

Machine Learning Apprentissage Machine INFO415

Dr. Amani RAAD

Amani.raad@ul.edu.lb

2023

Chapter 1

General Concepts

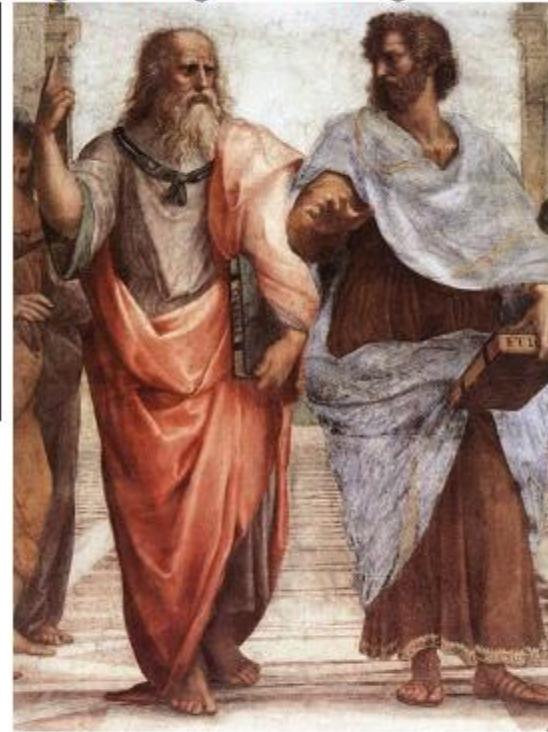
Introduction

Arithmetics, algorithms



An abacus

Logic and logical reasoning



Aristotle

Automata



A writing automaton

Programmable machines



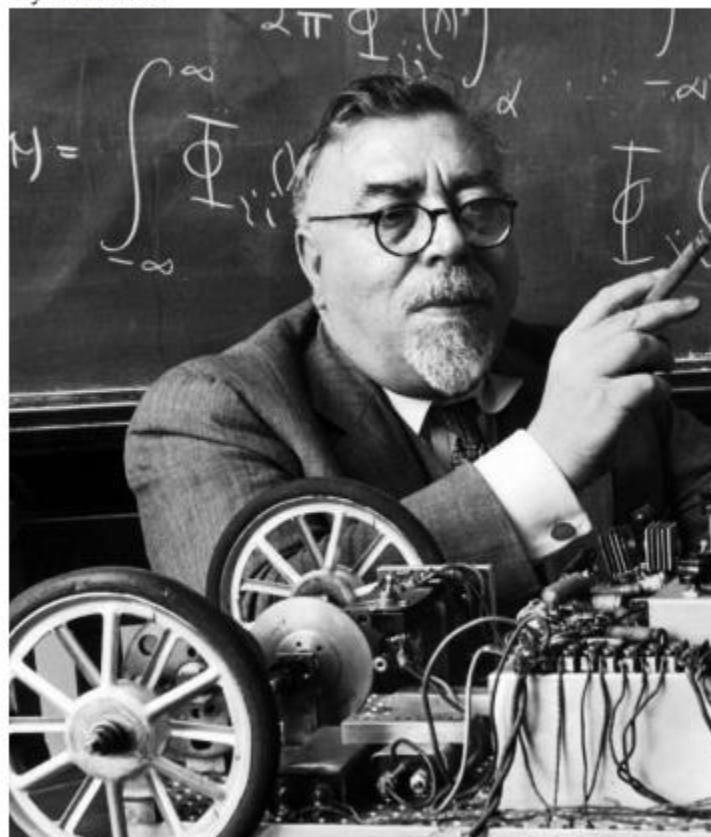
The Jacquard programmable loom

Turing and the science of computation



Alan M. Turing (1912 - 1954)

Cybernetics



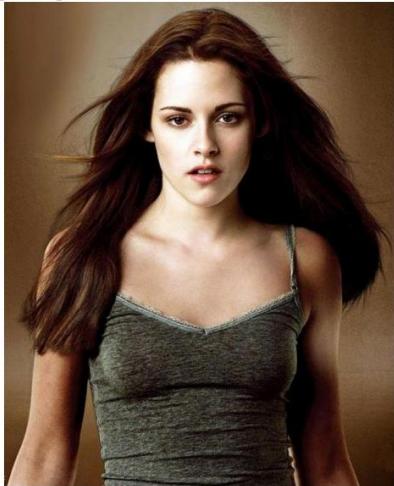
Norbert Wiener (1894 - 1964)

For these problems we do have algorithms

- Sorting → insert sort, bubble sort, Shell sort, radix sort, heapsort, bogosort. . .
- Spectral analysis of periodic signals → FFT
- Database filtering → SQL SELECT

We have no algorithms for

Recognizing a face



We have no algorithms for

Distinguishing between genuine works by a given author and fakes

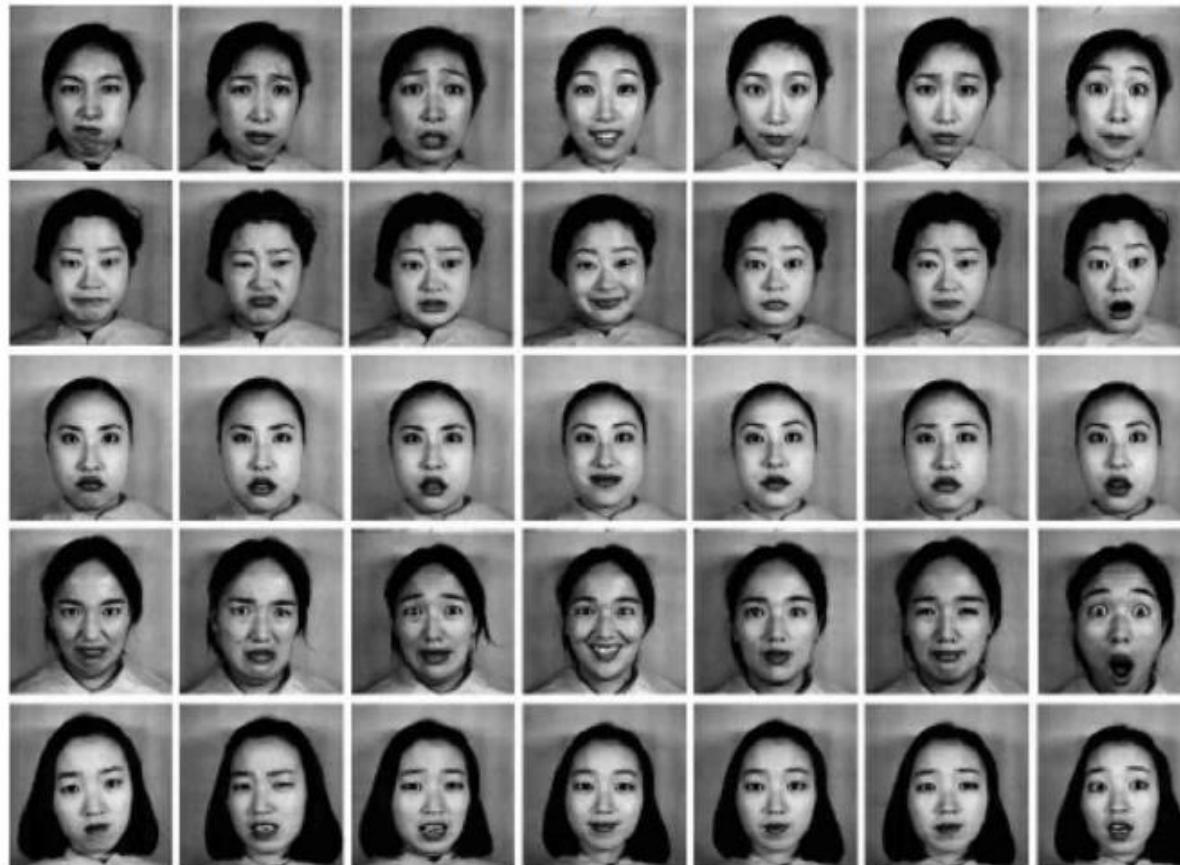


You are authorized
to offer my autobiography
to the great Tell or
any other publisher
on the terms to which
we have agreed.

