

10-315Machine LearningProblem Formulation

Instructor: Pat Virtue

Today

Autoencoder (Aliens) (previous slides)

Features

ML Problem Formulation

- Task input and output
- Task, Performance, Experience
- Data and notation
- Examples: Iris Classification and Car Price Regression

ML Training and Models

- Linear
- Memorization
- Nearest Neighbor



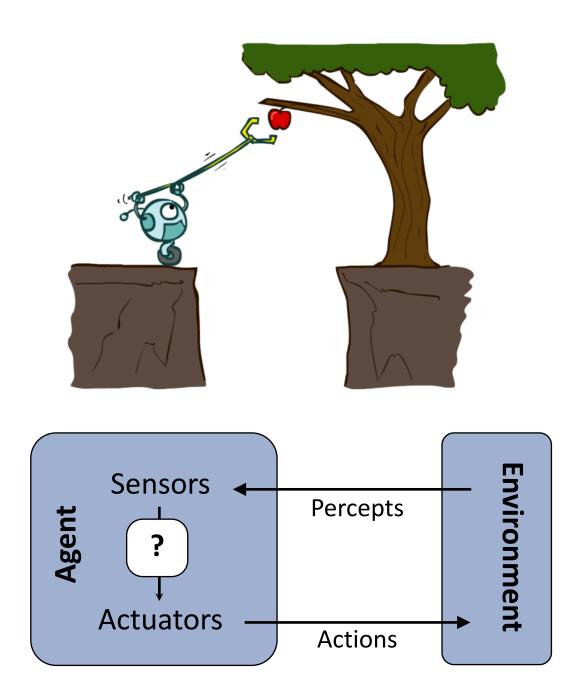
ML Problem Formulation

Agents

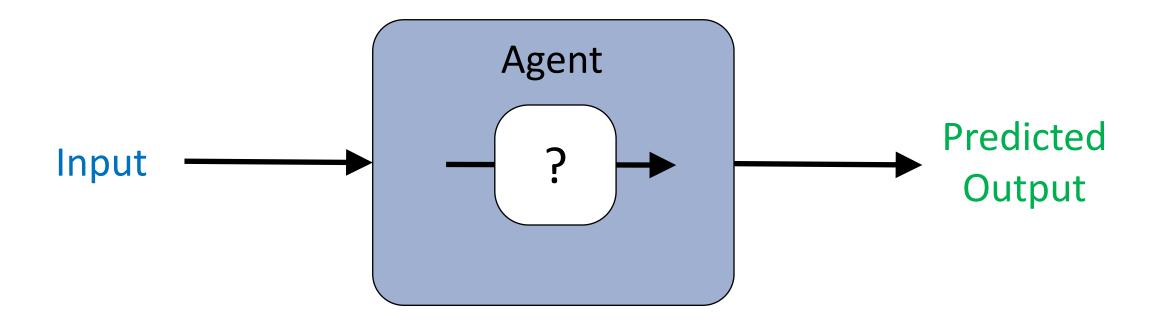
An **agent** is an entity that *perceives* and *acts*.

Actions can have an effect on the environment.

The specific sensors and actuators affect what the agent is capable of perceiving and what actions it is capable of taking



Agent: Simple Input/Output Task



Task Input and Output

Input	Task	Output
Petal measurements	Iris classification	Category
Time of day	Traffic prediction	Traffic Volume
Image	Image classification	Category
Image	Image denoising	Image
Text	Text to image generation	Image
???	Face generation	Image

Task: Face Generation

https://thispersondoesnotexist.com/



Machine Learning Problem Formulation

Three components <*T,P,E>*:

- 1. Task, *T*
- 2. Performance measure, P
- 3. Experience, E

Definition of learning:

A computer program **learns** if its performance at tasks in T, as measured by P, improves with experience E

Machine Learning Problem Formulation

Task

Formalize the task as a mapping from input to output

Experience

Data! Task experience examples will usually be pairs: (input, measured output)

Performance measure

Objective function that gives a single numerical value representing how well the system performs for a given dataset

- Classification: error rate
- Regression: mean squared error

Notation

$$h(x) \to \hat{y}$$

$$\mathcal{D} = \{ (x^{(i)}, y^{(i)}) \}_{i=1}^{N}$$

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(y^{(i)} \neq \hat{y}^{(i)})$$

$$\frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \hat{y}^{(i)})^2$$

Slide: CMU ML, Tom Mitchel and Roni Rosenfeld

ML Problem Formulation

Task

Formalize the task as a mapping from input to output

Experience

Data! Task experience examples will usually be pairs: (input, measured output)

