

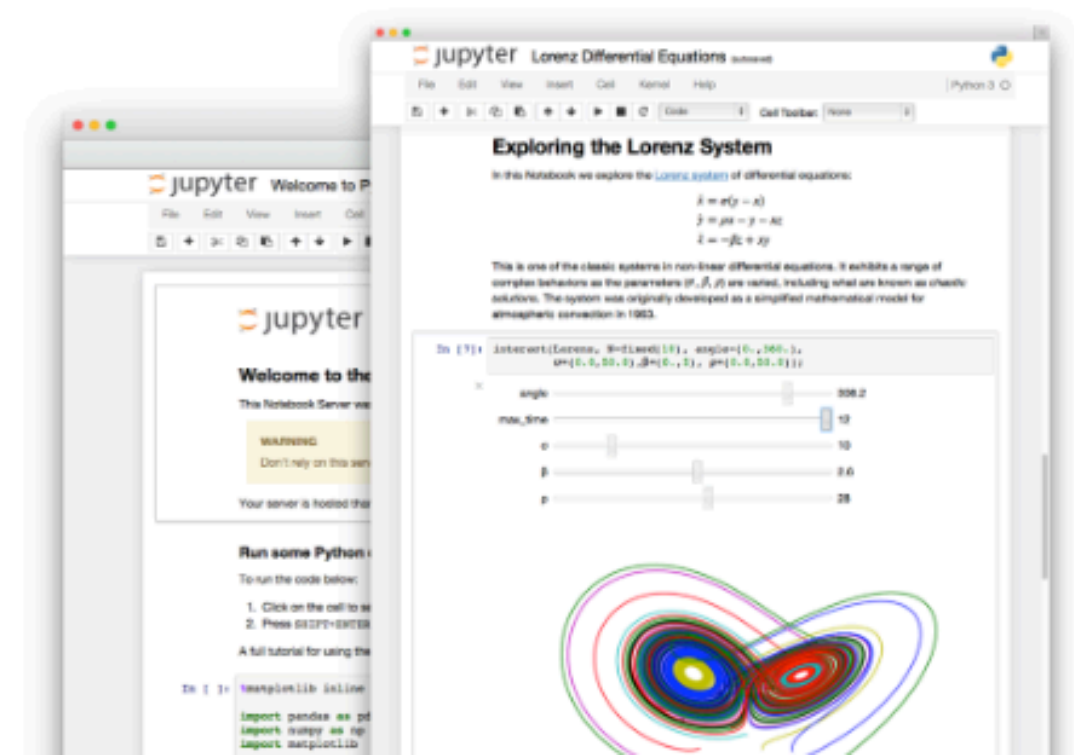
# Linear Regression with Python

# Tools



## Jupyter Notebook

The Jupyter Notebook is a web-based interactive computing platform that allows users to author data- and code-driven narratives that combine live code, equations, narrative text, visualizations, interactive dashboards and other media.



Projects/jupyter-notebook/

Sebastian

localhost:8888/tree/Projects/jupyter-notebook#notebooks

jupyter

Logout

Files

Running

Clusters

Select items to perform actions on them.

Upload

New

☐

/ Projects / jupyter-notebook

Name

Last Modified

..

seconds ago

☐

my-first-notebook.ipynb

6 days ago

Projects/jupyter-notebook/Untitled

localhost:8888/notebooks/Projects/jupyter-notebook/Untitled.ipynb?kernel\_name=python3

jupyter

Untitled

Last Checkpoint: a few seconds ago (unsaved changes)

Logout

FileEditViewInsertCellKernelWidgetsHelp

TrustedPython 3

In [ ]:

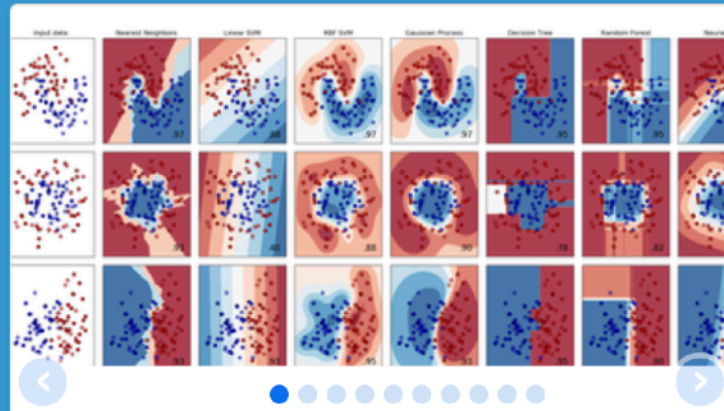
```
In [1]: print('Hello World')
```

```
Hello World
```

```
In [2]: i = 1  
while i <= 10:  
    print(i)  
    i = i + 1
```

```
1  
2  
3  
4  
5  
6  
7  
8  
9  
10
```

```
In [ ]:
```



# scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

## Classification

Identifying to which category an object belongs to.

**Applications:** Spam detection, Image recognition.

**Algorithms:** SVM, nearest neighbors, random forest, ... — [Examples](#)

## Regression

Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.

**Algorithms:** SVR, ridge regression, Lasso, ... — [Examples](#)

## Clustering

Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes

**Algorithms:** k-Means, spectral clustering, mean-shift, ... — [Examples](#)

## Dimensionality reduction

Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency

**Algorithms:** PCA, feature selection, non-negative matrix factorization. — [Examples](#)

## Model selection

Comparing, validating and choosing parameters and models.

**Goal:** Improved accuracy via parameter tuning

**Modules:** grid search, cross validation, metrics. — [Examples](#)

## Preprocessing

Feature extraction and normalization.

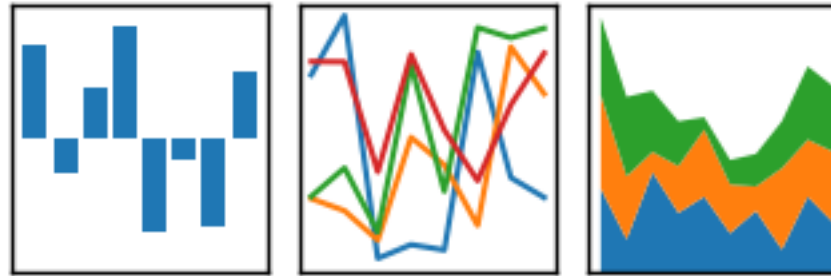
**Application:** Transforming input data such as text for use with machine learning algorithms.

**Modules:** preprocessing, feature extraction. — [Examples](#)

<http://scikit-learn.org/stable/>

# pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



*pandas*

- open source library
- high-performance, easy-to-use data structures and data analysis tools

<https://pandas.pydata.org/>





[Scipy.org](#)

# NumPy

NumPy is the fundamental package for scientific computing with Python. It contains among other things:

- a powerful N-dimensional array object
- sophisticated (broadcasting) functions
- tools for integrating C/C++ and Fortran code
- useful linear algebra, Fourier transform, and random number capabilities

Besides its obvious scientific uses, NumPy can also be used as an efficient multi-dimensional container of generic data. Arbitrary data-types can be defined. This allows NumPy to seamlessly and speedily integrate with a wide variety of databases.

NumPy is licensed under the [BSD license](#), enabling reuse with few restrictions.

## Getting Started

- [Getting NumPy](#)
- [Installing the SciPy Stack](#)
- [NumPy and SciPy documentation page](#)
- [NumPy Tutorial](#)
- [NumPy for MATLAB® Users](#)
- [NumPy functions by category](#)
- [NumPy Mailing List](#)

For more information on the SciPy Stack (for which NumPy provides the fundamental array data structure), see [scipy.org](#).