Most sincerely yours,
Howard R. Hughes

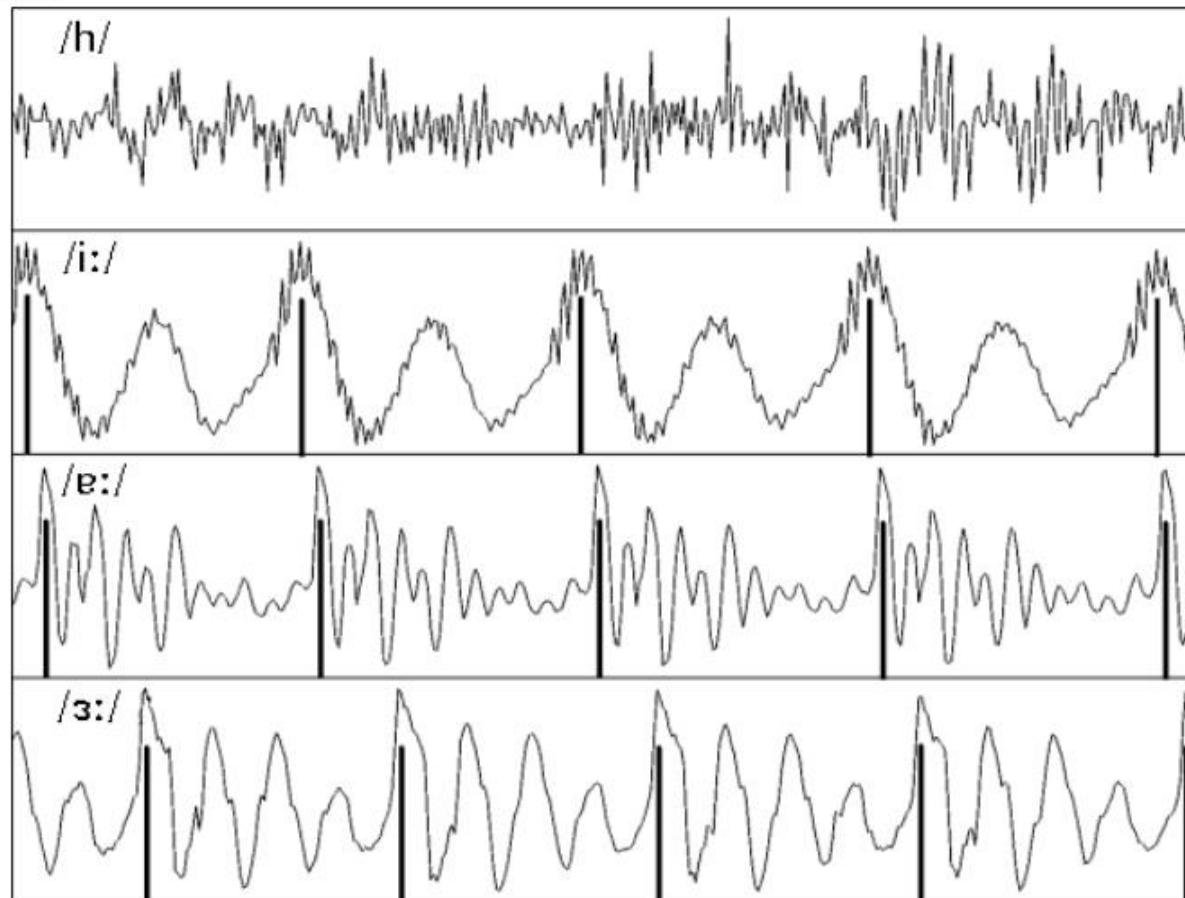
We have no algorithms for

Recognizing emotions from gestures or facial expressions



We have no algorithms for

Understanding speech



We have no algorithms for machine translation

This text has been [automatically translated](#) from Arabic:

Moscow stressed tone against Iran on its nuclear program. He called Russian Foreign Minister Tehran to take concrete steps to restore confidence with the international community, to cooperate fully with the IAEA. Conversely Tehran expressed its willingness

Translate text

شددت موسكو لهجتها ضد إيران بشأن برنامجه النووي. ودعا وزير الخارجية الروسي طهران إلى اتخاذ خطوات ملموسة لاستئناف اللغة مع المجتمع الدولي والتعاون الكامل مع الوكالة الذرية. بالمقابل أبدت طهران استعدادها لاستئناف المسار بعمليات التفتيش المفاجئة بشرط إسقاط مجلس الأمن ملفها النووي.

from Arabic to English BETA

We have no algorithms for

Spam filtering

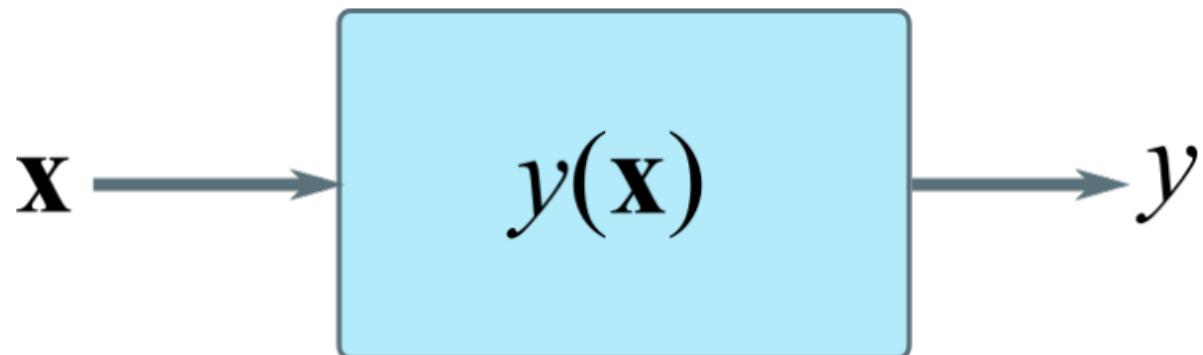


- Many interesting problems are too complex to admit an algorithmic solution or even a complete description
 - For these problems, only data are available
 - Nowadays we have LOTS OF DATA FROM LOTS OF SOURCES
-
- Machine learning is about using data to solve problem

- Instead of writing a program by hand for each specific task, we collect lots of examples that specify the correct output for a given input.
- A machine learning algorithm then takes these examples and produces a program that does the job.
- Massive amounts of computation are now cheaper than paying someone to write a task-specific program.

General structure of a learning machine

A rule that receives an observation \mathbf{x}
and emits an output $y = y(\mathbf{x})$

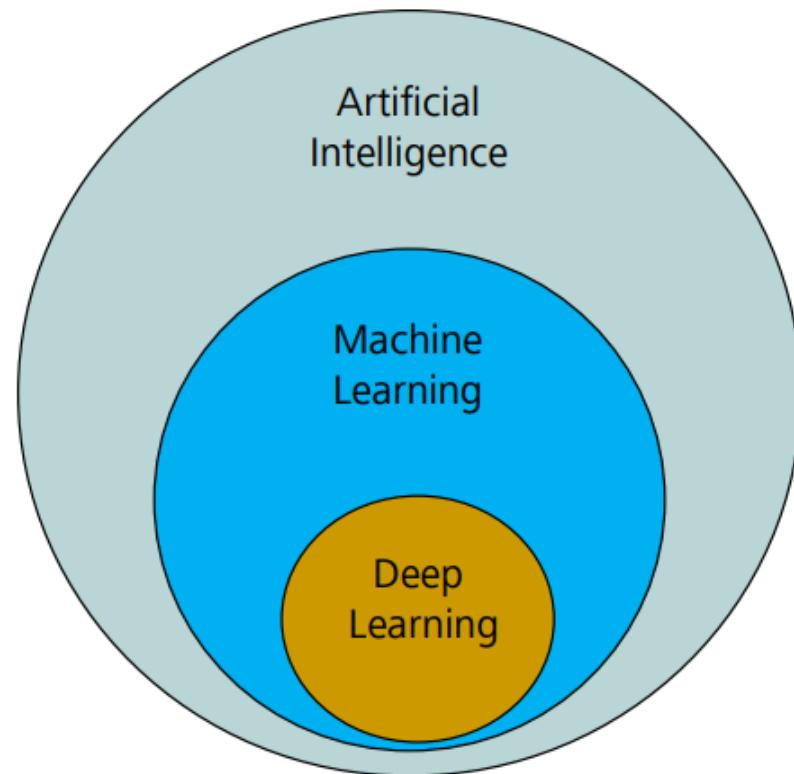


- The rule can be implemented as a software, as a procedure, as a hardware machine, as an embedded computer...
- The rule normally has a predefined structure, but depends on a set of parameters.
- Learning consists in fitting these parameters to a specific problem.

Machine Learning

- Machine Learning is a subset of Artificial Intelligence et is coined by **Arthur Samuel** in 1959, an American pioneer in the field of computer gaming and artificial intelligence
- It stated that
- **“it gives computers the ability to learn without being explicitly programmed”.**

Relationship between AI, Machine Learning and Deep Learning



- **Learning machine** – also “learner”.
Not necessarily a real machine! Maybe a program
- **Task** – a problem to be solved.
We don’t have a description of the problem, but **data**
- **Learning** – adjusting quantities “inside the machine”
(e.g., algorithm parameters) to do a certain task

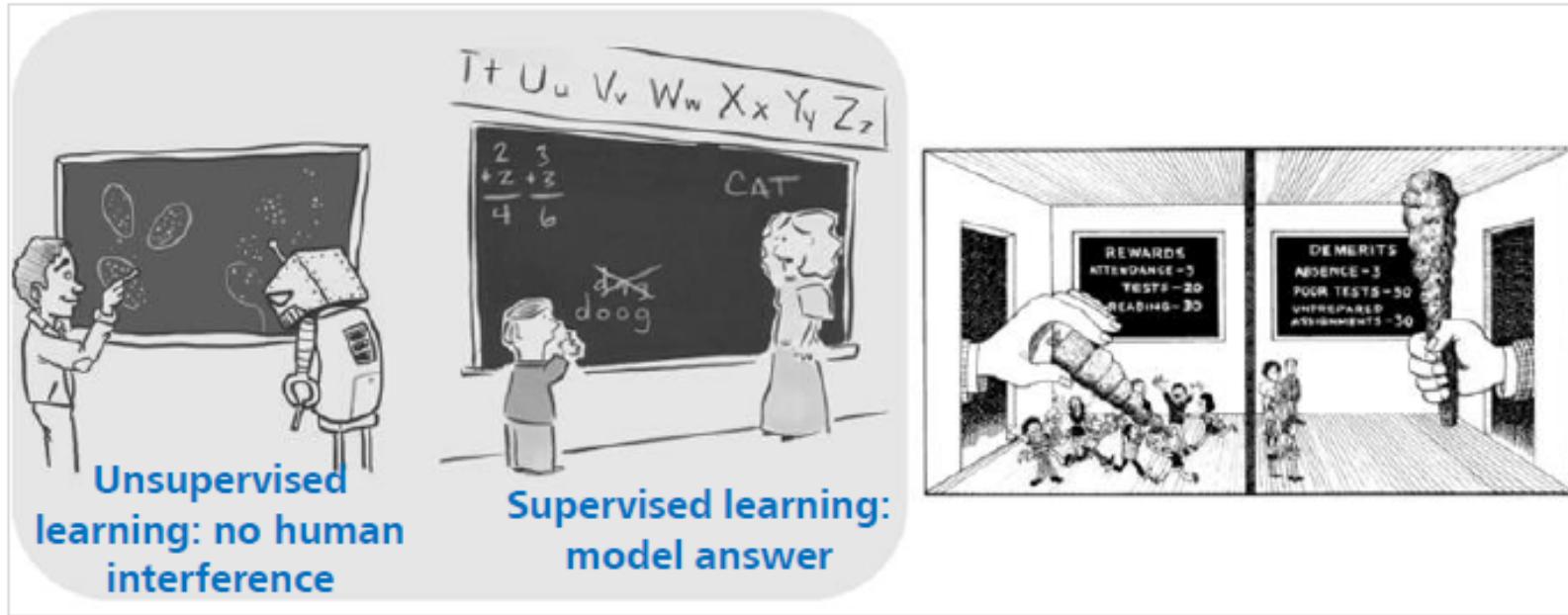
Tasks

- Mapping a stimulus to a response:
supervised learning
- Describing the data:
unsupervised learning