Performance measure

Objective function that gives a single numerical value representing how well the system performs for a given dataset

- Classification: error rate
- Regression: mean squared error

Notation alert: Indicator function

$$\mathbb{I}(z) = \mathbf{1}(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$h(x) \rightarrow \hat{y}$$

$$\mathcal{D} = \{ (x^{(i)}, y^{(i)}) \}_{i=1}^{N}$$

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(y^{(i)} \neq \hat{y}^{(i)})$$

$$\frac{1}{N} \sum_{i=1}^{N} (y^{(i)} - \hat{y}^{(i)})^2$$

Slide: CMU ML, Tom Mitchel and Roni Rosenfeld

Experience: Data and Notation

Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
2	5.9	3.0	5.1	1.8



Assume samples in data are i.i.d.

from sklearn import datasets

iris = datasets.load_iris()

X = iris.data

y = iris.target

Dataset notation

$$\mathcal{D} = \left\{ \left(y^{(i)}, \mathbf{x}^{(i)} \right) \right\}_{i=1}^{N}$$

$$= \left\{ \left(y^{(i)}, x_1^{(i)}, x_2^{(i)}, x_3^{(i)}, x_4^{(i)} \right) \right\}_{i=1}^{N}$$

Linear algebra can represent all data

$$\mathbf{y} \in \{0,1,2\}^N$$

 $X \in \mathbb{R}^{N \times 4}$ (design matrix)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
2	5.9	3.0	5.1	1.8

Assume samples in data are i.i.d.

from sklearn import datasets

iris = datasets.load_iris()

X = iris.data

y = iris.target

Dataset notation

$$\mathcal{D} = \{ (y^{(i)} | \mathbf{x}^{(i)}) \}_{i=1}^{N}$$

$$= \{ (y^{(i)}, x_1^{(i)}, x_2^{(i)}, x_3^{(i)}, x_4^{(i)}) \}_{i=1}^{N}$$

Data point $i = 6: (y^{(6)}, \mathbf{x}^{(6)})$

Spe	cies	Sepal Length	Sepal Width	Petal Length	Petal Width
C)	4.3	3.0	1.1	0.1
C)	4.9	3.6	1.4	0.1
C)	5.3	3.7	1.5	0.2
1	L	4.9	2.4	3.3	1.0
1	L	5.7	2.8	4.1	1.3
1	L	6.3	3.3	4.7	1.6
2	<u>)</u>	5.9	3.0	5.1	1.8

Assume samples in data are i.i.d.

from sklearn import datasets

iris = datasets.load_iris()

X = iris.data

y = iris.target

Dataset notation

$$\mathcal{D} = \left\{ \left(y^{(i)}, \mathbf{x}^{(i)} \right) \right\}_{i=1}^{N}$$

$$= \left\{ \left(y^{(i)}, x_1^{(i)}, x_2^{(i)}, x_3^{(i)}, x_4^{(i)} \right) \right\}_{i=1}^{N}$$

Linear algebra can represent all data

$$\mathbf{y} \in \{0,1,2\}^N$$
 $X \in \mathbb{R}^{N \times 4}$ (design matrix)

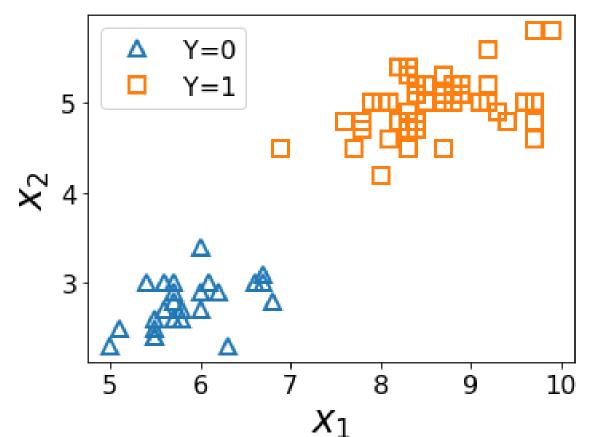
Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
2	5.9	3.0	5.1	1.8

Task: Classification

ML Task: Classification

Predict species label from first two input measurements

$$h(\mathbf{x}) \to \hat{y}$$





Species	Sepal Length	Sepal Width
0	4.3	3.0
0	4.9	3.6
0	5.3	3.7
1	4.9	2.4
1	5.7	2.8
1	6.3	3.3

Classification

Iris data example

$$\mathbb{I}(z) = \mathbf{1}(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^{N}$$
, where $\mathbf{x}^{(i)} \in \mathbb{R}^4$, $y^{(i)} \in \{0, 1, 2\}$

Predict species label from input measurements

$$h(\mathbf{x}) \to \hat{y}$$

Performance measure?

