

COSC 325: Introduction to Machine Learning

Dr. Hector Santos-Villalobos

Dr. Santos



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

Lecture 03: Scientific Computing with Python



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE



Class Announcements

Homework:

First homework due Sunday 09/01.
The due date may shift according to
the material covered.
TAs will not troubleshoot your code.

Quizzes:

Office hours question out of date.

Course Project:

Pick teammates by 08/29

Today's Topics

Notation



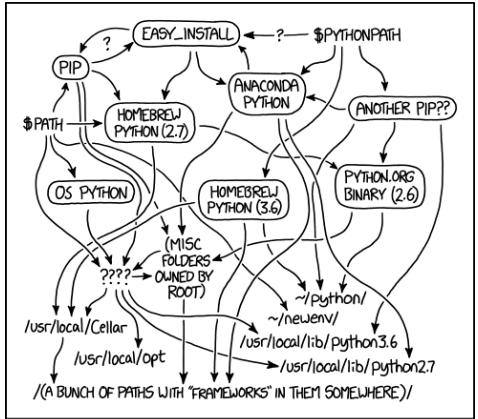
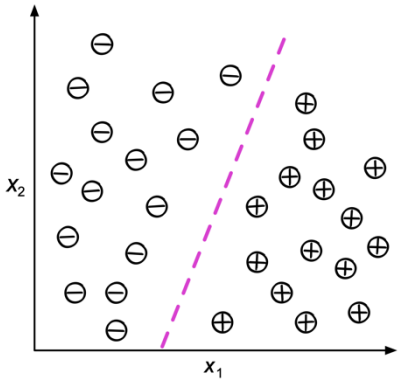
Scientific Computing in Python



Last Lecture

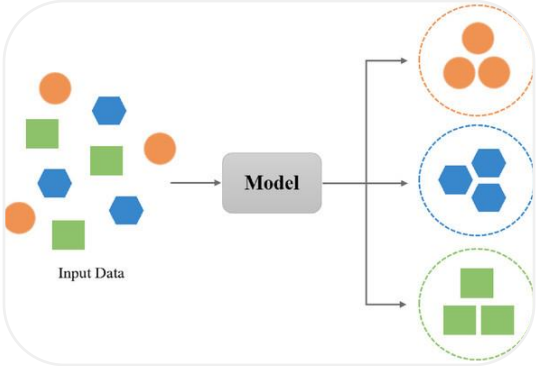


- Machine learning
 - A subfield of artificial intelligence
 - Models need to generalize (i.e., learn)
 - Task, Experience, Performance
 - Different learning categories
- Programming
 - Python: flexible, efficient, collaborative, powerful
 - Always work from a dockerized or virtual environment
 - Practice



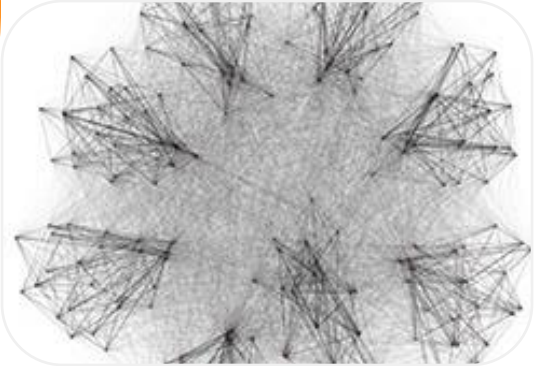
MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.

Machine Learning Categories



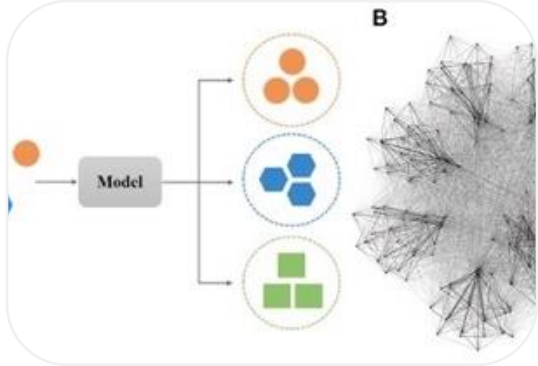
Supervised Learning

- Trained on “Labeled dataset”
- Needs pairs of inputs and outputs (ground truth)



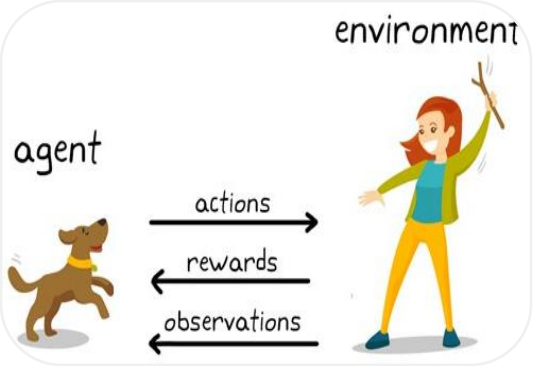
Unsupervised Learning

- The algorithm discovers patterns and relationships in unlabeled data
- It needs inputs only and, most of the time, some context. (e.g., number of unique labels)