Tasks

- In **supervised tasks**
the data contain both input patterns and output values
- In **unsupervised tasks**
the data contain just input patterns

Types of Machine learning



Traditional machine learning

Reinforcement learning:
reward and punishment

Tasks

- Mapping a stimulus to a response (from input to output):
 - Classification
 - Regression
- Describing the data (from input to a more compact representation of the input itself):
 - Clustering
 - Mapping in lower dimensionality

Supervised learning

- This approach is indeed similar to human learning under the supervision of a teacher. The teacher provides good examples for the student to memorize, and the student then derives general rules from these specific examples.
- In Supervised learning, an AI system is presented with data which is labeled, which means that each data tagged with the correct label.
- Categories: classification, regression

Classification

- When inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more of these classes. This is typically tackled in a supervised way.
- Example: Spam filtering is an example of classification, where the inputs are email (or other) messages and the classes are “spam” and “not spam”.
- Example: Medical and industrial diagnosis : predict heart disease or predict the failure of a machine in industry in the case of labelled data.

Example

Iris (flower) recognition

- The Iris dataset has been in use since 1936
- Collected by botanist Edgar Anderson in 1935
- Used by statistician Sir Ronald A. Fisher in 1936

Sources:

Edgar Anderson (1935). “The irises of the Gaspe Peninsula”. Bulletin of the American Iris Society 59: 25.

Fisher, R.A. (1936). “The Use of Multiple Measurements in Taxonomic Problems”. Annals of Eugenics 7: 179-188.

Download it from the University of California - Irvine repository at:

<http://archive.ics.uci.edu/ml/datasets/Iris>



Iris Setosa



Iris Virginica



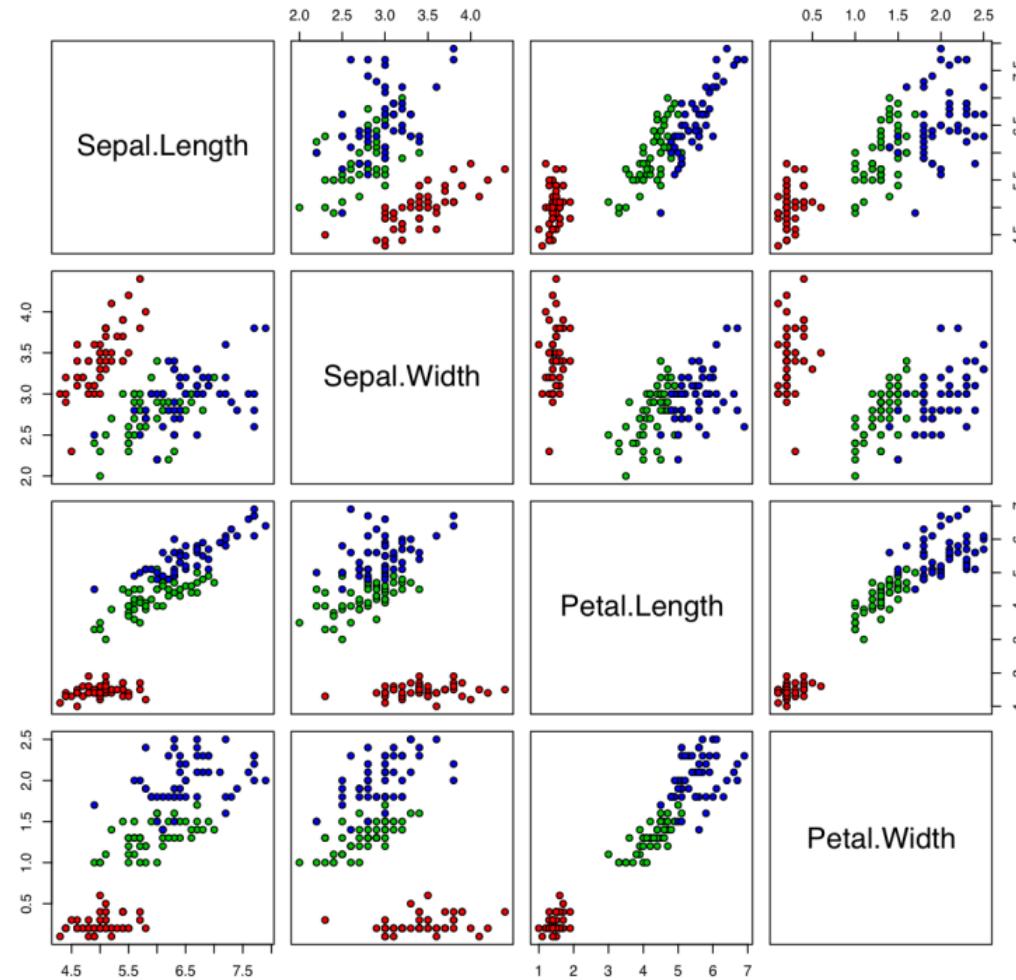
Iris Versicolor

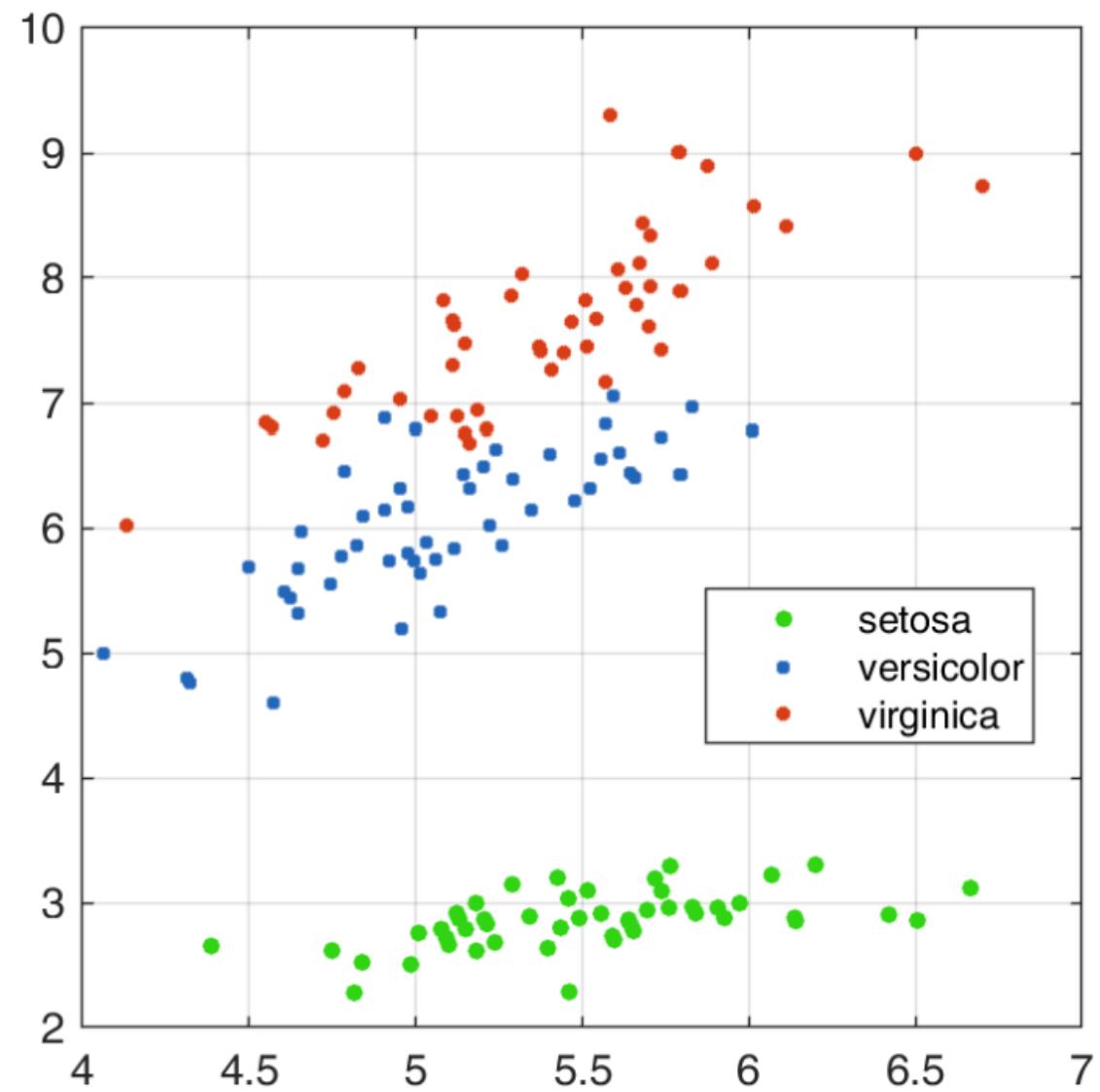
180 MULTIPLE MEASUREMENTS IN TAXONOMIC PROBLEMS