Classification error rate

- Fraction of times $y \neq \hat{y}$ in a given dataset

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
2	5.9	3.0	5.1	1.8

ML Tasks

Supervised learning: Pairs of input and output in training data

$$\mathcal{D} = \left\{ \left(\mathbf{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^{N} \qquad h(\mathbf{x}) \to \hat{y}$$

Classification

- Output labels
- $y \in \mathcal{Y}$, where \mathcal{Y} is discrete and order of values has no meaning

Regression

- Output values
- $y \in \mathcal{Y}$, where \mathcal{Y} is usually continuous, order of values has meaning

Unsupervised Tasks

ML Tasks

Unsupervised learning

$$\mathcal{D} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^{N} \quad h(\mathbf{x}) \to ???$$

- Training data has no output values
- Tasks can vary
- Often used to organize data for future (minimally) supervised learning

Task: Face Generation

https://thispersondoesnotexist.com/



ML Tasks

Unsupervised learning

$$\mathcal{D} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^{N} \quad h(\mathbf{x}) \to ???$$

- Training data has no output values
- Tasks can vary
- Often used to organize data for future (minimally) supervised learning

Example: Unsupervised autoencoder \rightarrow Random image generation

$$\mathbf{x} \to \boxed{h(\mathbf{x})} \to \hat{\mathbf{x}}$$

$$\mathbf{x} \to \boxed{f(\mathbf{x})} \to \mathbf{z} \to \boxed{g(\mathbf{z})} \to \hat{\mathbf{x}}$$

$$\mathbf{z} \to \boxed{g(\mathbf{z})} \to \hat{\mathbf{x}}$$

ML Tasks

Unsupervised learning

$$\mathcal{D} = \left\{ \mathbf{x}^{(i)} \right\}_{i=1}^{N} \quad h(\mathbf{x}) \to ???$$

- Training data has no output values
- Tasks can vary
- Often used to organize data for future (minimally) supervised learning

Example: Text Generation

Vocab pause

Task

- Prediction
- Inference
- Hypothesis function
- Classification
- Regression

Experience/Data

Input

- Input feature
- Measurement
- Attribute

Output

- Target
- Class/category/label
- True output
- Measured output
- Predicted output

