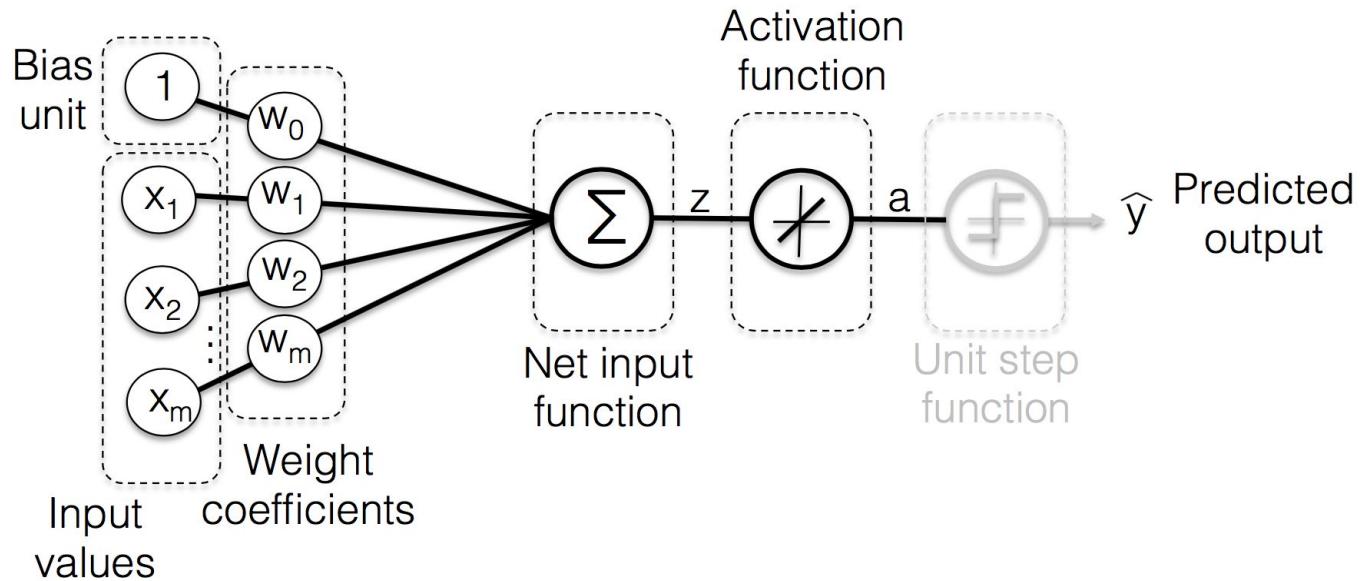
setosa
setosa
versicolor
.
.
.
virginica

Note about Non-Stratified Splits

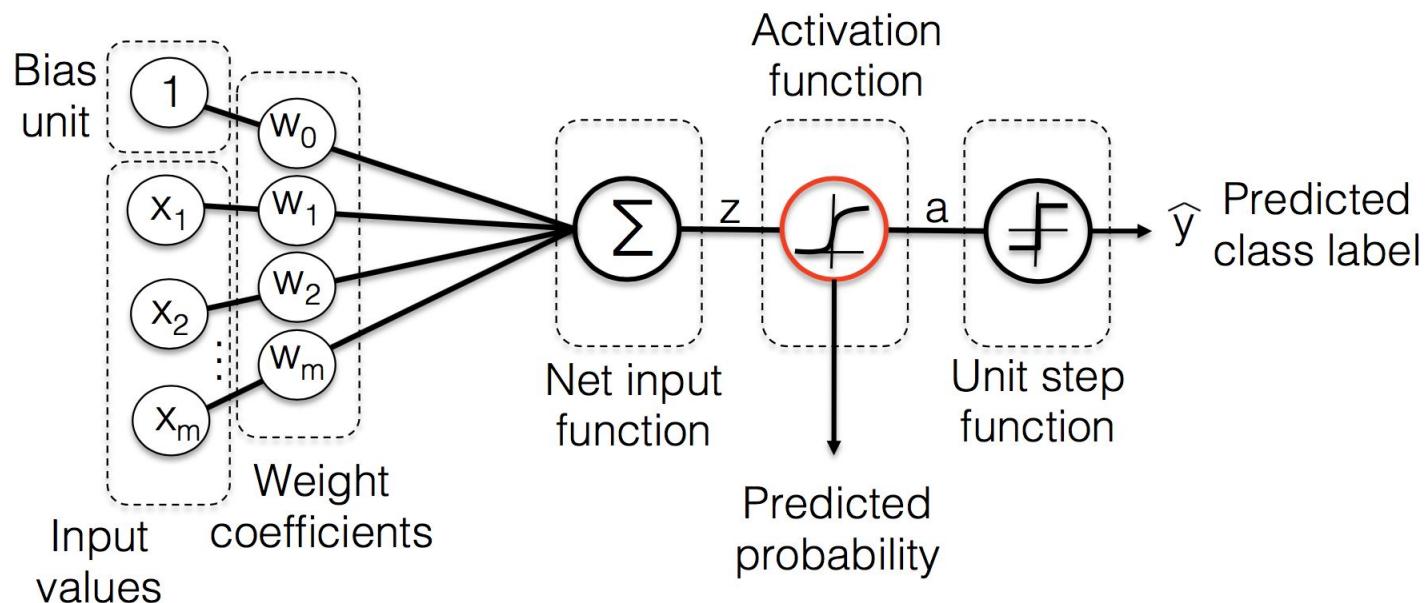


- training set → 38 x Setosa, 28 x Versicolor, 34 x Virginica
- test set → 12 x Setosa, 22 x Versicolor, 16 x Virginica