Table I

<i>Iris setosa</i>				<i>Iris versicolor</i>				<i>Iris virginica</i>			
Sepal length	Sepal width	Petal length	Petal width	Sepal length	Sepal width	Petal length	Petal width	Sepal length	Sepal width	Petal length	Petal width
5.1	3.5	1.4	0.2	7.0	3.2	4.7	1.4	6.3	3.3	6.0	2.5
4.9	3.0	1.4	0.2	6.4	3.2	4.5	1.5	5.8	2.7	5.1	1.9
4.7	3.2	1.3	0.2	6.9	3.1	4.9	1.5	7.1	3.0	5.9	2.1
4.6	3.1	1.5	0.2	5.5	2.3	4.0	1.3	6.3	2.9	5.6	1.8
5.0	3.6	1.4	0.2	6.5	2.8	4.6	1.5	6.5	3.0	5.8	2.2
5.4	3.9	1.7	0.4	5.7	2.8	4.5	1.3	7.6	3.0	6.6	2.1
4.6	3.4	1.4	0.3	6.3	3.3	4.7	1.6	4.9	2.5	4.5	1.7
5.0	3.4	1.5	0.2	4.9	2.4	3.3	1.0	7.3	2.9	6.3	1.8
4.4	2.9	1.4	0.2	6.6	2.9	4.6	1.3	6.7	2.5	5.8	1.8
4.9	3.1	1.5	0.1	5.2	2.7	3.9	1.4	7.2	3.6	6.1	2.5
5.4	3.7	1.5	0.2	5.0	2.0	3.5	1.0	6.5	3.2	5.1	2.0
4.8	3.4	1.6	0.2	5.9	3.0	4.2	1.5	6.4	2.7	5.3	1.9
4.8	3.0	1.4	0.1	6.0	2.2	4.0	1.0	6.8	3.0	5.5	2.1
4.3	3.0	1.1	0.1	6.1	2.9	4.7	1.4	5.7	2.5	5.0	2.0
5.8	4.0	1.2	0.2	5.6	2.9	3.6	1.3	5.8	2.8	5.1	2.4
5.7	4.4	1.5	0.4	6.7	3.1	4.4	1.4	6.4	3.2	5.3	2.3
5.4	3.9	1.3	0.4	5.6	3.0	4.5	1.5	6.5	3.0	5.5	1.8
5.1	3.5	1.4	0.3	5.8	2.7	4.1	1.0	7.7	3.8	6.7	2.2
5.7	3.8	1.7	0.3	6.2	2.2	4.5	1.5	7.7	2.6	6.9	2.3
5.1	3.8	1.5	0.3	5.6	2.5	3.9	1.1	6.0	2.2	5.0	1.5
5.4	3.4	1.7	0.2	5.9	3.2	4.8	1.8	6.9	3.2	5.7	2.3
5.1	3.7	1.5	0.4	6.1	2.9	4.9	1.2	5.6	2.8	4.9	2.0

Iris Data (red=setosa,green=versicolor,blue=virginica)





Example 2- MNIST DATA

Recognizing handwritten digits

Images of size 28x28 pixels, grayscale

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Example 3

- A wall-following robot has to make decisions as to the direction to take, depending on a circular array of ultrasound sensors
- The robot has 24 such sensors evenly spread over 360 degrees.
- Possible directions are:

Sharp-Left-Turn	Slight-Left-Turn	Move-forward	Slight-Right-Turn	Sharp-Right-Turn
-----------------	------------------	--------------	-------------------	------------------

The scitos G5 robot is a multipurpose, modular platform for robotic research and development.

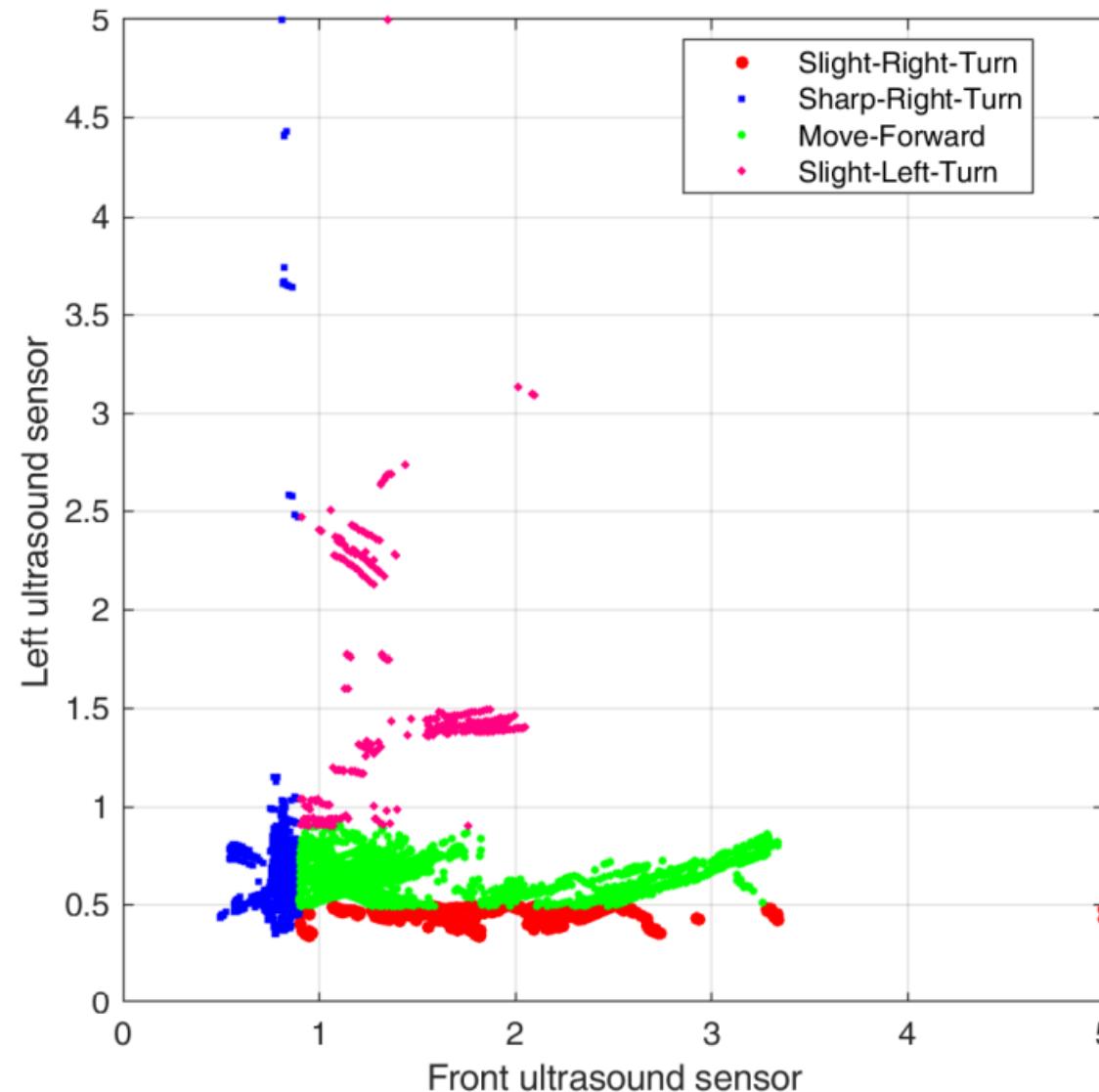


<https://www.metralabs.com/en/mobile-robot-scitos-g5/>

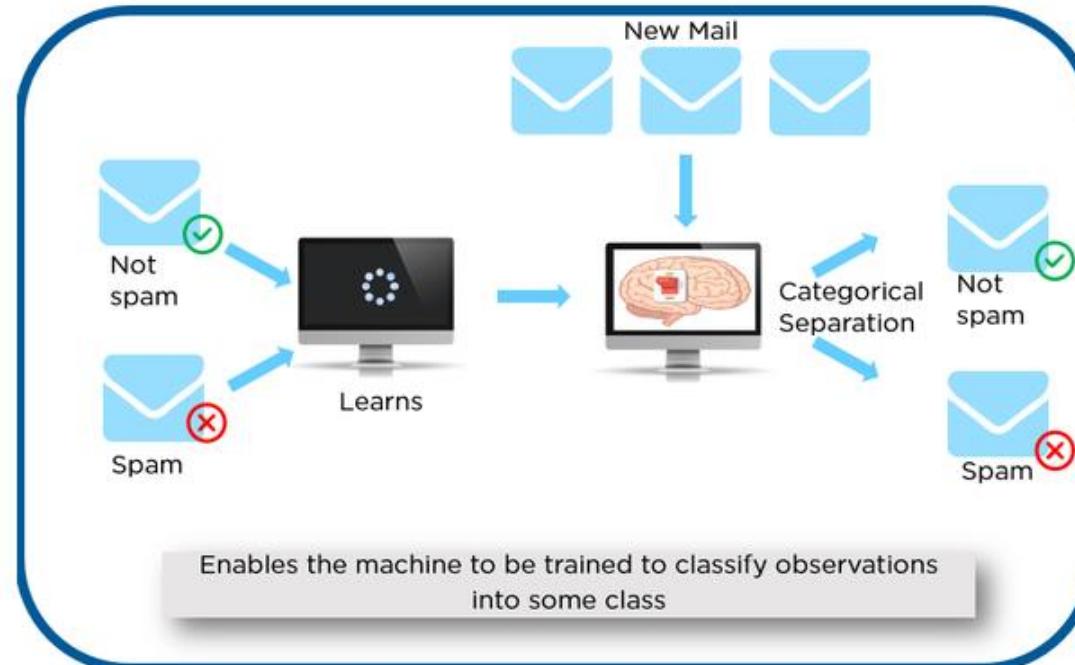
The ultrasonic sensor's output is available as a voltage in the range 0...5V



Minimum readings from two groups of sensors, on the forward and on the left.
Colors correspond to directions to take.