Supervised

Unsupervised

Performance Measure

Objective function

Classification

- Frror rate
- Accuracy rate

Regression

Mean squared error

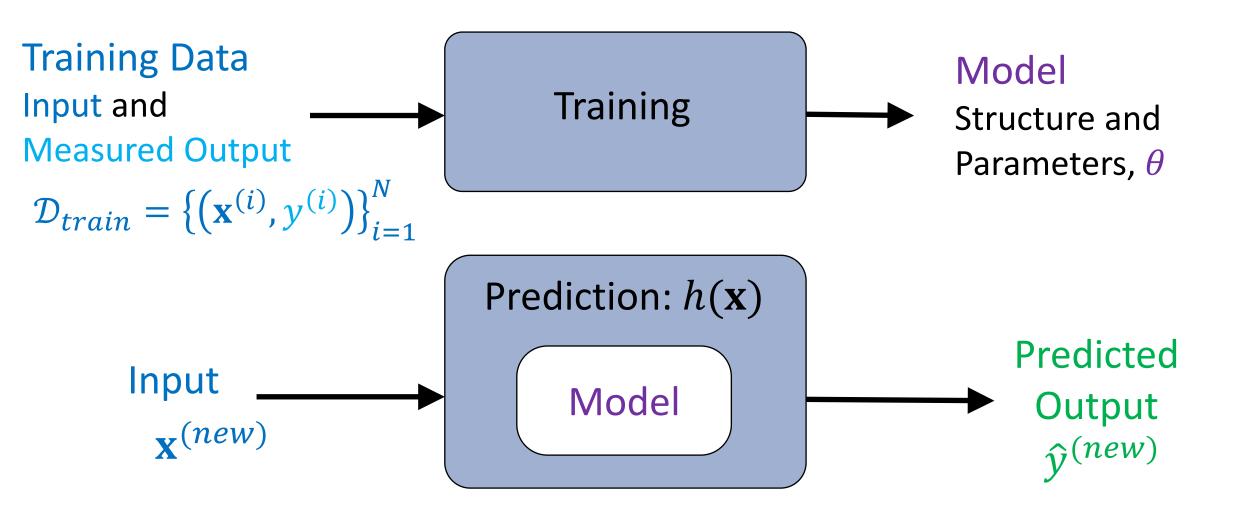
Training

- Model
- Model structure
- Model parameters

Training and ML Models

Machine Learning

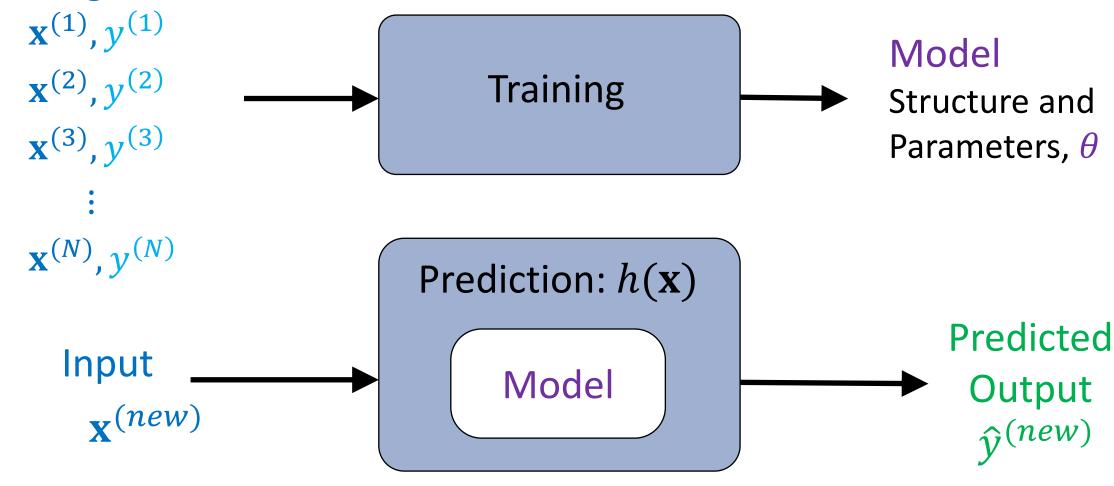
Using (training) data to learn a model that we'll later use for prediction



Machine Learning

Using (training) data to learn a model that we'll later use for prediction

Training Data

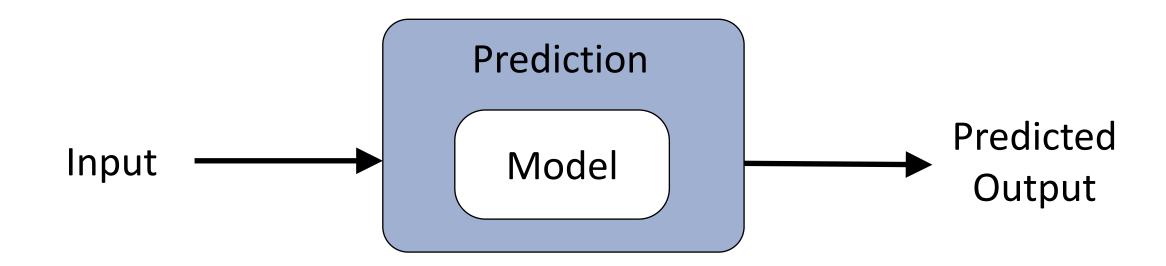


Task: Car Price Prediction

Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Example

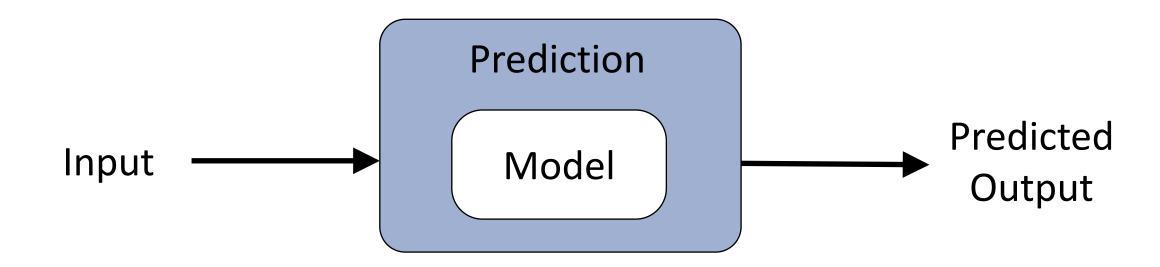
Trying to see how much I should sell my car for.