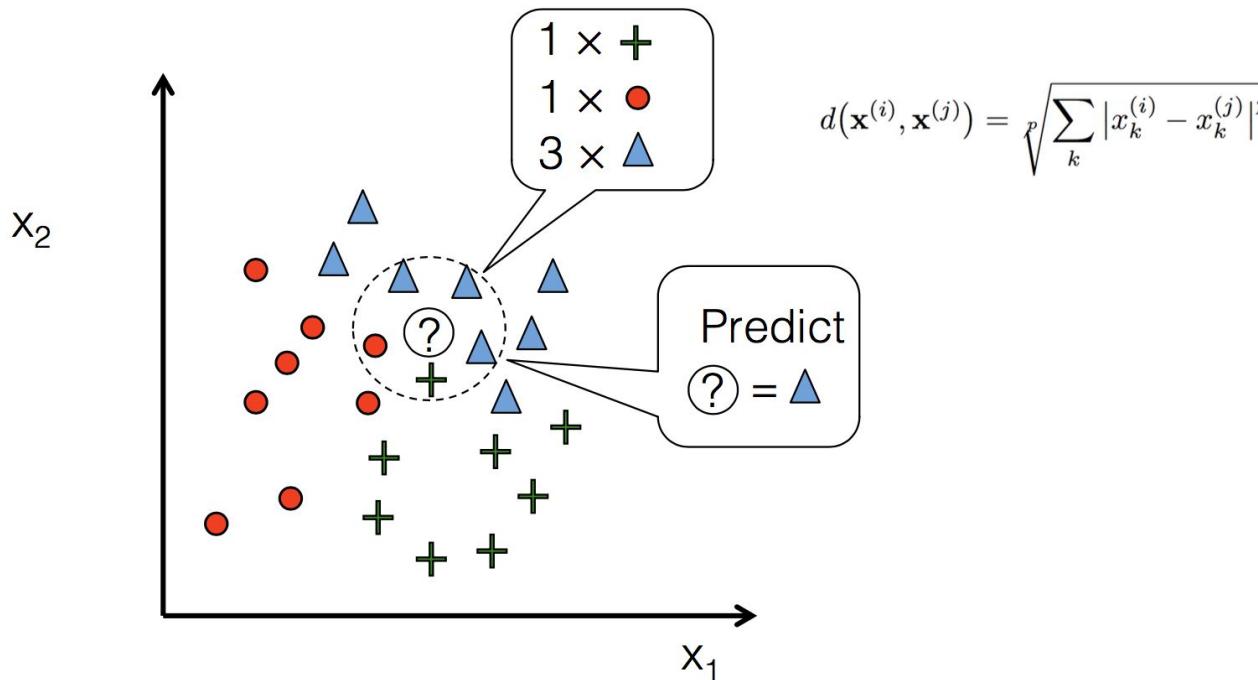
Linear Regression Recap



Logistic Regression, a Generalized Linear Model



A “Lazy Learner:” K-Nearest Neighbors Classifier



$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sqrt[p]{\sum_k |x_k^{(i)} - x_k^{(j)}|^p}$$

Coding Time!

The logo consists of the letters "MLDS" in a bold, white, sans-serif font. The letter "M" is slightly taller than the others. The letters are positioned above a faint, light-gray network graph of interconnected dots and lines.

MLDS

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. Model Evaluation & Hyperparameter Tuning

Categorical Variables

color	size	price	class label
red	M	\$10.49	0
blue	XL	\$15.00	1
green	L	\$12.99	1

Encoding Categorical Variables



Feature Normalization

Min-max scaling

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Z-score standardization

$$z = \frac{x - \mu}{\sigma}$$

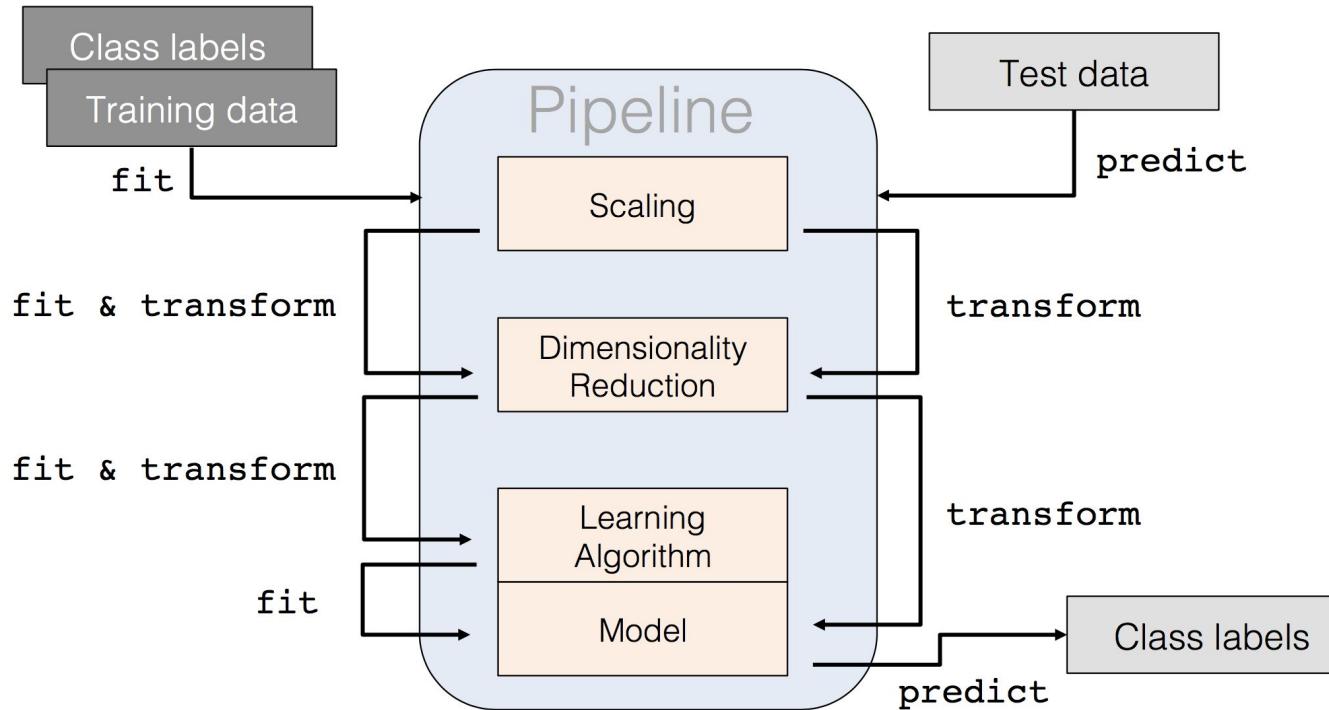
feature	minmax	z-score
1.0	0.0	-1.46385
2.0	0.2	-0.87831
3.0	0.4	-0.29277
4.0	0.6	0.29277
5.0	0.8	0.87831
6.0	1.0	1.46385