Semi-Supervised Learning

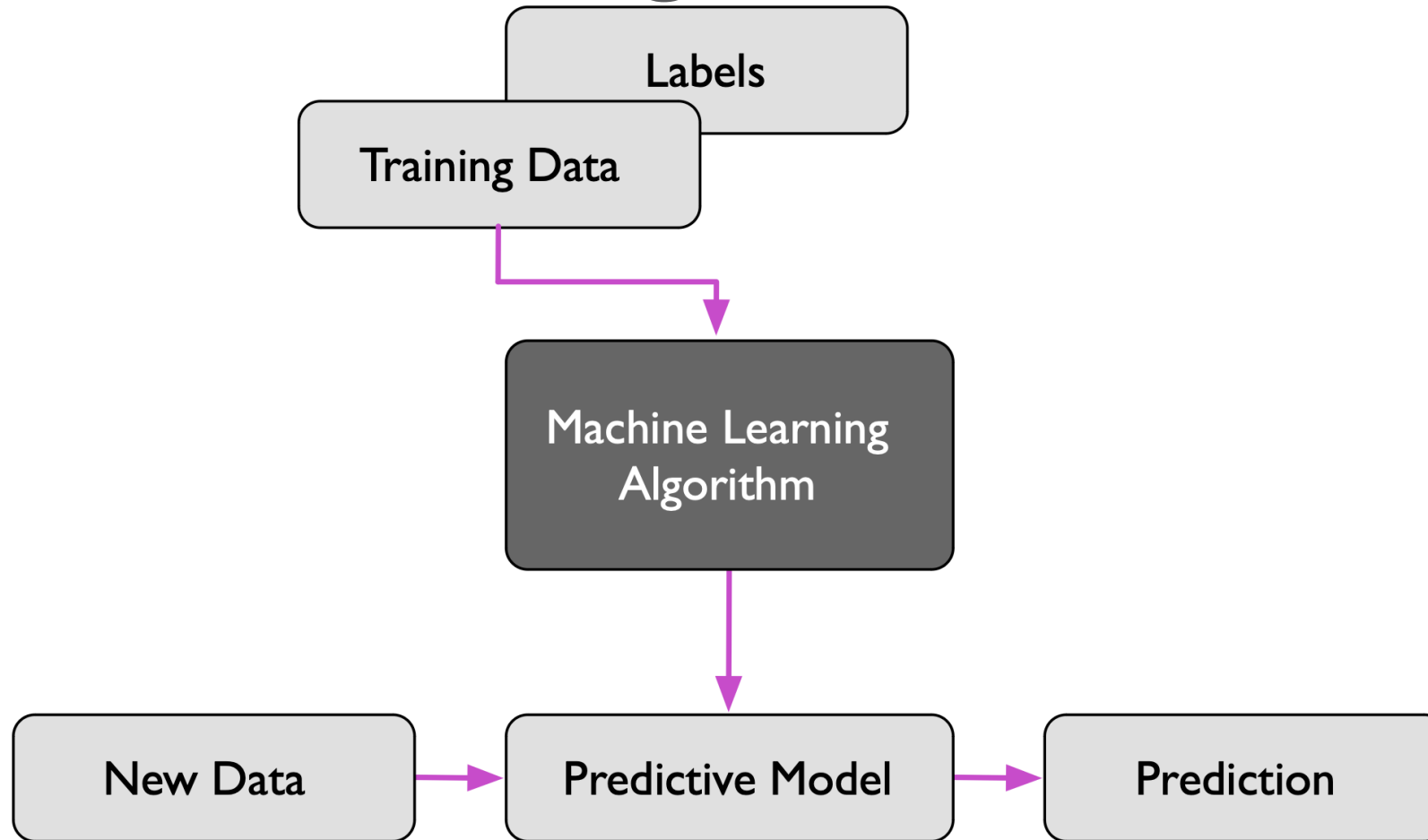
- Combines SL and UL
- E.g., a small subset of labeled data is used to label unlabeled data.
- E.g., Generative Adversarial Network mapping blonde to brunette.



Reinforcement Learning

- It learns by interacting with the environment.
- Trial, error, and delay
- Needs well-defined reward feedback.

Supervised Learning Workflow



Notation

Features

Targets/Labels

n training examples

Training set: $D = \{\langle x^{(i)}, y^{(i)} \rangle, i = 1, \dots, n\}$

Training sample #3: $\langle x^{(3)}, y^{(3)} \rangle$

This is the function/program our ML algorithm will generate.

Unknown function: $f(x) = y$

Hypothesis: $h(x) = \hat{y}$

Classification

$h: \mathbb{R}^m \rightarrow \text{_____}$

Features

Binary

Multi-class

$\{0, 1\}$

$\{0, 1, 3, 4, 5\}$

$\{ 'a', 'q', 'z' \}$

Notation

Features

Targets/Labels

n training examples

Training set: $D = \{\langle x^{(i)}, y^{(i)} \rangle, i = 1, \dots, n\}$

Training sample #3: $\langle x^{(3)}, y^{(3)} \rangle$

This is the function/program our ML algorithm will generate.

Unknown function: $f(x) = y$

Hypothesis: $h(x) = \hat{y}$

Classification

Regression

$h: \mathbb{R}^m \rightarrow \underline{\{0, 1, 3, 4, 5\}}$

$h: \mathbb{R}^m \rightarrow \underline{\mathbb{R}}$

Features

Inputs Representation

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

$m = n_x$

House Market Price



$$x^{(i)} = \begin{bmatrix} \text{Sqft size of house} \\ \text{zipcode} \\ \vdots \\ \text{no. bathrooms} \\ \text{no. bedrooms} \end{bmatrix}$$

List of floating or integer numbers

Inputs Representation

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

Dog detector
Input $x^{(i)}$

Pixel values for 8-bit grayscale image are between 0 and 255.