Task: Car Price Prediction

Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

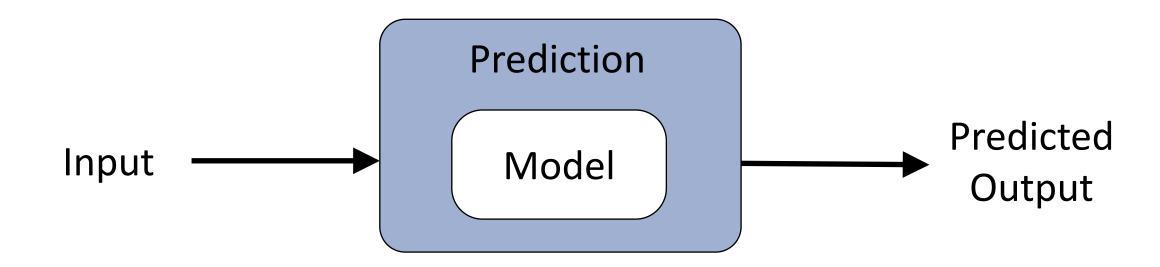
What input features should we use?



Poll 2

Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

What input features should we use?

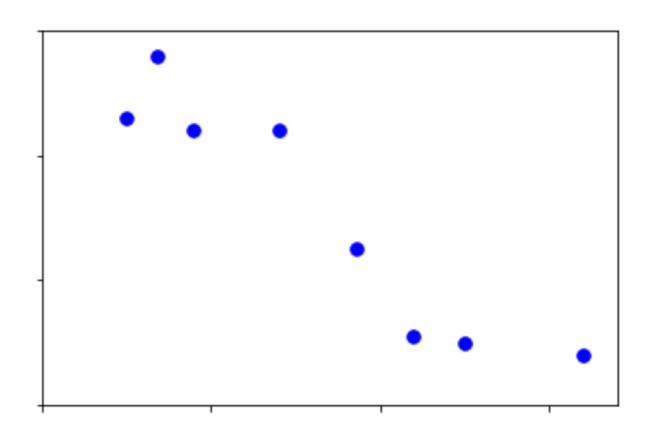


Regression

Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

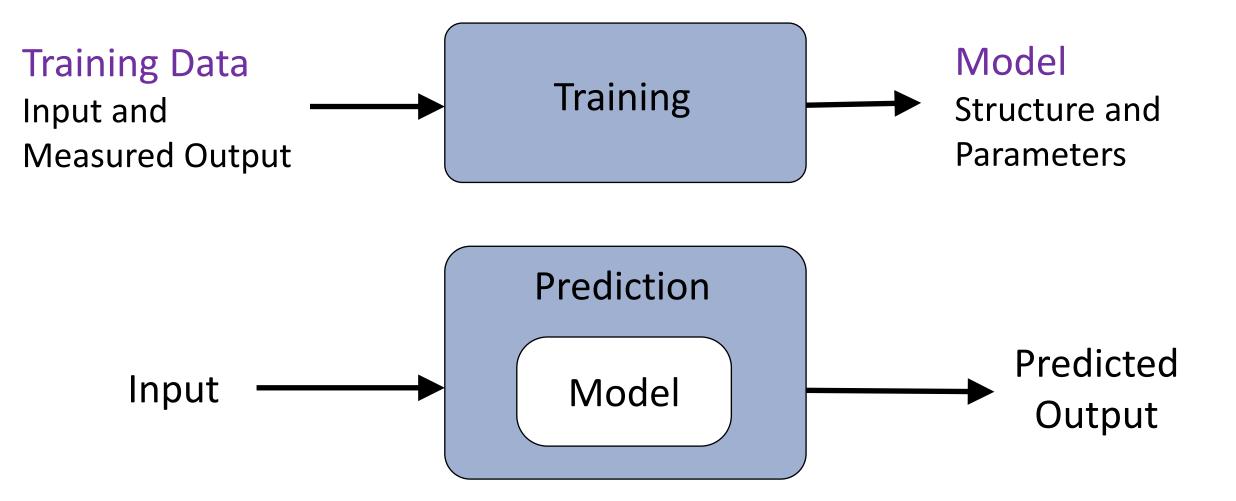
Example

Trying to see how much I should sell my car for.
Looking up data from car websites, I find the mileage for a set of cars and the selling price for each car.