Scikit-learn API

```
class UnsupervisedEstimator(...):
    def __init__(self, ...):
        ...
    def fit(self, X):
        ...
        return self
    def transform(self, X):
        ...
        return X_transf
    def predict(self, X):
        ...
        return pred
```

Scikit-learn Pipelines



Coding Time!

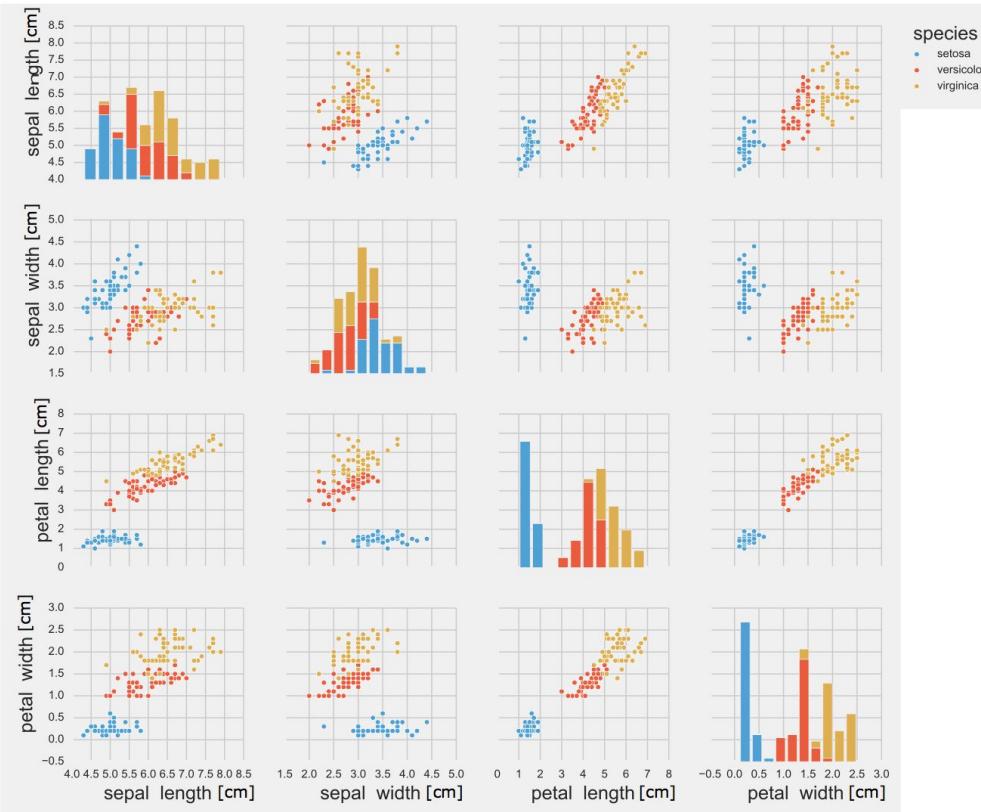


MLDS

Topics

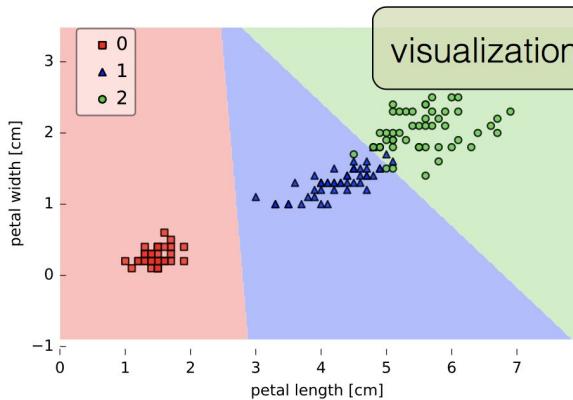
- 
1. Introduction to Machine Learning
 2. Linear Regression
 3. Introduction to Classification
 4. Feature Preprocessing & scikit-learn Pipelines
 - 5. Dimensionality Reduction: Feature Selection & Extraction**
 6. Model Evaluation & Hyperparameter Tuning

Dimensionality Reduction – why?



MLDS

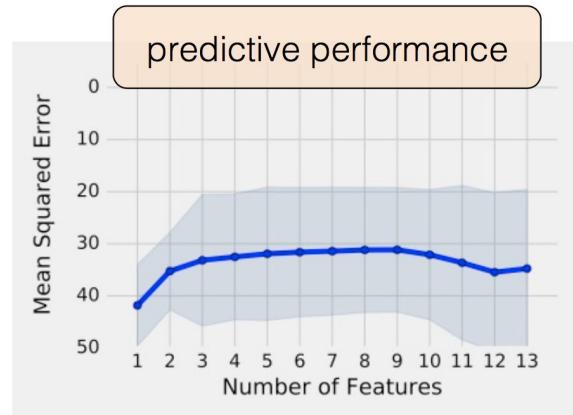
Dimensionality Reduction – why?



visualization & interpretability



storage & speed



predictive performance



MLDS

Recursive Feature Elimination

available features:

[f1 f2 f3 f4]