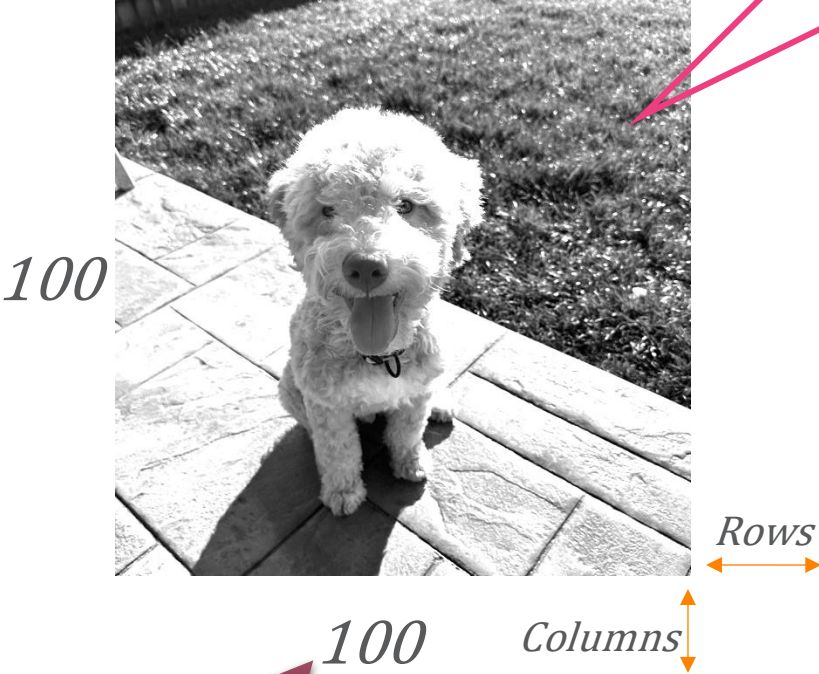


Image Dimensions

Input Notation

Input $x^{(i)}$



100

100

Columns

Rows

Image Dimensions

2D matrix $x^{(i)}$

0	5	5	0
25	4	30	120
135	130	200	230
145	160	210	250

⋮

Unroll

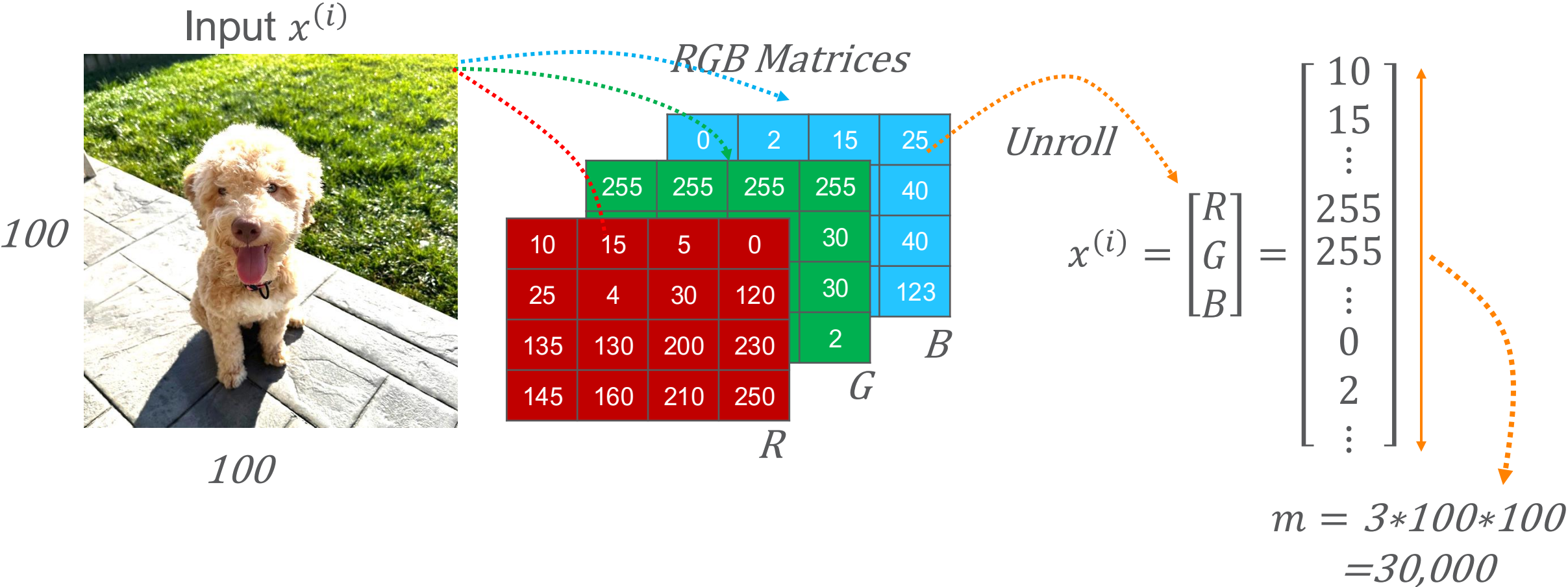
$x^{(i)} =$

0
5
5
0
⋮
25
4
30
⋮

$$m = n_x = 100 * 100 = 10,000$$

Number of features

What about color?



What about text?

Email Spam Detector



National Security Department

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

A vulnerability has been identified in the Apple Facetime mobile applications that allow an attacker to record calls and videos from your mobile device without your knowledge.

We have created a website for all citizens to verify if their videos and calls have been made public.

To perform the verification, please use the following link:

[Facetime Verification](#)

This website will be available for 72 hours.

National Security Department

Check this tutorial about a spam email classifier with logistic regression: <https://towardsdatascience.com/spam-detection-with-logistic-regression-23e3709e522>

What about text?

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector



Email Spam Detection
National Security Department

A vulnerability has been identified in the Apple Facetime mobile applications that allow an attacker to record calls and videos from your mobile device without your knowledge.

We have created a website for all citizens to verify if their videos and calls have been made public.

To perform the verification, please use the following link:

[Facetime Verification](#)

This website will be available for 72 hours.

National Security Department

Tokens or terms

Tokenization

Bag of Words

	identified	vulnerability	Urgent	Money	Transaction	Pay	link	Password	Verify	Hack
Doc 1	1	1	0	0	0	0	1	0	3	1
Doc 2	0	2	1	0	1	0	1	1	1	1
Doc 3	0	0	1	2	1	2	1	0	0	0

Frequency

Check this tutorial about a spam email classifier with logistic regression: <https://towardsdatascience.com/spam-detection-with-logistic-regression-23e3709e522>

What about text?