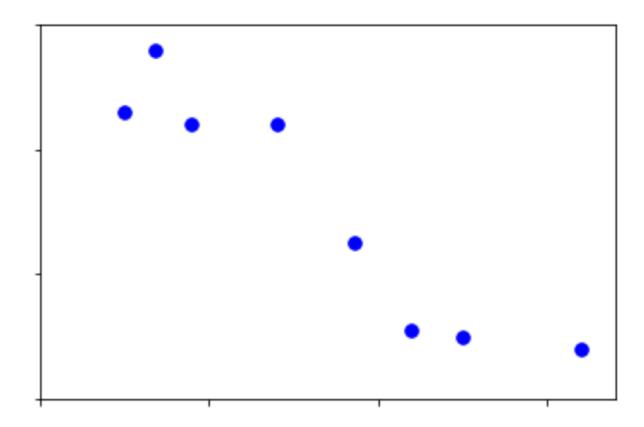
Machine Learning

Using (training) data to learn a model that we'll later use for prediction



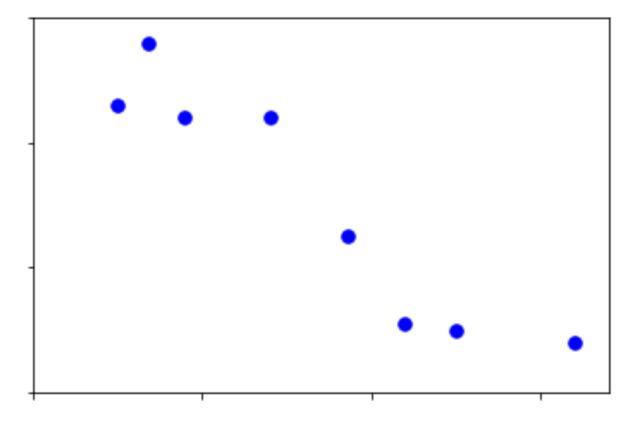
Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Model?



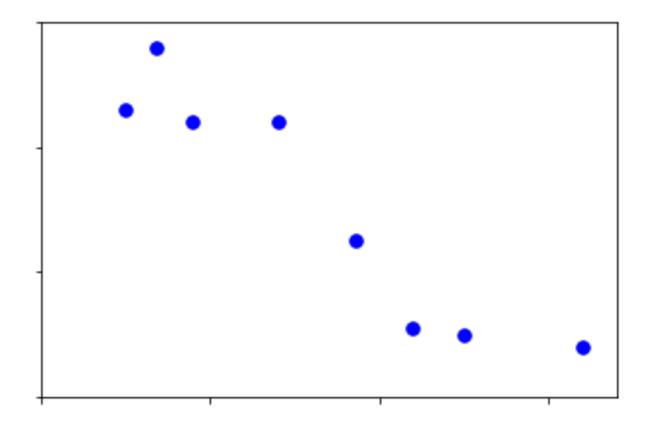
Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Model: Memorization



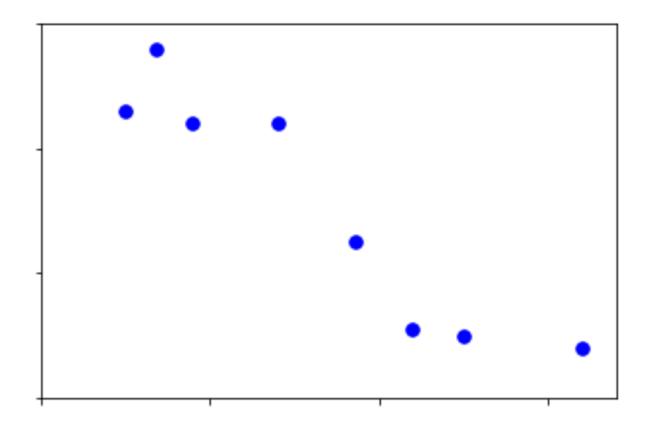
Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Model: Nearest neighbor



Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Model: Linear



Regression: learning a model to predict a numerical output (but not numbers that just represent categories, that would be classification)

Model?