[w1 w2 w3 w4]

fit model, remove lowest weight, repeat

[w1 w2 w4]

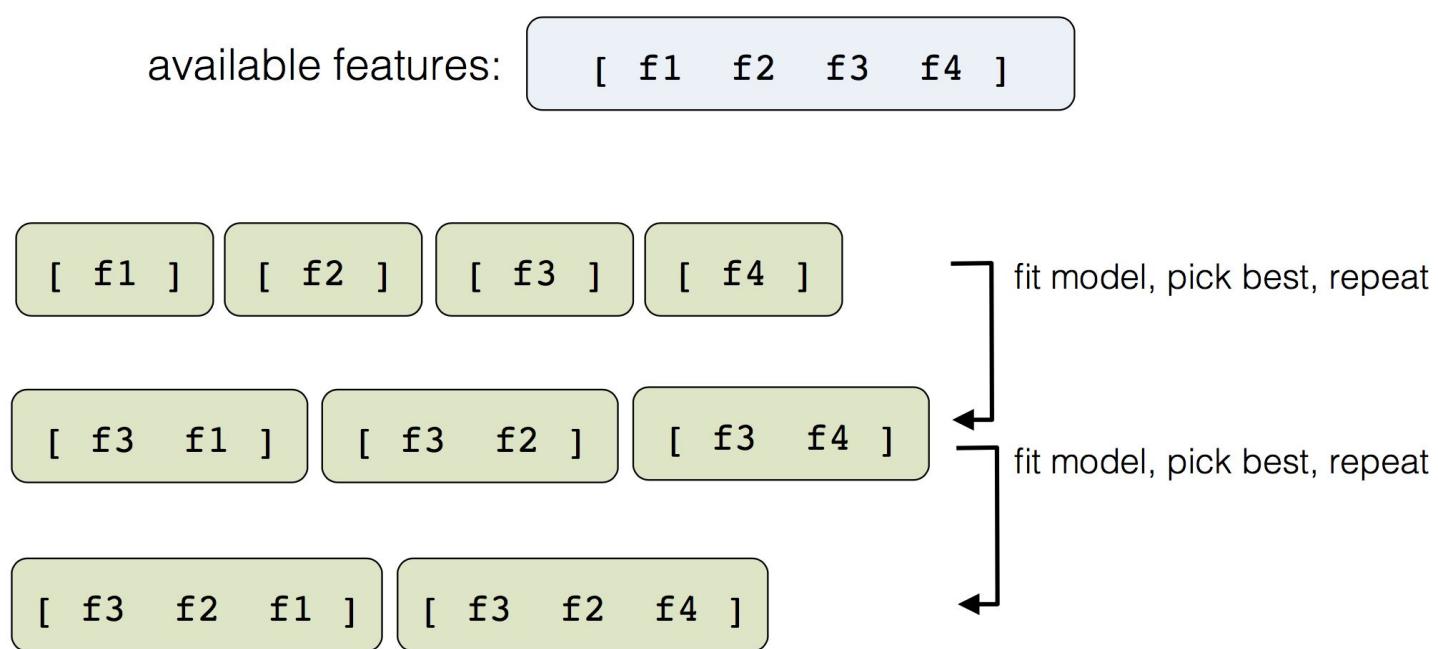
fit model, remove lowest weight, repeat

[w1 w4]

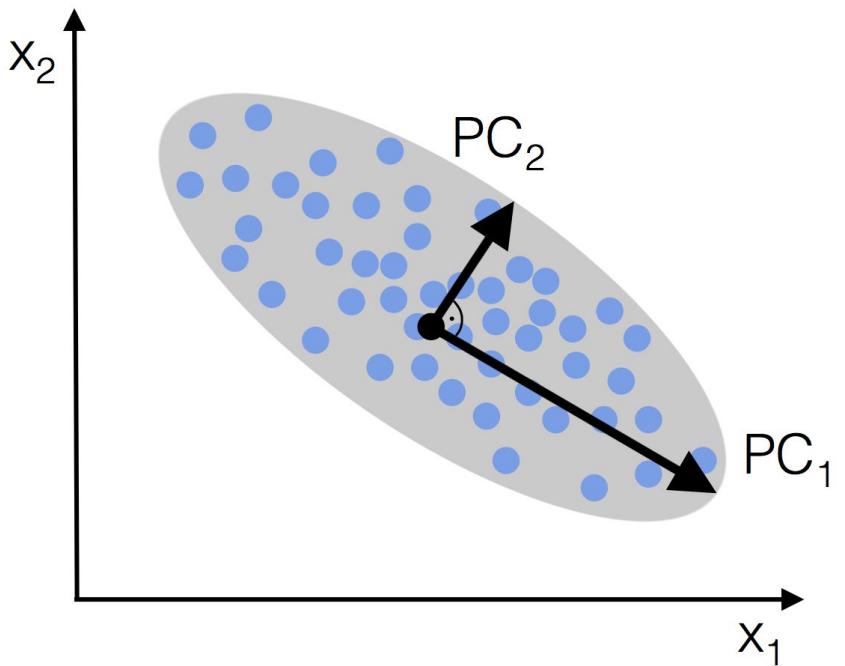
fit model, remove lowest weight, repeat

[w4]

Sequential Feature Selection



Principal Component Analysis



MLDS

Coding Time!