Bag of Words

	identified	vulnerability	Urgent	Money	Transaction	Pay	link	Password	Verify	Hack
Doc 1	1	1	0	0	0	0	1	0	3	1
Doc 2	0	2	1	0	1	0	1	1	1	1
Doc 3	0	0	1	2	1	2	1	0	0	0

Tokenization

Email Spam Detection



National Security Department

A vulnerability has been identified in the Apple Facetime mobile applications that allow an attacker to record calls and videos from your mobile device without your knowledge.

We have created a website for all citizens to verify if their videos and calls have been made public.

To perform the verification, please use the following link:

[Facetime Verification](#)

This website will be available for 72 hours.

National Security Department

Term Weighting

$$w_j^{[i]} = TF_j^{[i]} * \ln(n/n_j)$$

$$w_9^{[1]} = 3 * \ln(3/2)$$

documents

documents with term

	identified	vulnerability	Urgent	Money	Transaction	Pay	link	Password	Verify	Hack
Doc 1									0.53	
Doc 2										
Doc 3										

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

What about text?

Bag of Words

	identified	vulnerability	Urgent	Money	Transaction	Pay	link	Password	Verify	Hack
Doc 1	1	1	0	0	0	0	1	0	3	1
Doc 2	0	2	1	0	1	0	1	1	1	1
Doc 3	0	0	1	2	1	2	1	0	0	0

Tokenization

Email Spam Detection



National Security Department

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

A vulnerability has been identified in the Apple Facetime mobile applications that allow an attacker to record calls and videos from your mobile device without your knowledge.

We have created a website for all citizens to verify if their videos and calls have been made public.

To perform the verification, please use the following link:

[Facetime Verification](#)

This website will be available for 72 hours.

National Security Department

Term Weighting

$$w_j^{[i]} = TF_j^{[i]} * \ln(n/n_j)$$

	identified	vulnerability	Urgent	Money	Transaction	Pay	link	Password	Verify	Hack
Doc 1	0.48	0.18	0.00	0.00	0.00	0.00	0.00	0.00	0.53	0.18
Doc 2	0.00	0.35	0.18	0.00	0.18	0.00	0.00	0.48	0.18	0.18
Doc 3	0.00	0.00	0.18	0.95	0.18	0.95	0.00	0.00	0.00	0.00

The document features

Representing the dataset

A sample: (x, y) $x \in \mathbb{R}^m$, $y \in \{0,1\}$

Input

Ground Truth Label

Set of Samples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(n)}, y^{(n)})\}$

Matrix representation of inputs: $X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(n)T} \end{bmatrix}$

$X \in \mathbb{R}^{n \times m}$

Python $X.shape = (n, m)$

Number of rows

Number of columns

Matrix representation of labels: $Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

$Y \in \mathbb{R}^{n \times 1}$

$Y.shape = (n, 1)$

Data Representation

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

Feature Vector

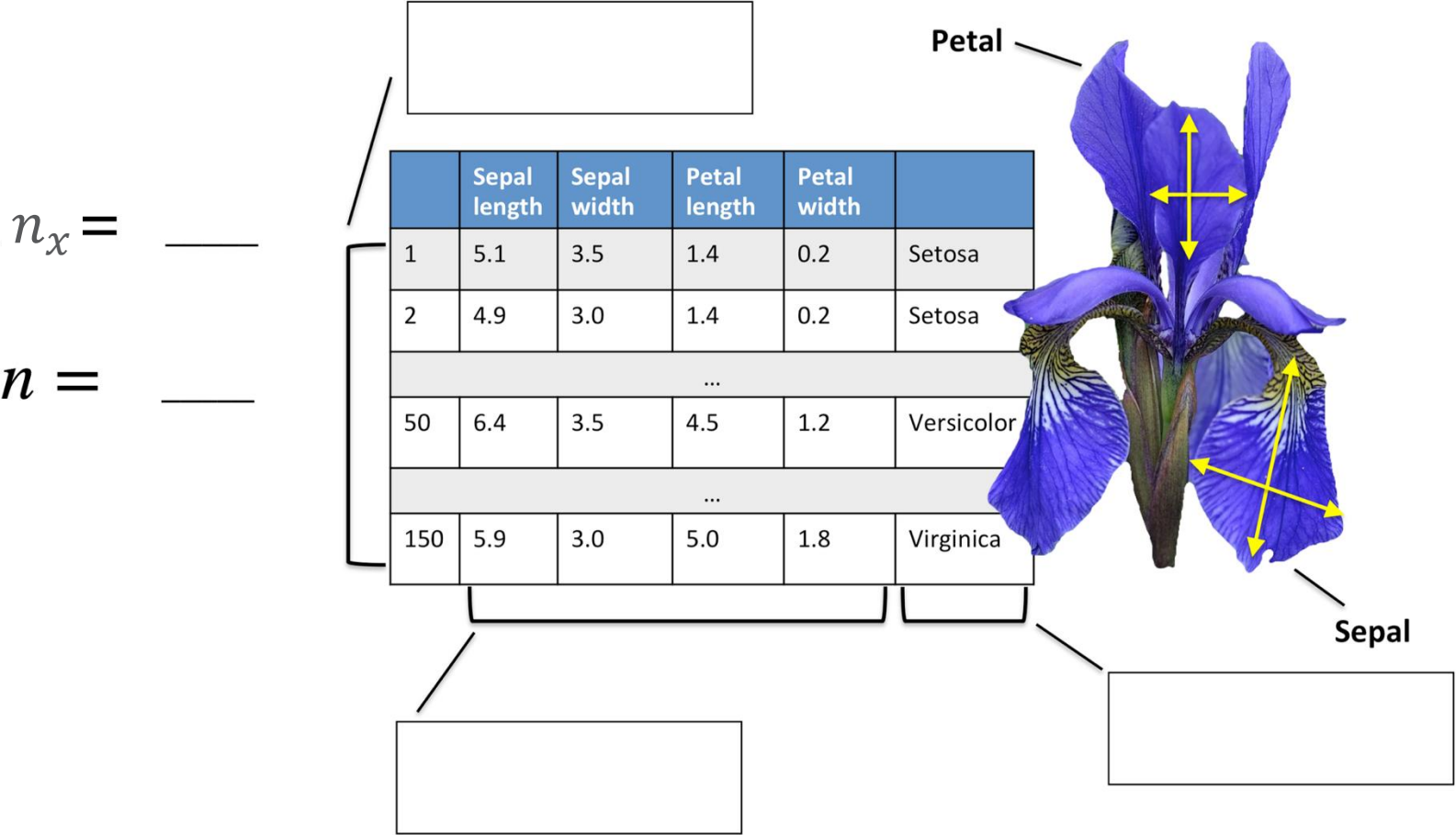
$$X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ \vdots \\ x^{(n)T} \end{bmatrix}$$