[Donate to Numpy](#)

[Open Hub](#)

[NumPy](#)

[About NumPy](#)

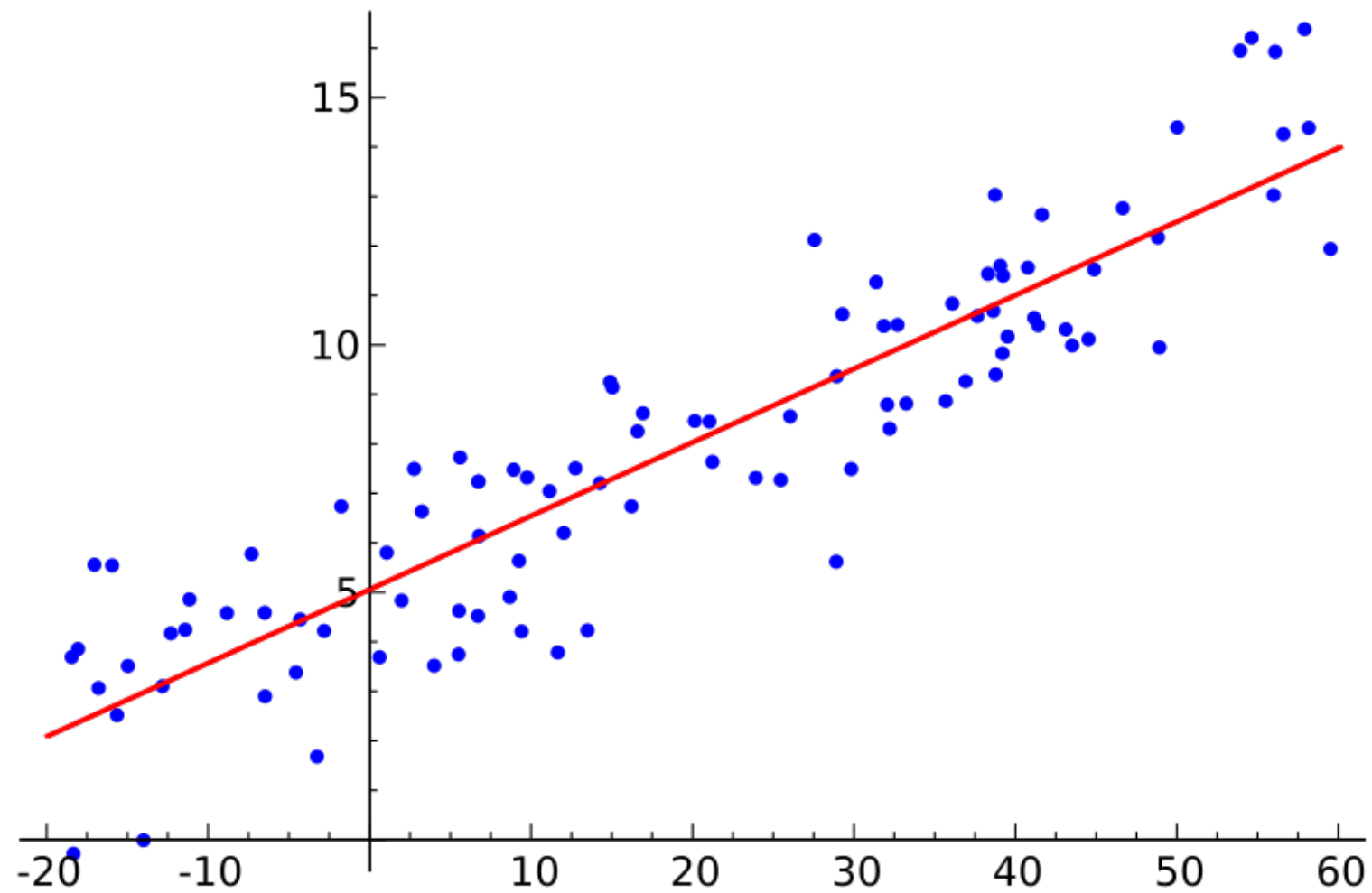
[License](#)

[Old array packages](#)