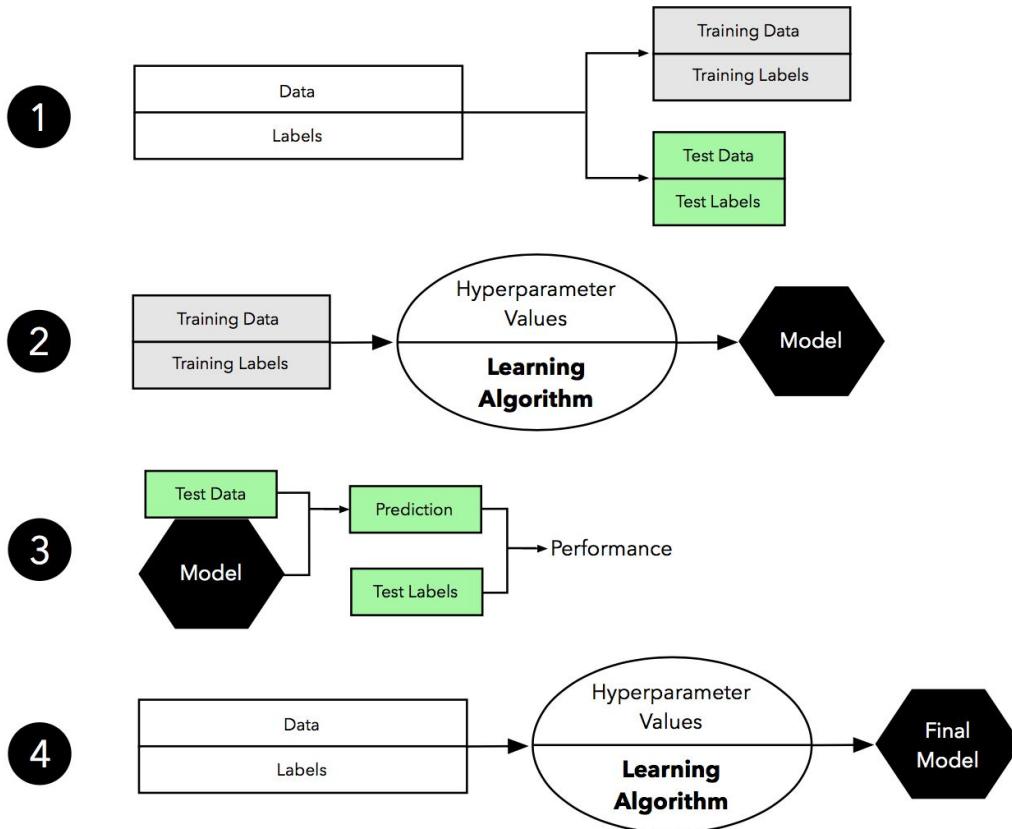


MLDS

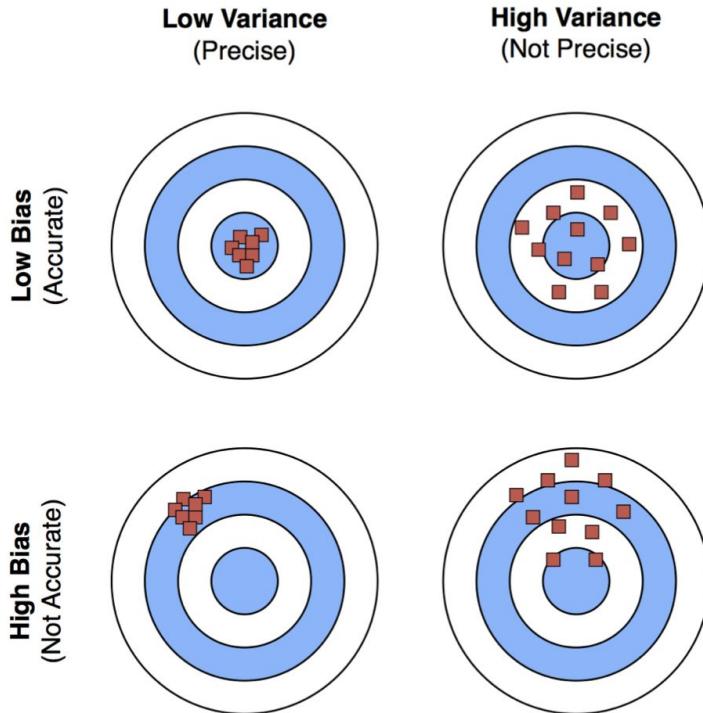
Topics

- 
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. Model Evaluation & Hyperparameter Tuning**