Input matrix of size $\langle n, m \rangle$

What is the value of sample 13 feature #5?

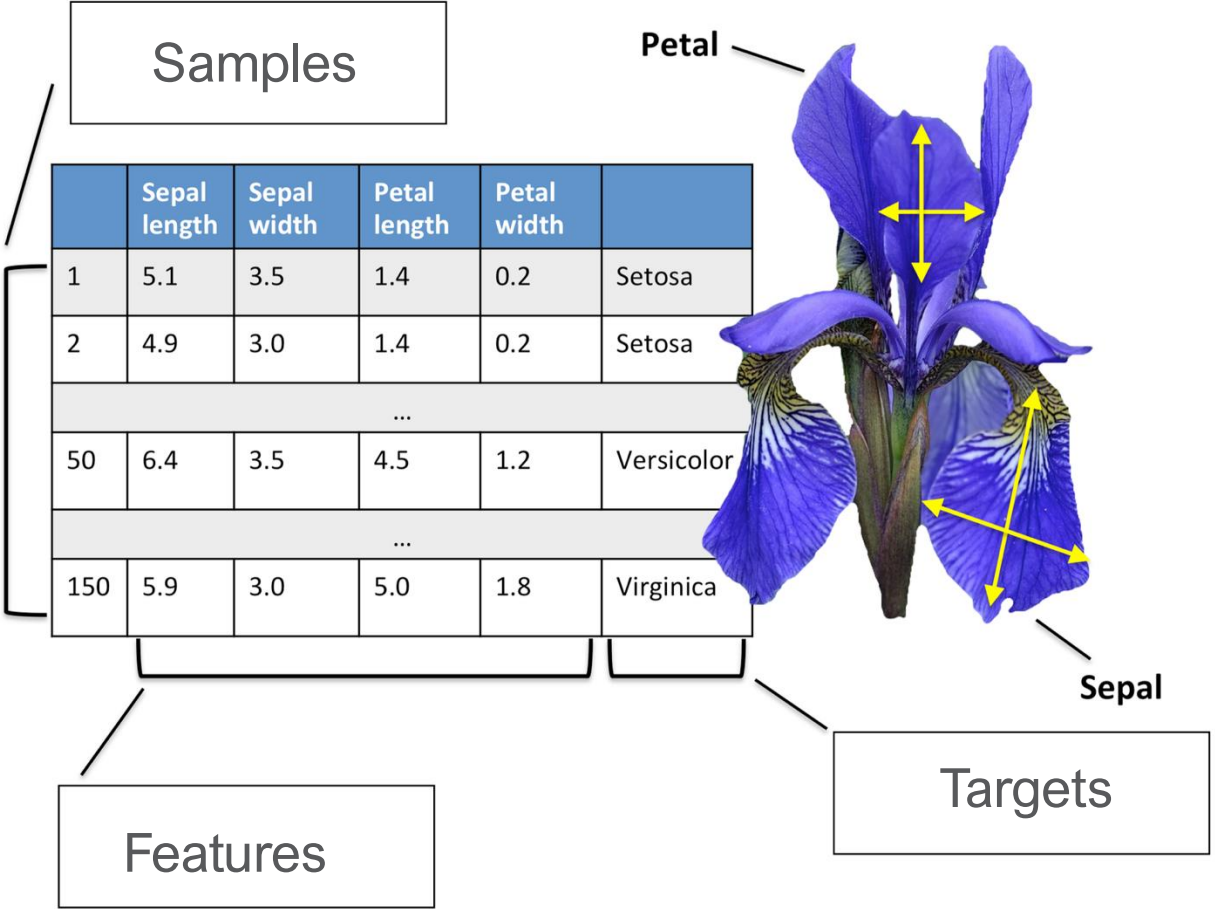
$$x_5^{(13)}$$

Data Representation



Data Representation

$$n_x = \frac{4}{n = 150}$$



Pop Quiz #1

1 | MULTIPLE CHOICE

POINTS: 1 |  Edit 

Below is an input matrix X with six samples and features. Following the notation discussed in class, what is sample's three feature vector?

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

- A. 21
- B. [3, 9, 15, 21, 27, 33]
- C. [13, 14, 15, 16, 17, 18]
- D. [4, 10, 16, 22, 28, 34]

Pop Quiz #1

1 | MULTIPLE CHOICE

POINTS: 1 |  Edit 

Below is an input matrix X with six samples and features. Following the notation discussed in class, what is sample's three feature vector?

$x^{(1)}$	1	2	3	4	5	6
$x^{(2)}$	7	8	9	10	11	12
$x^{(3)}$	13	14	15	16	17	18
$x^{(4)}$	19	20	21	22	23	24
$x^{(5)}$	25	26	27	28	29	30
$x^{(6)}$	31	32	33	34	35	36

A. 21

B. [3, 9, 15, 21, 27, 33]

C. [13, 14, 15, 16, 17, 18]

D. [4, 10, 16, 22, 28, 34]

Note

- Notation can change
- Samples are also represented as x_i or $x^{[i]}$
 - Make sure you understand the notation used in a book or paper
- Dimensions of the X matrix can be different
 - E.g., if $X = [x^{(1)}, x^{(2)}, \dots, x^{(n)}]$, then, Y is a row vector.