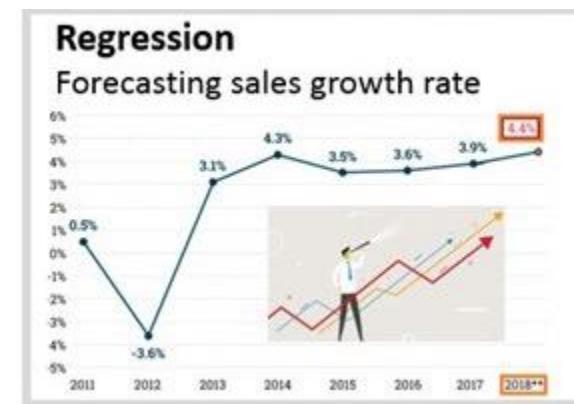
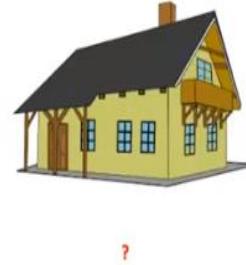
Example 4



Example of Supervised Learning

Regression

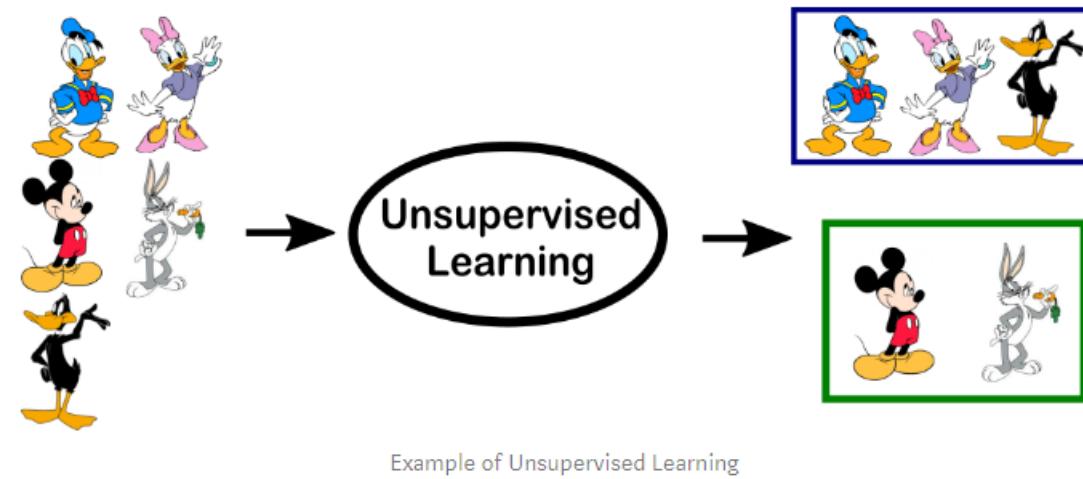
- **Regression** : Which is also a supervised problem, A case when the outputs are continuous rather than discrete.
- Example: to predict the price of a house



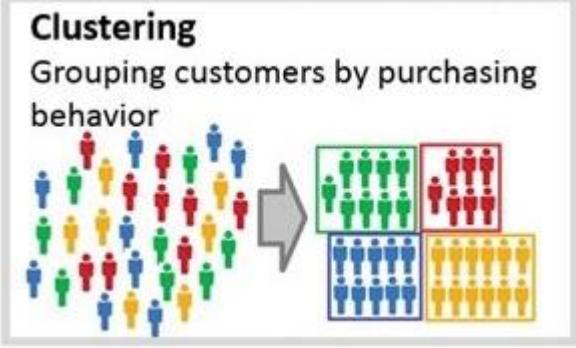
Unsupervised learning

- In unsupervised learning, an AI system is presented with unlabeled, uncategorized data and the system's algorithms act on the data without prior training.
- Category: Clustering, Association

Unsupervised learning



Clustering

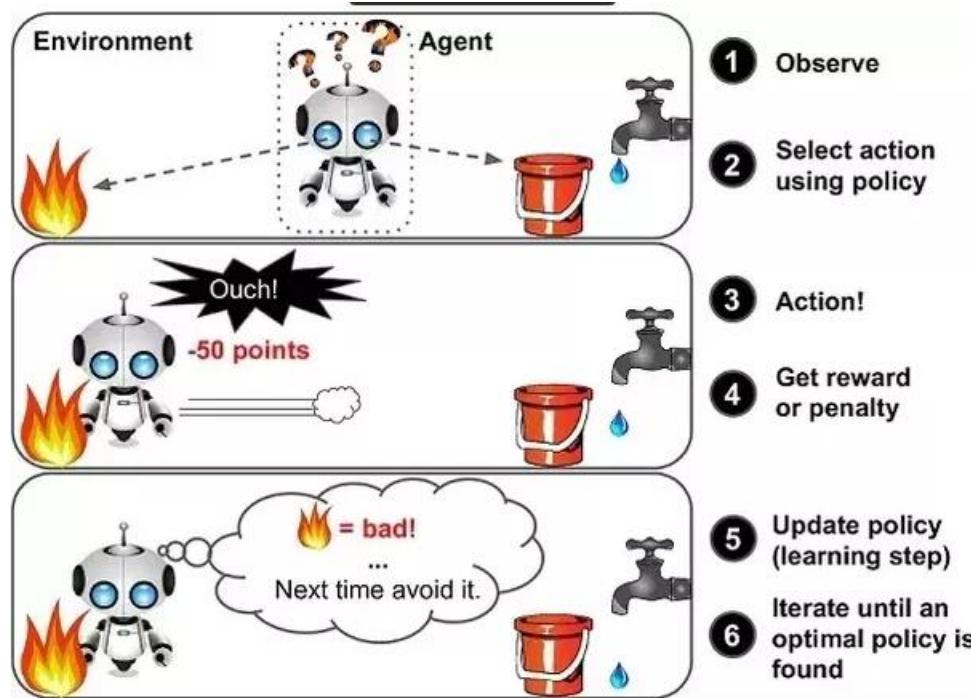


- When a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task.
- **Example : YouTube:** uses our search history or watched history and suggests videos we might like. Feature data set for Facebook contains is people we follow, pages we follow, comments we input, photos or videos we like, pictures or photos we tag at.
- **Example: Banking:** With the help of clustering, insurance companies can find fraud, acknowledge customers about it and understand policies brought by the customer.
- **Example: Google** is one of the search engine people uses. Let's take an example when we search for some information like pet store in the area, Google will provide us different options.

Reinforcement learning

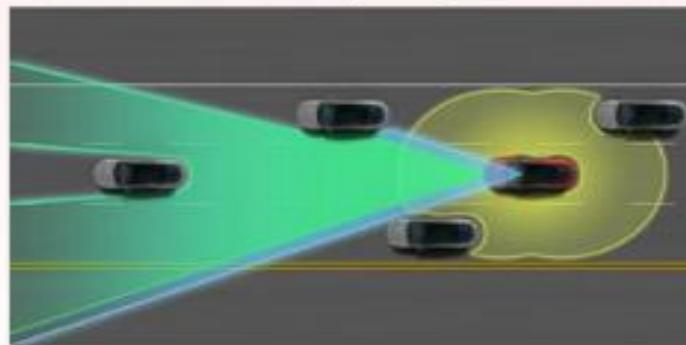
- A reinforcement learning algorithm, or agent, learns by interacting with its environment. The agent receives rewards by performing correctly and penalties for performing incorrectly. The agent learns without intervention from a human by maximizing its reward and minimizing its penalty. It is a type of dynamic programming that trains algorithms using a system of reward and punishment.
- In the human world, it is just like learning by trial and error. Errors help you learn because they have a penalty added (cost, loss of time, regret, pain, and so on), teaching you that a certain course of action is less likely to succeed than others.
- An interesting example of reinforcement learning occurs when computers learn to play video games by themselves.

Reinforcement learning



Self Driving cars

- Autonomous cars(self-driving car)



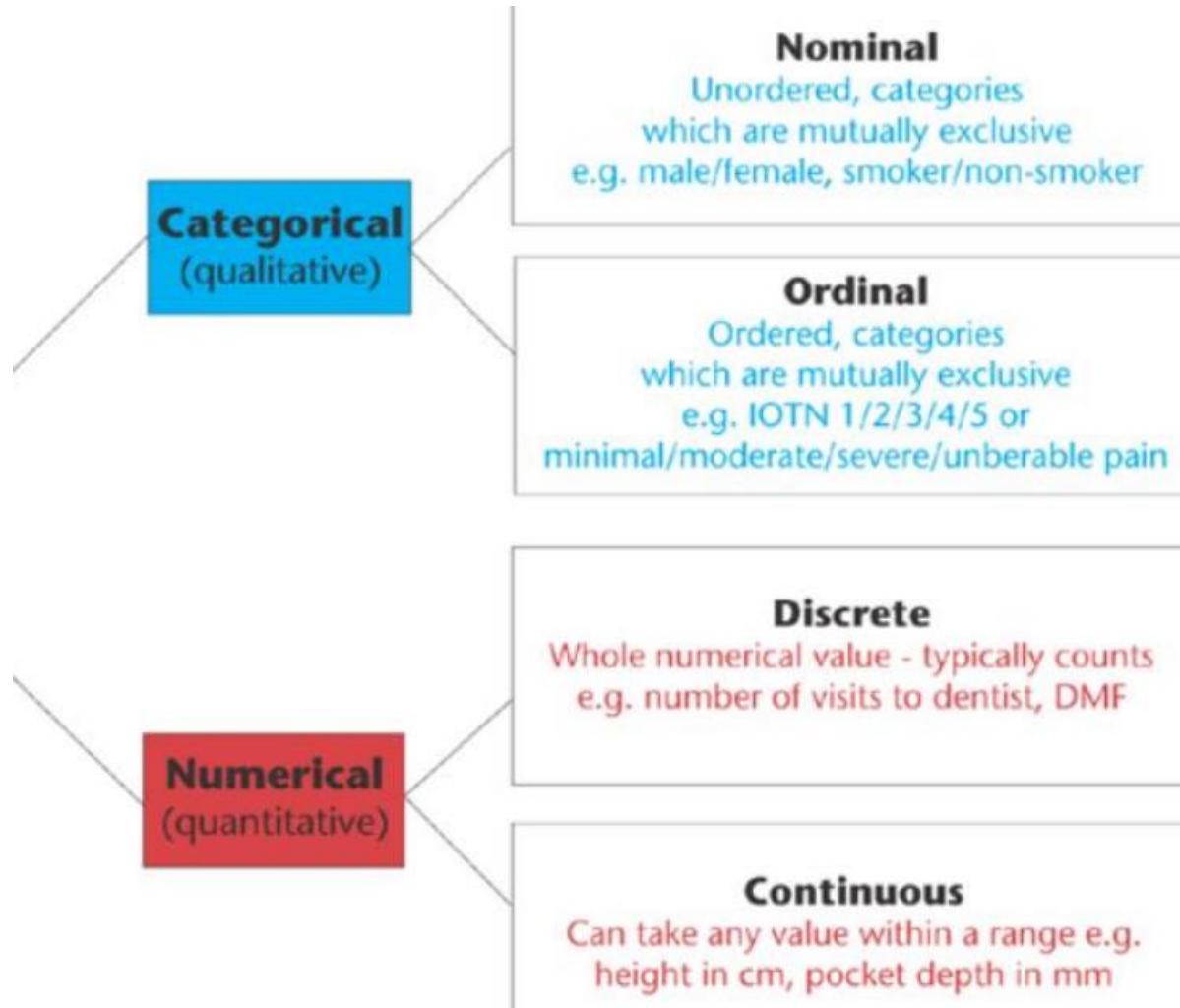
- Google Self-driving cars
- Tesla Model S – Autopilot mode