“Basic” Supervised Learning Workflow



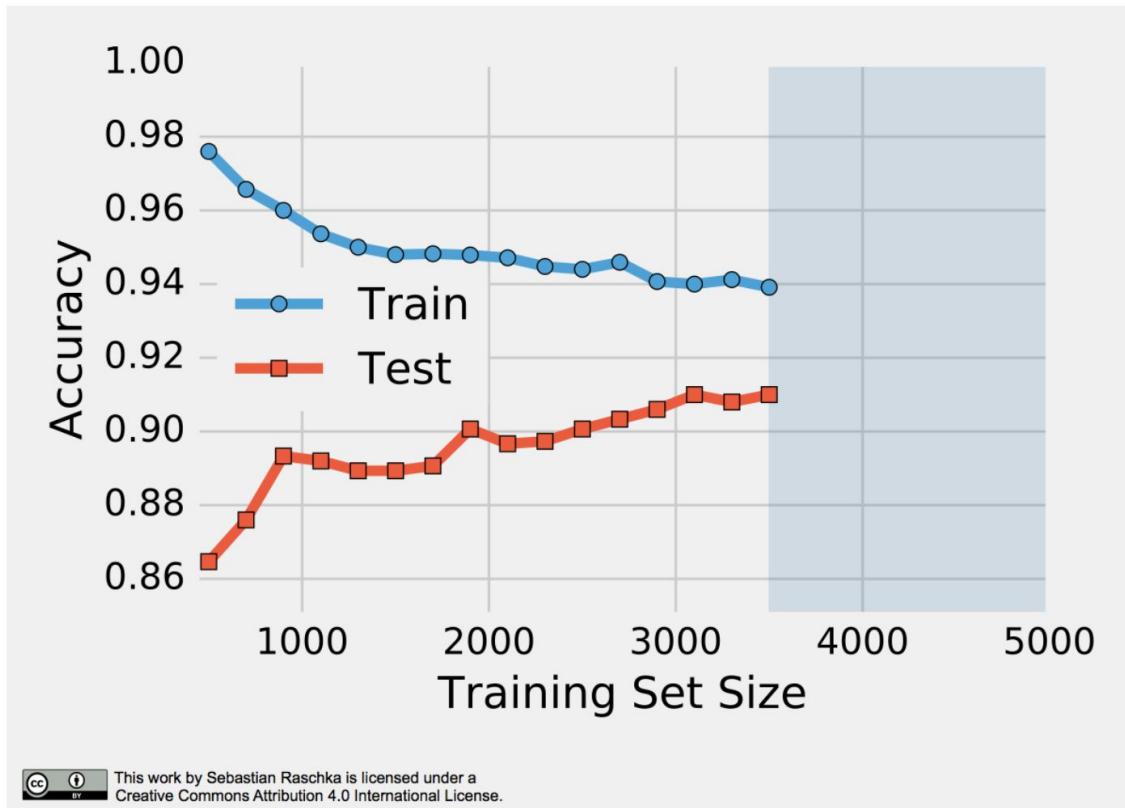
Bias and Variance



This work by Sebastian Raschka is licensed under a
Creative Commons Attribution 4.0 International License.

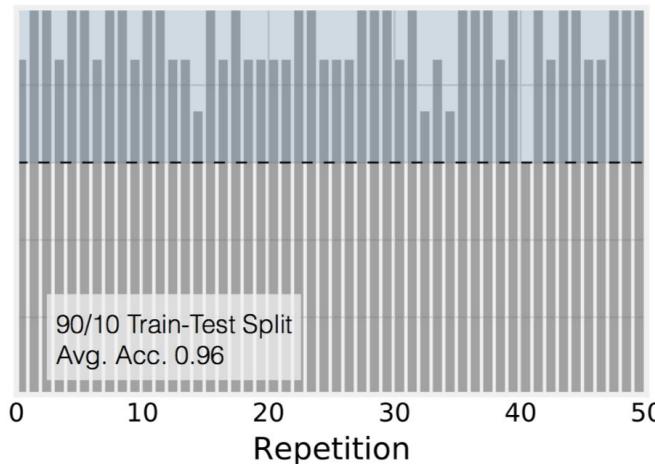
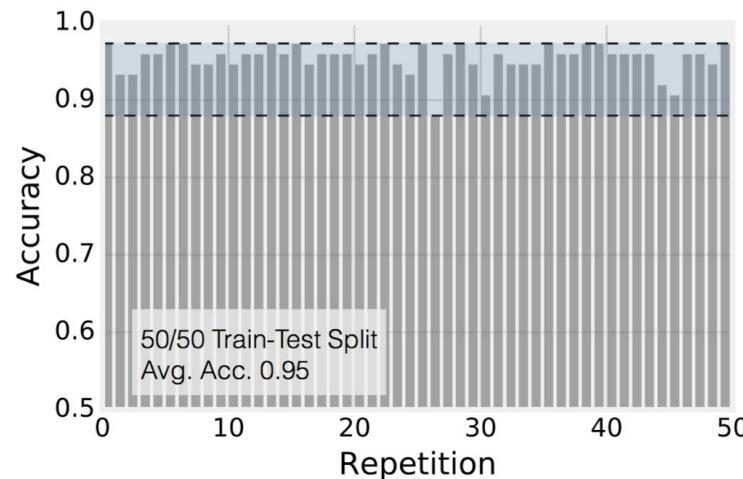
MLDS