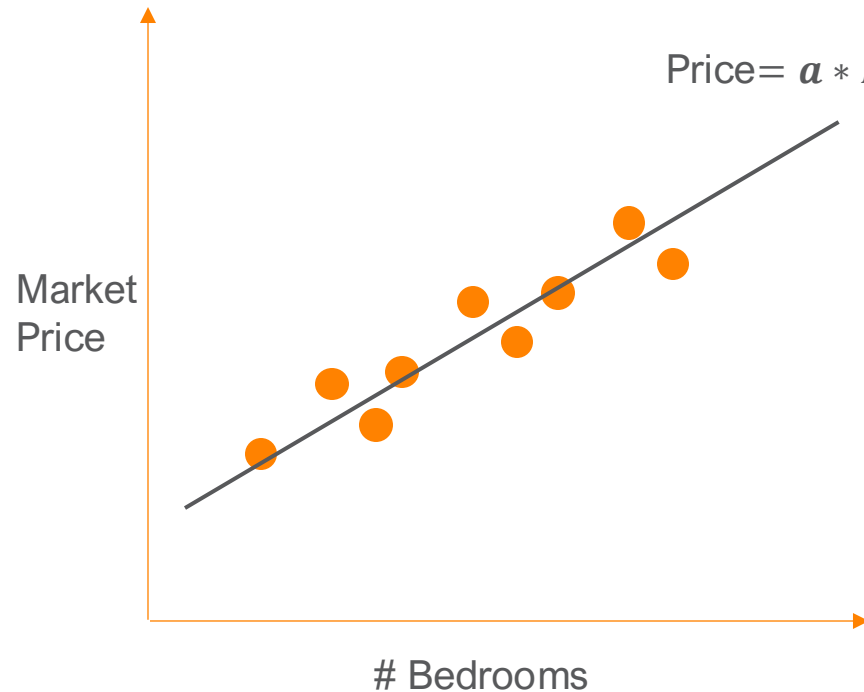
Recap

- **Training example:** A row in the table representing the dataset. Synonymous to an observation, training record, training instance, training sample (in some contexts, sample refers to a collection of training examples)
- **Feature:** a column in the table representing the dataset. Synonymous to predictor, variable, input, attribute, covariate.
- **Targets/Labels:** What we want to predict. Synonymous to outcome, output, ground truth, response variable, dependent variable, (class) label (in classification).
- **Output/prediction:** use this to distinguish from targets; here, means output from the model.

Let's talk again about lines.

How do computers learn?

Housing Market Price Prediction



- Model parameters: a and b
- Find the values for a and b that best fit the data.
 - E.g., linear regression

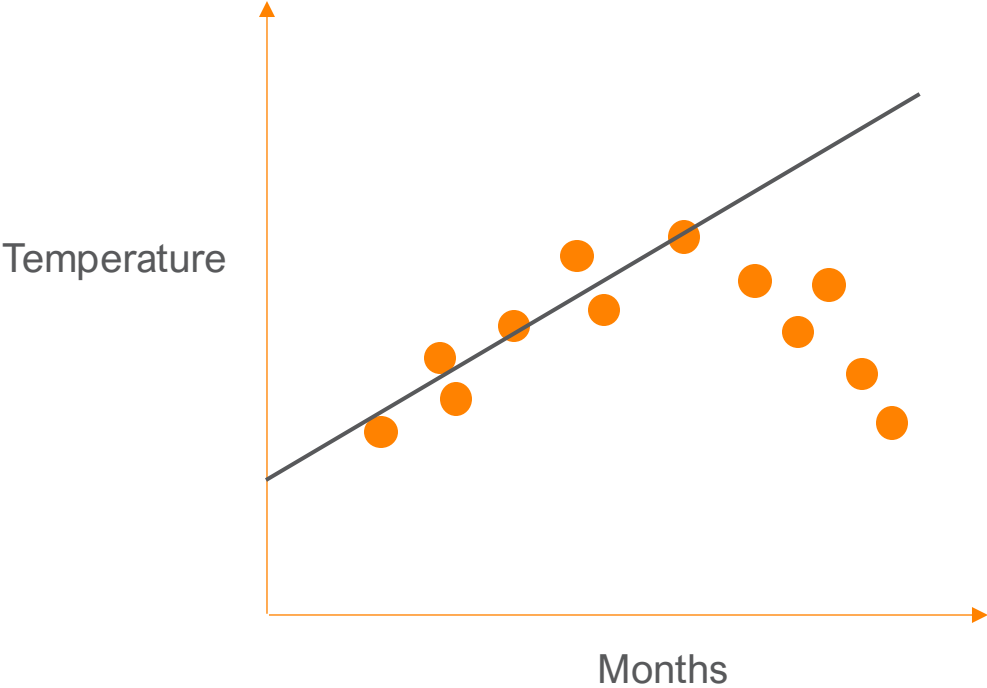
This is a machine-learning model!

How do computers learn?