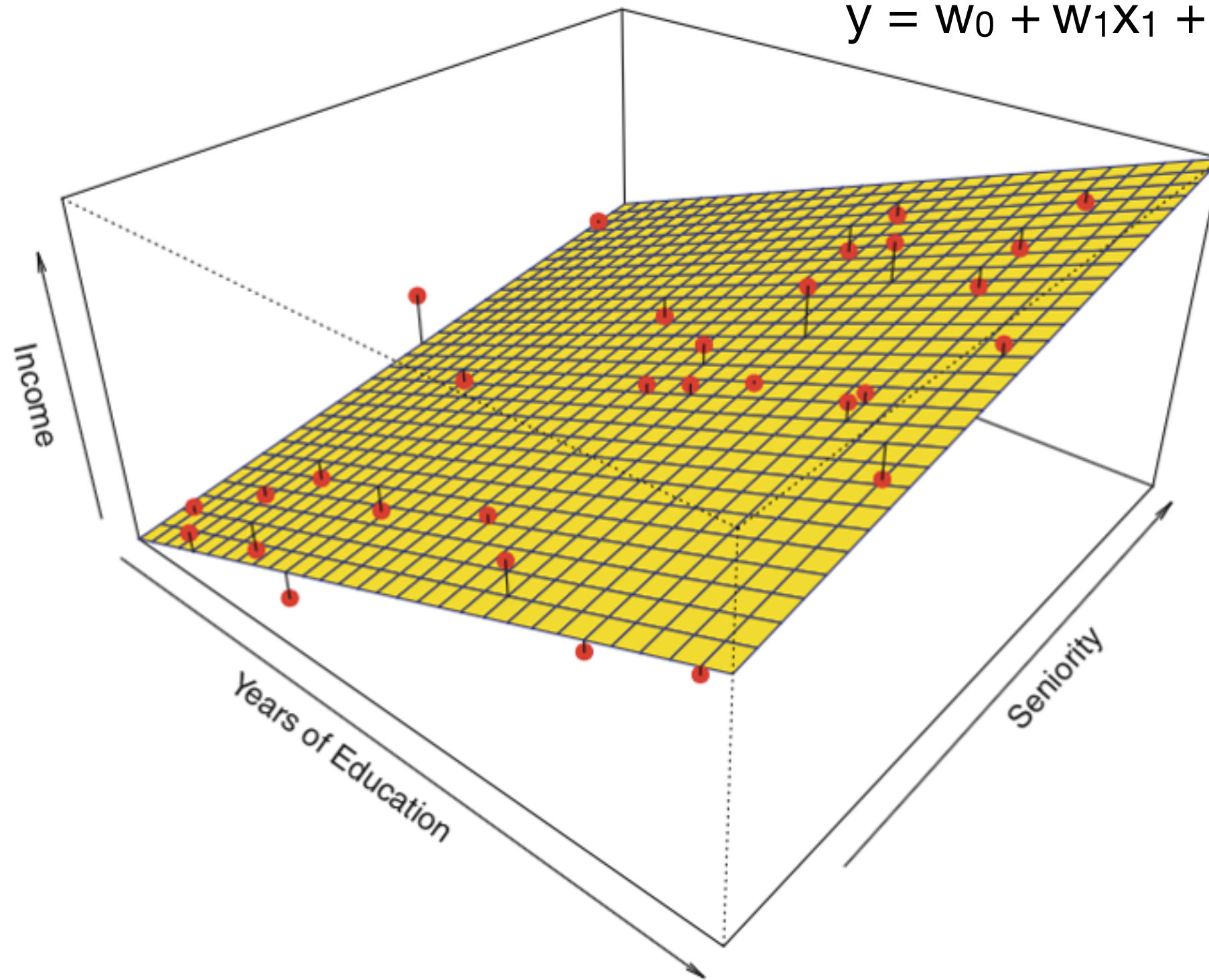
<http://www.numpy.org/>

# Linear Regression

$$y = w_1 x + w_0$$



$$y = w_0 + w_1X_1 + w_2X_2$$



**General formulation:**

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

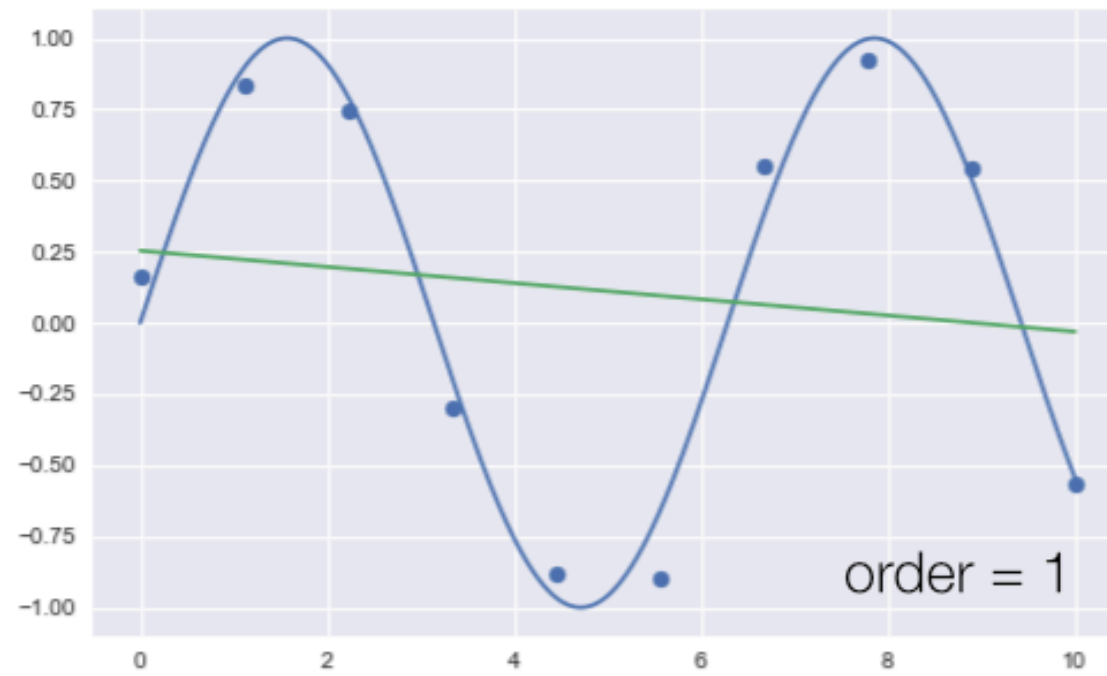
**General formulation:**

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j x_j$$

**General formulation:**

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$