Learning Curves



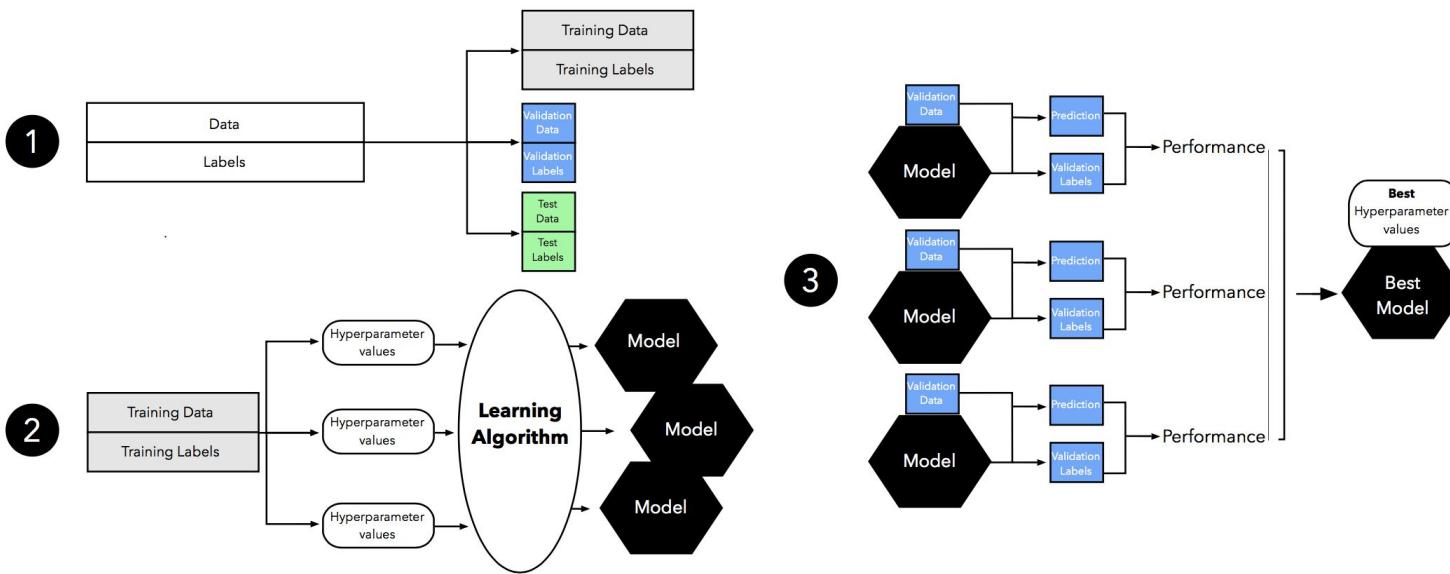
MLDS

Repeated Holdout



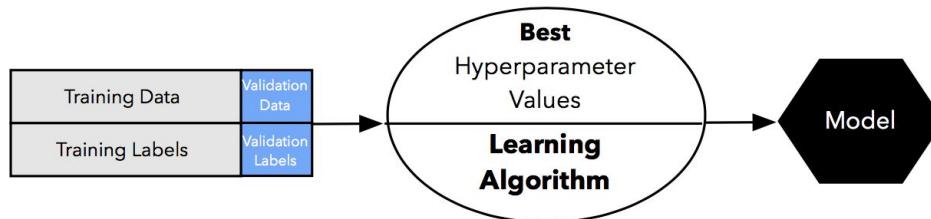
This work by Sebastian Raschka is licensed under a
Creative Commons Attribution 4.0 International License.