Housing Market Price Prediction



Season Temperature Prediction

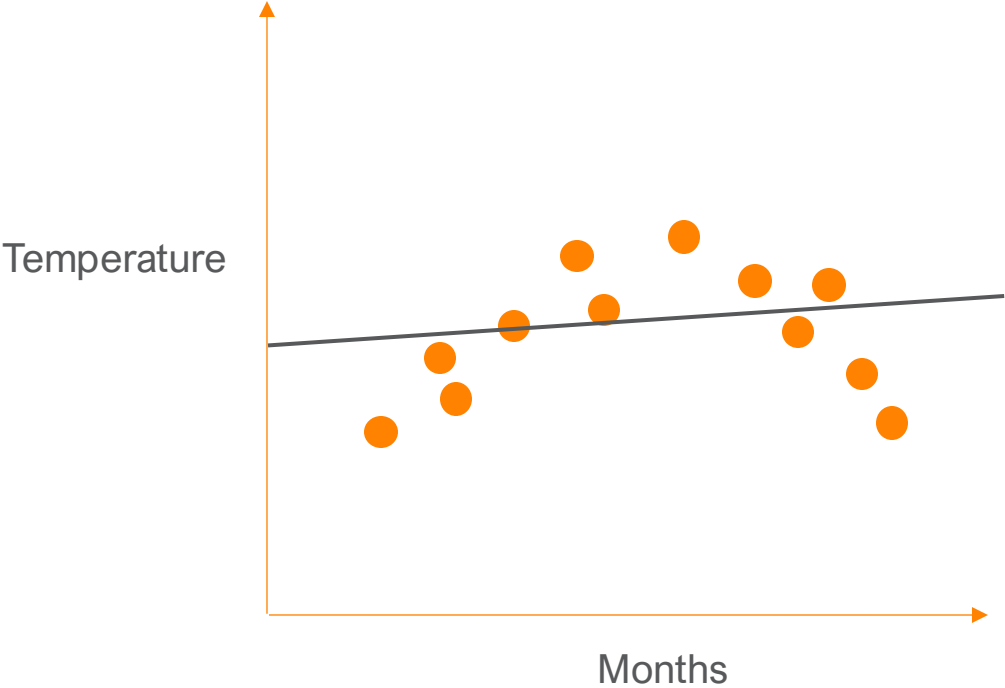


How do computers learn?

Housing Market Price Prediction



Season Temperature Prediction

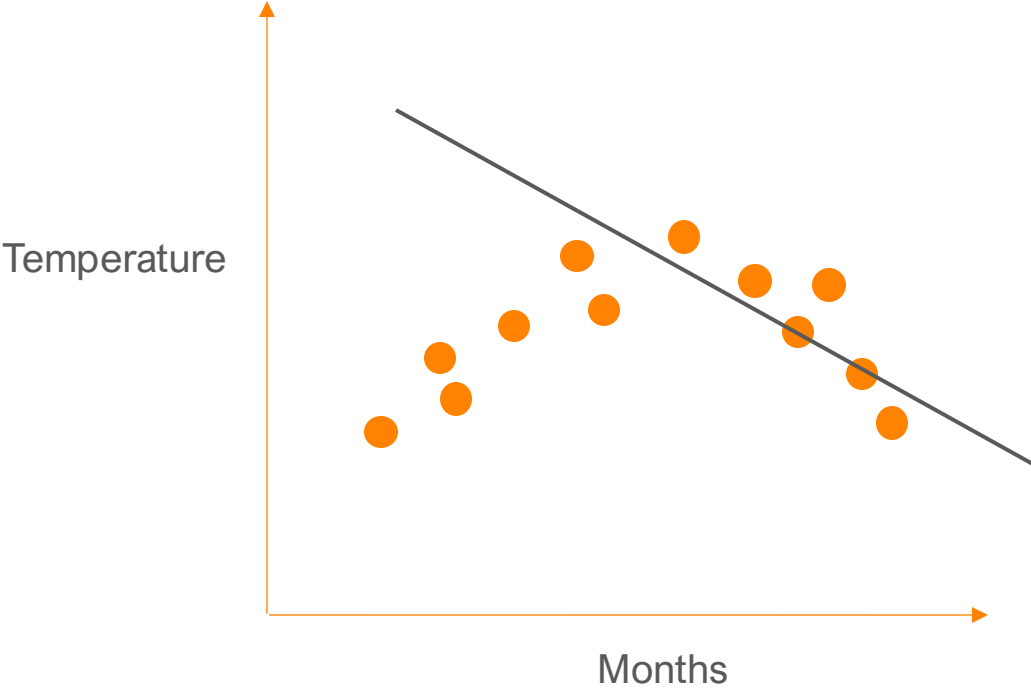


How do computers learn?