?

The key property of this model is that it is a linear function of the parameters !

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j x_j \quad \text{Linear in } \mathbf{x} \text{ and } \mathbf{w}$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) \quad \text{Linear in } \mathbf{w} \text{ and } \phi(\mathbf{x})$$

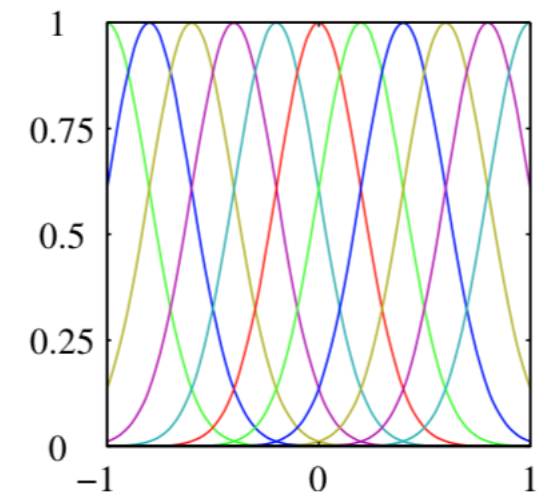
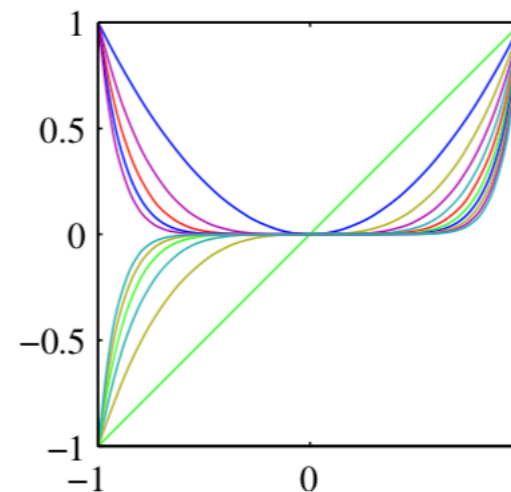