Gaming bots

- **Deep blue** : Chess playing computer developed by IBM. It is the masterpiece of a truly intelligent player. It **won against** Garry Kasparov(Chess grandmaster) in a match becoming the first computer system to defeat a reigning world champion in a match under standard chess tournament time controls.



Types of variables



examples of categorical information:

- colour = {red, green, blue, cyan, magenta, yellow, black}
- name = {Socrates, Plato}
- truth value = {true, false}

		type of output	
		quantitative	nominal
supervised	YES	REGRESSION	CLASSIFICATION
	NO	LOW-DIMENSIONAL MAPPING	CLUSTERING

A generic Python machine learning library

Scikit-learn

<https://scikit-learn.org/>

- A simple-to-use, yet fairly complete, machine learning library
- Built upon other standard libraries:
 - numpy (Matlab-like vectors, arrays, linear algebra)
 - scipy (computational methods)
 - matplotlib (Matlab-like plotting)

A python data analysis library

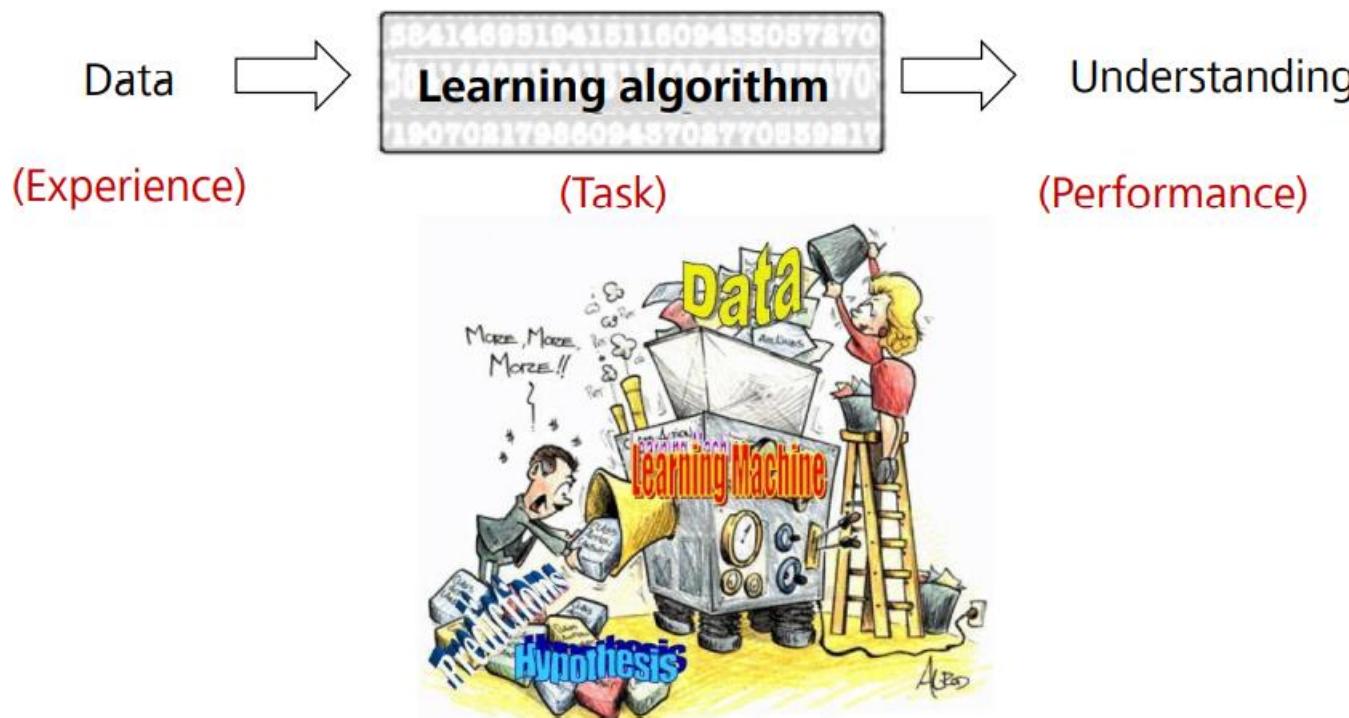
Pandas, a Python data analysis library

```
import pandas as pd
```

- Suitable for tabular data
- Series and DataFrame data structures
- Heterogeneous columns, missing data
- Intrinsic data alignment (data with their labels)
- Works with rows, columns, slices, boolean indexing
- Data types interoperable with numpy

Learning algorithms

- ◆ Machine learning (including deep learning) is a study of learning algorithms. A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .



Basic Terms and Concepts

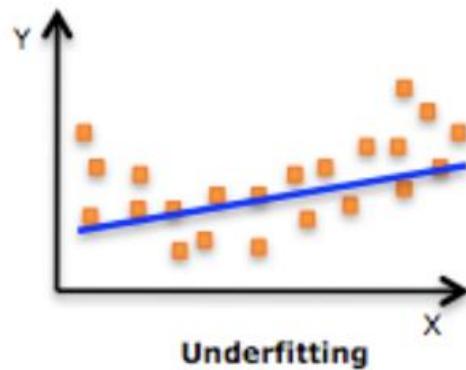
- ◆ **Dataset:** refers to a set of data used in machine learning tasks. Each piece of data is called a sample. The event or attribute that reflects the performance or nature of a sample in a certain aspect is called a **feature**.
- ◆ **Training set:** refers to a dataset used in the training process, where each sample is referred to as a training sample. The process of learning a model from data is called **learning (training)**.
- ◆ **Test set:** Test refers to the process of using the learnt model for prediction. The dataset used is called a test set, and each sample is called a test sample.

Basic Terms and Concepts

- ◆ **Generalization capability:** The goal of machine learning is that the learnt model should perform well on new samples, not just those on which the model has been trained. The ability to perform well on new samples is called generalization capability.
- ◆ **Error:** refers to the difference between the sample result predicted by the learnt model and the actual sample result.
 - Training error: error of the model on the training set
 - Generalization error: error on the new sample. Obviously, we prefer a model with a smaller generalization error.
- ◆ **Underfitting:** occurs when the training error is too large.
- ◆ **Overfitting:** occurs when the training error of the learnt model is small but the generalization error is large (weak generalization capability).

Basic Terms and Concepts

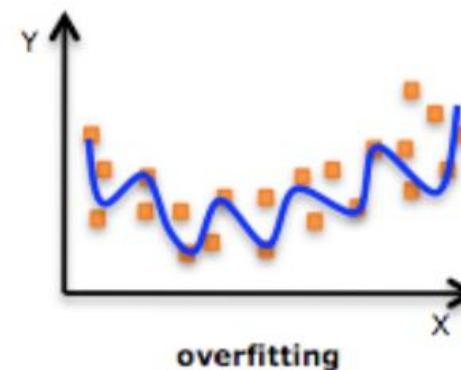
- ◆ **Capacity of a model:** refers to the ability to fit a wide variety of functions. Machine learning algorithms will generally perform best when their capacity is appropriate for the true complexity of the task they need to perform and the amount of training data they are provided with. Models with an insufficient capacity are unable to solve complex tasks. Models with a high capacity can solve complex tasks, but when their capacity is higher than that is needed to solve the present task, they may overfit.