Housing Market Price Prediction



Season Temperature Prediction

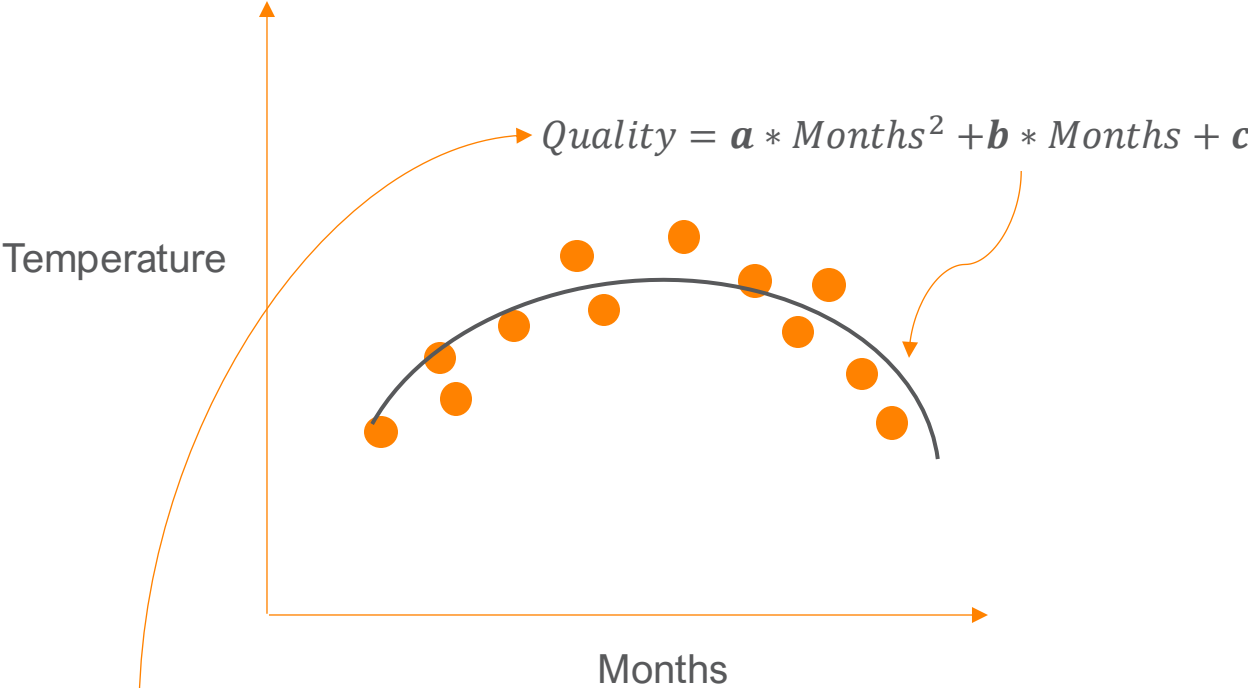


How do computers learn?

Housing Market Price Prediction

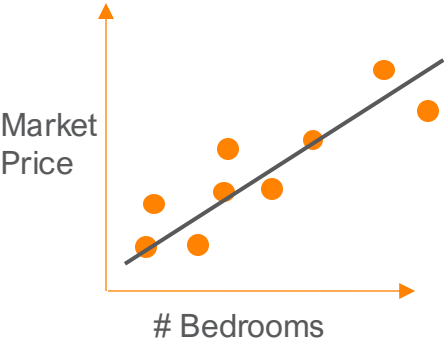


Season Temperature Prediction

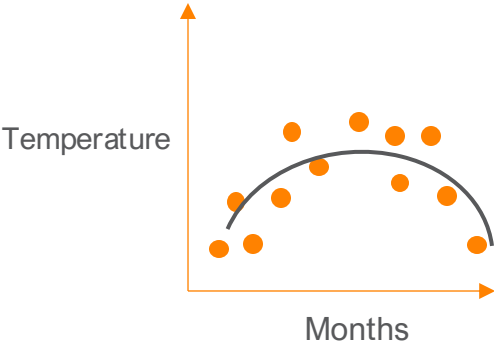


Find the values for parameters **a**, **b**, and **c** that best fit the data.

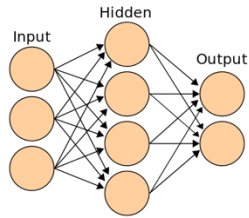
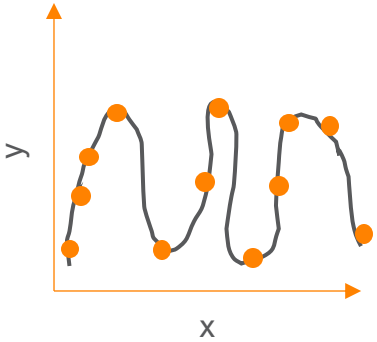
Parameters and Capacity



$$\text{Price} = a * \text{Beds} + b$$



$$\text{Quality} = a * \text{Months}^2 + b * \text{Months} + c$$



Points to take home:

- The computer tries multiple curve parameters to search for the parameter set that best fits the data.
 - Minimizes or maximizes an objective function
- The more parameters, the better we can fit the data... more capacity.

As we increase the number of parameters...

Numpy

<https://www.numpy.org>

Numpy (<https://numpy.org/>)

- General-purpose open-source array-processing library
- High-performance N-dimensional array objects
 - Optimized C code
- Comprehensive built-in functions and random generators
 - Statistics, linear algebra, Fourier transform, etc.
- Fundamental package for scientific computing in Python

NumPy

[Install](#) [Documentation](#) [Learn](#) [Community](#) [About Us](#) [News](#) [Contribute](#) [English](#) ▾



The fundamental package for scientific computing with Python

LATEST RELEASE: NUMPY 2.0. [VIEW ALL RELEASES](#)

NumPy 2.0 released! [2024-06-17](#)

Powerful N-dimensional arrays

Fast and versatile, the NumPy vectorization, indexing, and broadcasting concepts are the de-facto standards of array computing today.

Numerical computing tools

NumPy offers comprehensive mathematical functions, random number generators, linear algebra routines, Fourier transforms, and more.

Open source

Distributed under a liberal [BSD license](#), NumPy is developed and maintained [publicly on GitHub](#) by a vibrant, responsive, and diverse [community](#).

Interoperable

NumPy supports a wide range of hardware and computing platforms, and plays well with distributed, GPU, and sparse array libraries.

Performant

The core of NumPy is well-optimized C code. Enjoy the flexibility of Python with the speed of compiled code.


Easy to use

NumPy's high level syntax makes it accessible and productive for programmers from any background or experience level.

<https://numpy.org/>

Scipy

Install Documentation Community About Us Contribute

SciPy  Fundamental algorithms for scientific computing in Python

[GET STARTED](#)

SciPy 1.14.0 released! [2024-06-24](#)

Fundamental algorithms

SciPy provides algorithms for optimization, integration, interpolation, eigenvalue problems, algebraic equations, differential equations, statistics and many other classes of problems.

Broadly applicable

The algorithms and data structures provided by SciPy are broadly applicable across domains.

Foundational

Extends NumPy providing additional tools for array computing and provides specialized data structures, such as sparse matrices and k-dimensional trees.

Performant

SciPy wraps highly-optimized implementations written in low-level languages like Fortran, C, and C++. Enjoy the flexibility of Python with the speed of compiled code.

Easy to use

SciPy's high level syntax makes it accessible and productive for programmers from any background or experience level.

Open source

Distributed under a liberal [BSD license](#), SciPy is developed and maintained [publicly on GitHub](#) by a vibrant, responsive, and diverse [community](#).

<https://scipy.org/>

Notebook Time

Next Lecture

- We will build our first machine learning pipeline/model with Scikit-Learn