The key property of this model is that it is a linear function of the parameters !

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j x_j \quad \text{Linear in } \mathbf{x} \text{ and } \mathbf{w}$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) \quad \text{Linear in } \mathbf{w} \text{ and } \phi(\mathbf{x})$$

**Basis functions**

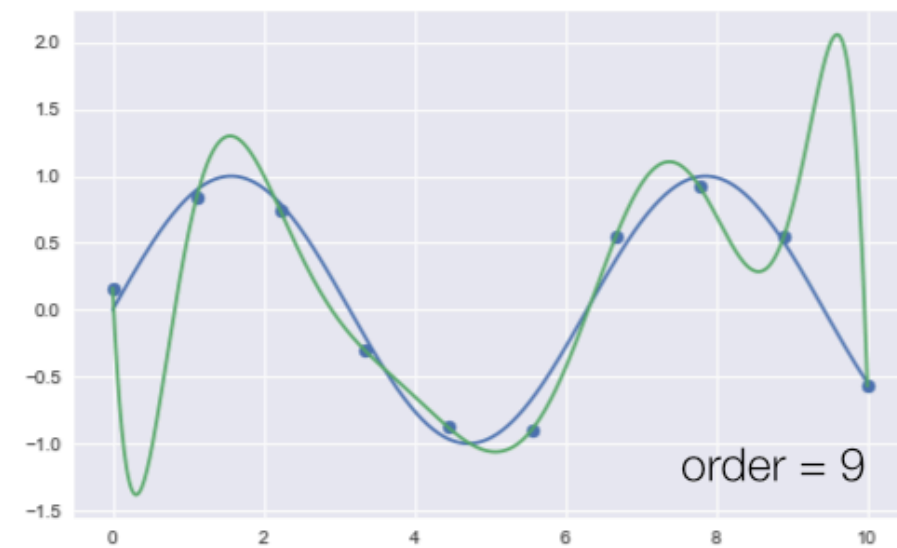
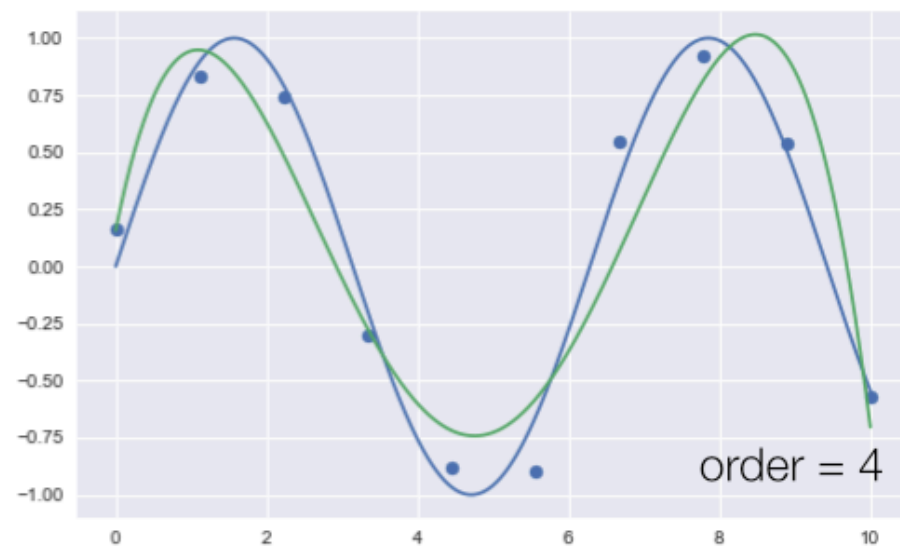
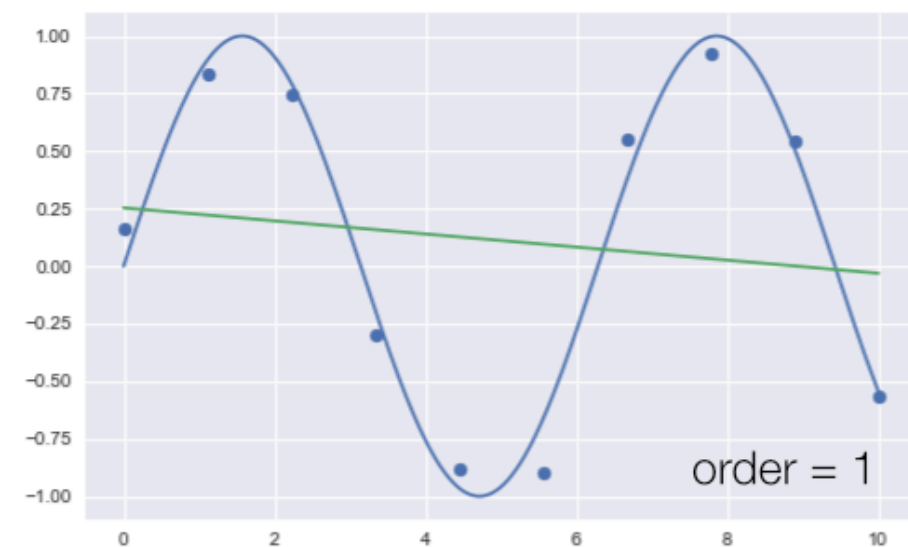
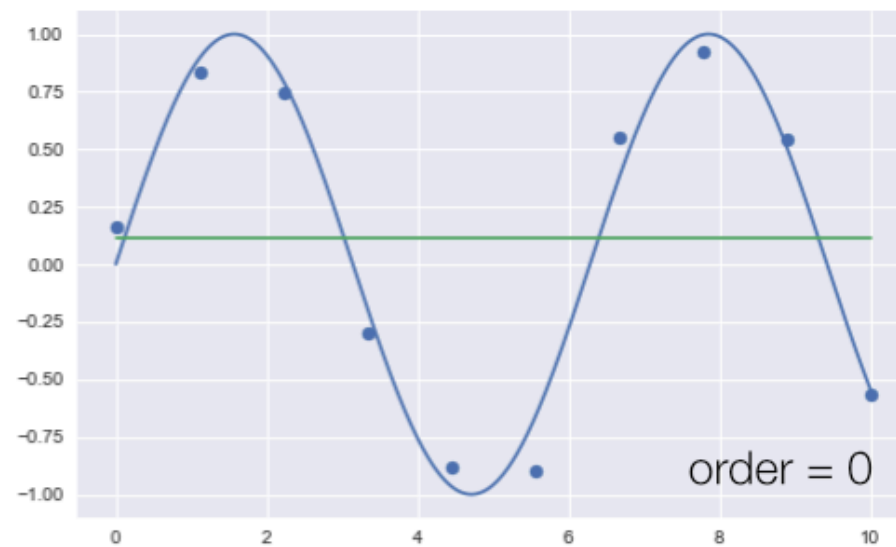


$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

**Linear in  $\mathbf{w}$  and  $\phi(\mathbf{x})$**

## Basis functions

fitting polynomial models of increasing complexity



$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j x_j$$

- **goal:** Predicting a continuous  $y$  as a function of the variables  $\mathbf{x}$  and the parameters  $\mathbf{w}$
- **Model assumption:**
  - data generated by  $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$
  - Error  $\epsilon \sim$  normal distribution with zero mean
- **Loss function:** sum-of-squares error between  $t$  and prediction  $y$  (least squares)