Underfitting
Failing to learn features

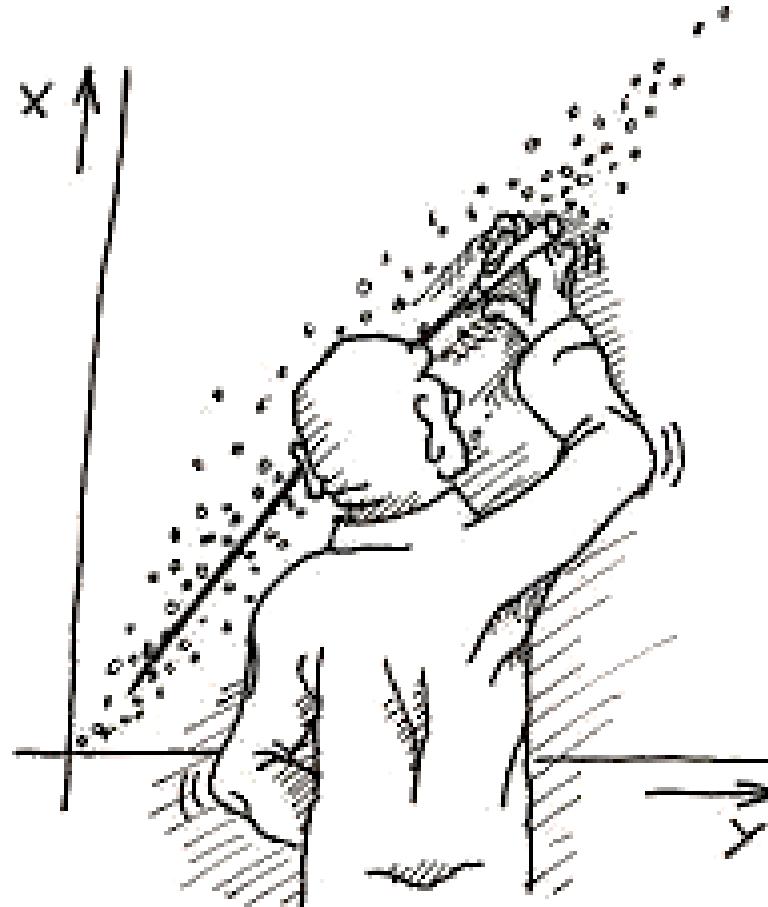


Just right!



overfitting
Learning noise

Régression



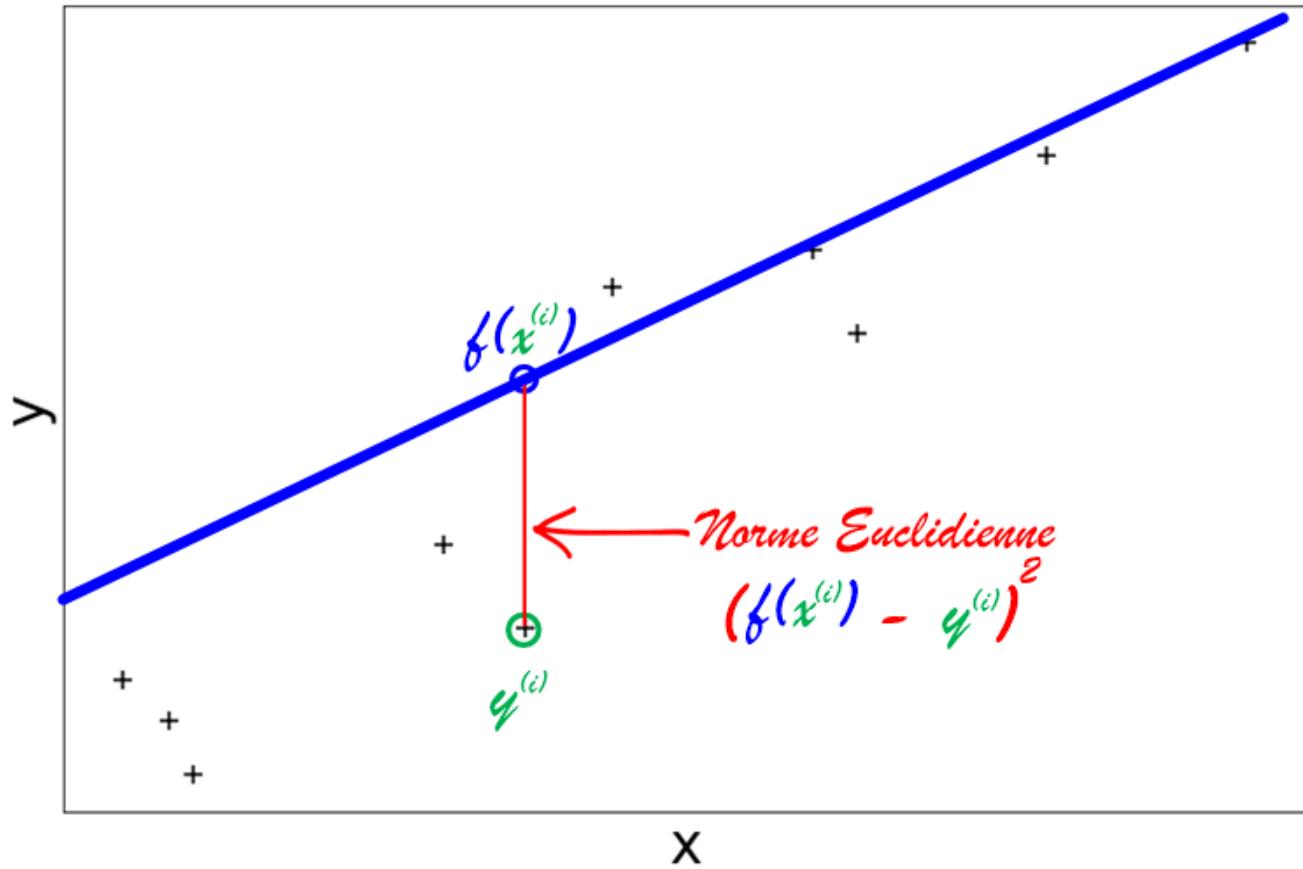
Régression

- La régression sert à trouver la relation d'**une** variable par rapport à **une** ou **plusieurs** autres. Le but est d'estimer une valeur (numérique) de sortie à partir des valeurs d'un ensemble de caractéristiques en entrée. Par exemple, estimer le prix d'une maison en se basant sur sa surface, nombre des étages, son emplacement, etc.

Donc, le problème revient à:

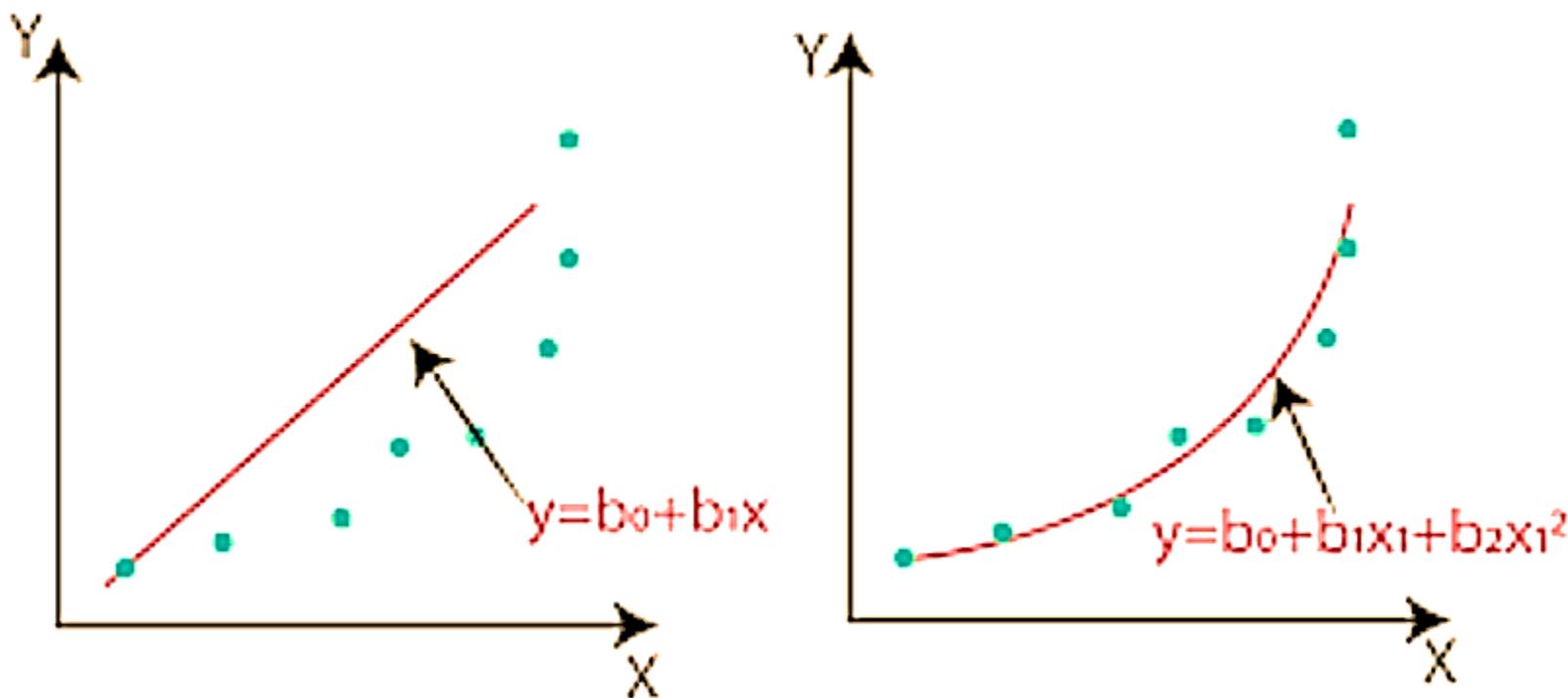
- **Estimer une fonction de calcul (fonction d'hypothèse)** en se basant sur des données d'entraînement.
- **Définir une fonction coût** : c'est une fonction mathématique qui mesure l'erreur que nous commettons en approximant les données. Nous parlons aussi d'erreur induite par la modélisation.

- **Minimiser cette fonction coût** : il faut trouver les bons paramètres de notre modèle pour minimiser l'erreur de modélisation.
- **Choisir une méthode de résolution du problème.** Il existe deux méthodes :
 - une méthode de résolution numérique, la descente de gradient,
 - une méthode analytique, la méthode des moindres carrés.