Holdout and Hyperparameter Tuning I

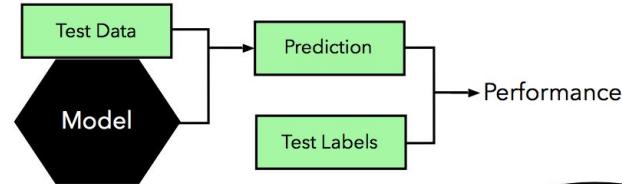


Holdout and Hyperparameter Tuning II

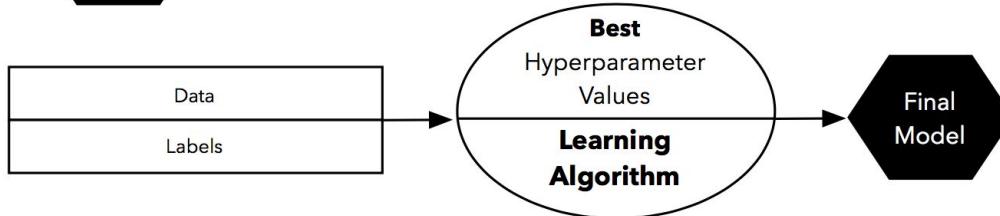
4



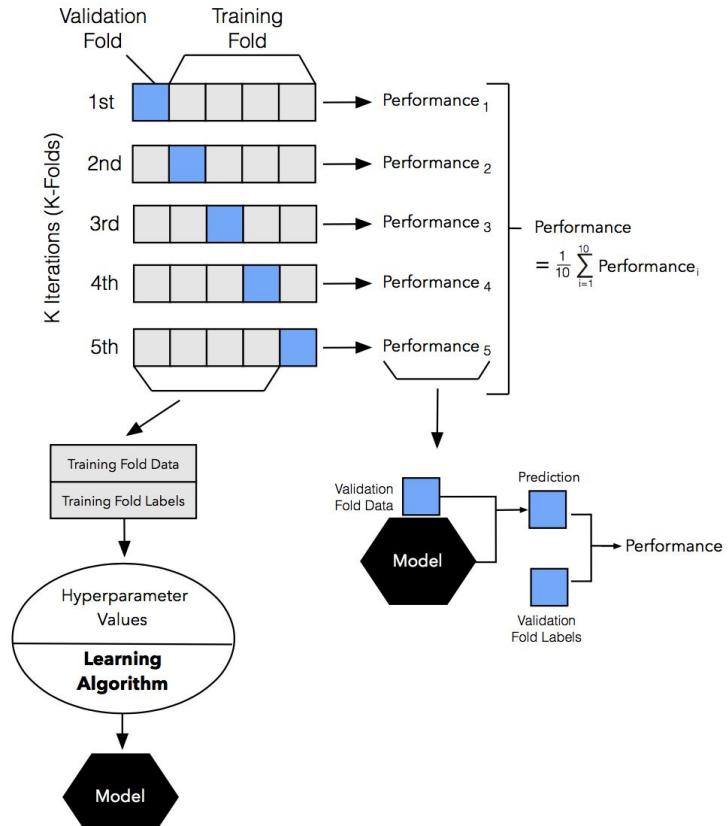
5



6

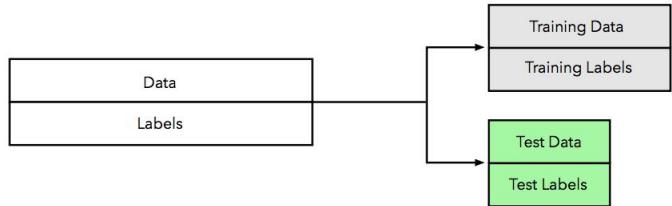


K-fold Cross-Validation

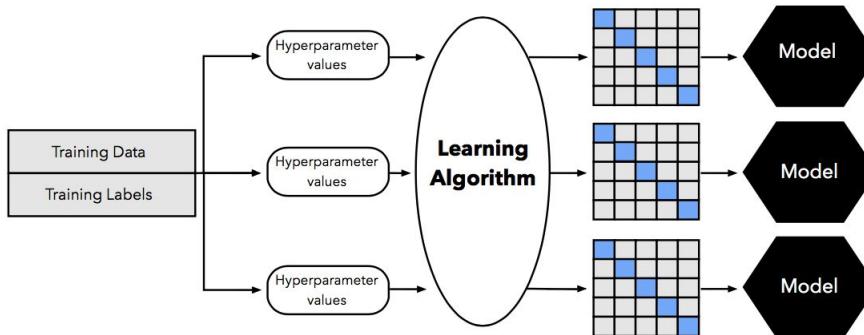


K-fold Cross-Validation Workflow I

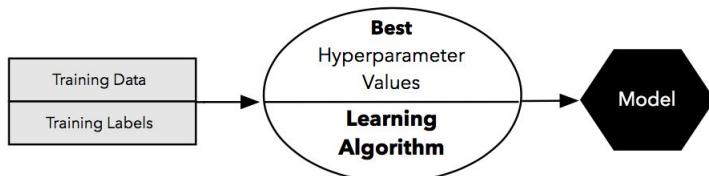
1



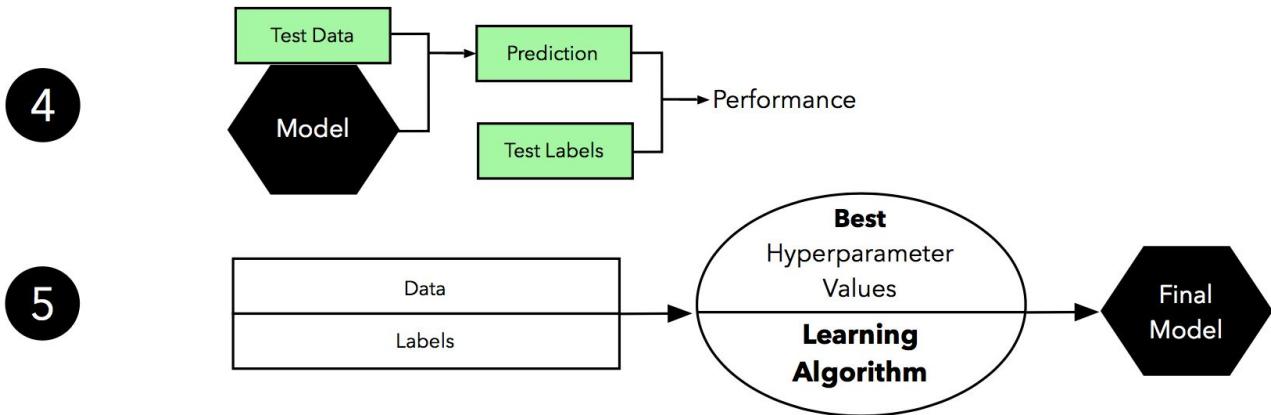
2



3



K-fold Cross-Validation Workflow II



MLDS

Coding Time!

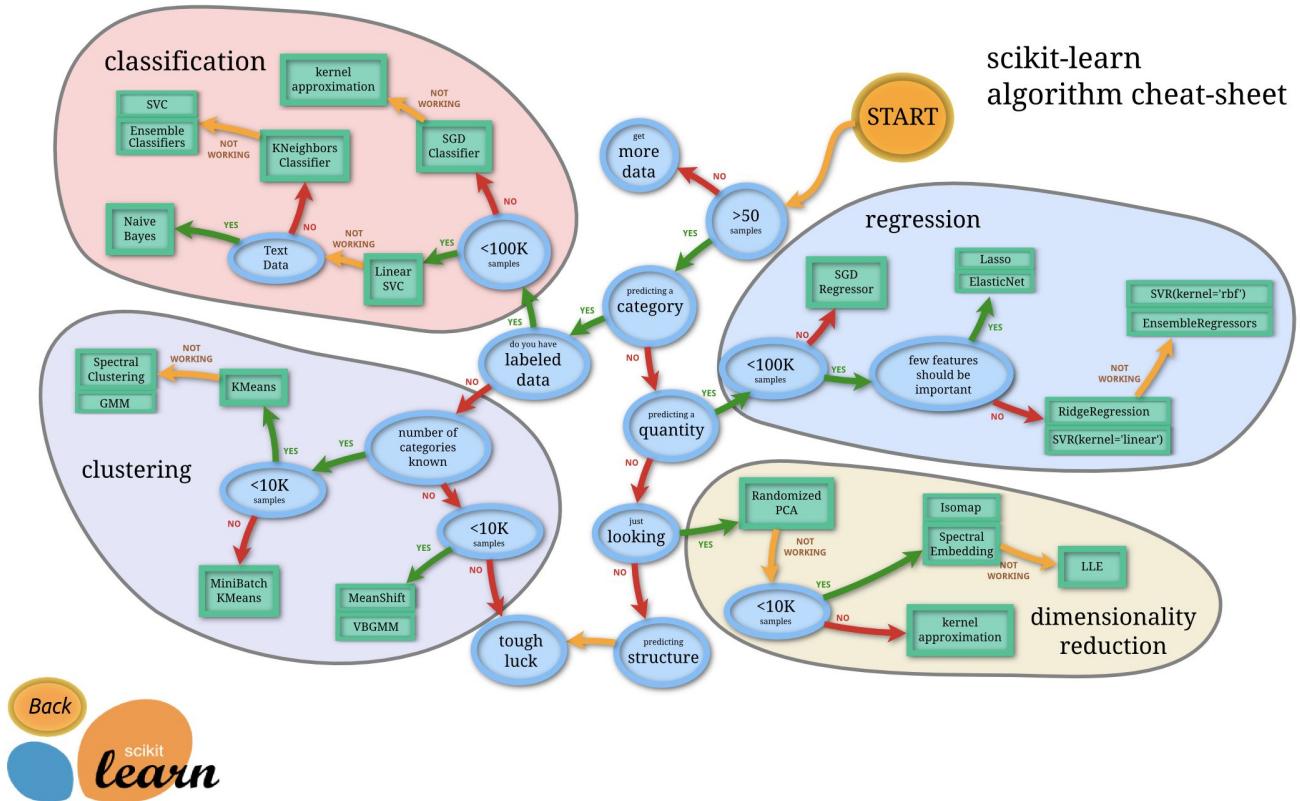


MLDS

Performance Metrics

http://scikit-learn.org/stable/modules/model_evaluation.html

Further Resources - Andreas' “cheat sheet”



Thanks!

Machine Learning for Data Science Interest Group
Advanced Informatics School
Universiti Teknologi Malaysia

@utmmlds
ais.utm.my/mlds

May 2017



MLDS