- Il existe plusieurs algorithmes pour la régression:
- Régression linéaire: Dans le cadre d'un modèle linéaire simple, on peut représenter graphiquement la relation entre x et y à travers un nuage de points. L'estimation du modèle linéaire permet de tracer la droite de régression, d'équation: $y=w_1X+w_0$. Le paramètre w_0 représente l'ordonnée à l'origine et w_1 le coefficient directeur de la droite.
- Régression polynomiale: La régression polynomiale est une approche statistique qui est employée pour modéliser une relation de forme **non-linéaire** entre la réponse (y) et les variables explicatives (x). Pour prendre en charge cette forme non-linéaire de la relation entre y et x , ces modèles de régression intègrent des polynômes dans leurs équations : $y=w_2X^2+w_1X+w_0$

Regression lineaire et polynomiale

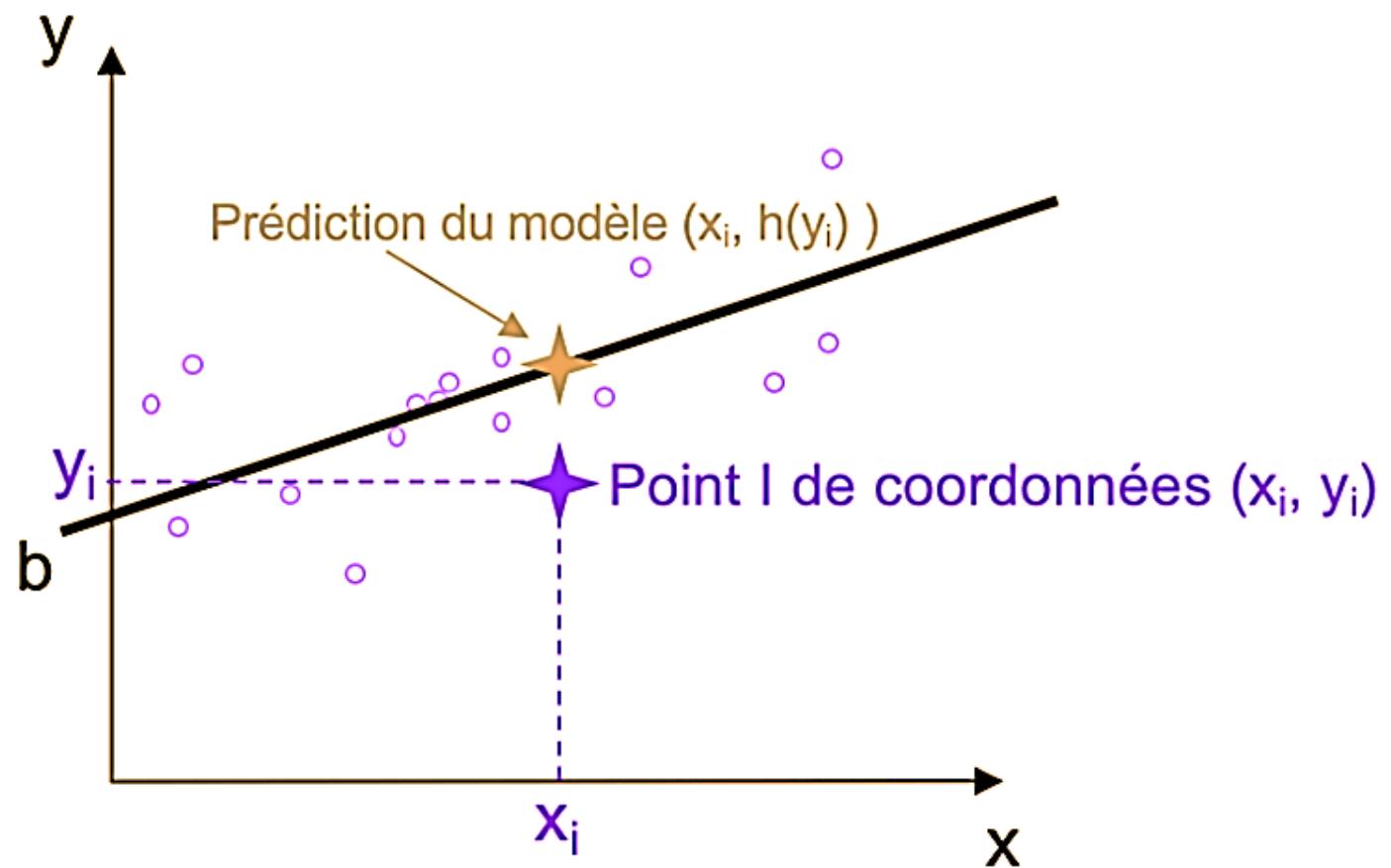


Régression linéaire

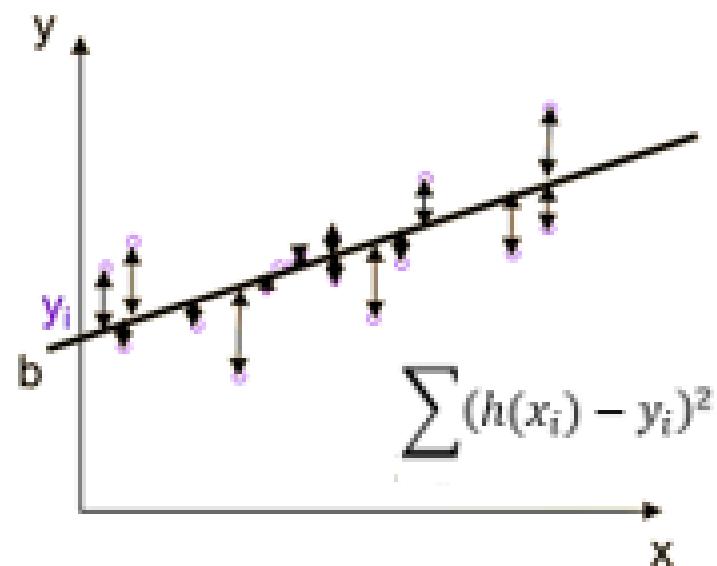
Pour chaque point x_i , la fonction hypothèse associe une valeur définie par $h(x_i)$ qui est plus ou moins proche de la variable cible y_i .

l'erreur unitaire pour x_i est *définie par* $h(x_i) - y_i$.

Or, chaque erreur unitaire pouvant être positive ou négative.. Il faut donc faire en sorte que la contribution de chaque erreur soit systématiquement pénalisante. Alors, on élève au carré l'erreur unitaire



).



On somme l'ensemble des erreurs unitaires pour l'ensemble des points de données : on parle alors d'erreur quadratique (erreur au carré) :

$$\sum (h(x_i) - y_i)^2$$

La fonction de coût s'écrit alors en pondérant la somme des erreurs quadratiques par le nombre de points dans la base d'apprentissage :

$$\frac{1}{n} \sum (h(x_i) - y_i)^2$$

Régression linéaire univariée

Dans notre cas de régression linéaire univariée, la fonction de coût devient :

$$\frac{1}{n} \sum (w_0 + w_1 x - y_i)^2$$

Déterminer les meilleurs paramètres (w_0 , w_1) pour la fonction hypothèse h revient à trouver la meilleure **droite**, celle qui minimise la somme de toutes les erreurs unitaires.

Du point de vue arithmétique, il s'agit de trouver le minimum de la fonction de coût. Pour cela il existe plusieurs méthodes:

- **Méthode Analytique par OLS** (ordinary least square): l'idée est de réduire le carré des erreurs et estimer les coefficients b_0 et b_1 par des relations mathématiques fermées. (Il n'y a pas besoin de faire d'algorithmes compliqués qui sont itératifs)
- **Méthode numérique par gradient descent** (gradient descendant). C'est la plus utilisée en machine learning pour résoudre les problèmes. C'est un algorithme qui procède par itération pour déterminer les valeurs des coefficients.

Linear regression

Regression means:

approximating a functional dependency based on measured data.

“learning a function” or “predicting continuous values”

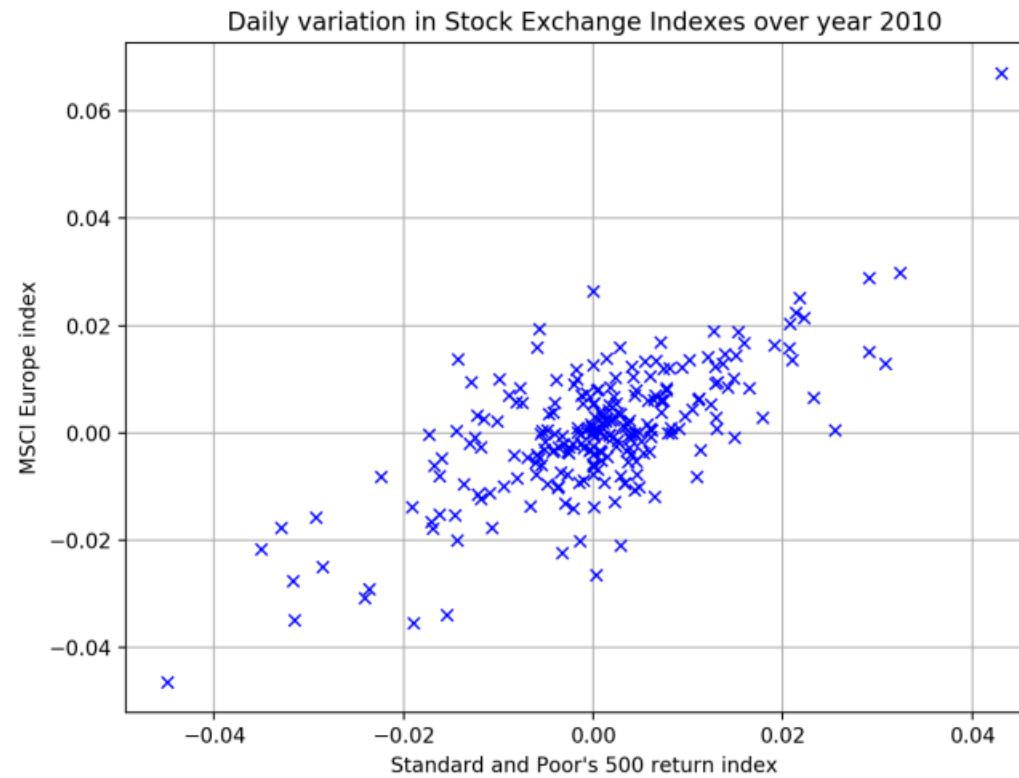
a typical supervised problem

Data

$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_N \end{pmatrix} \quad \begin{pmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \\ \mathbf{t}_3 \\ \vdots \\ \mathbf{t}_N \end{pmatrix}$$

Observations Target

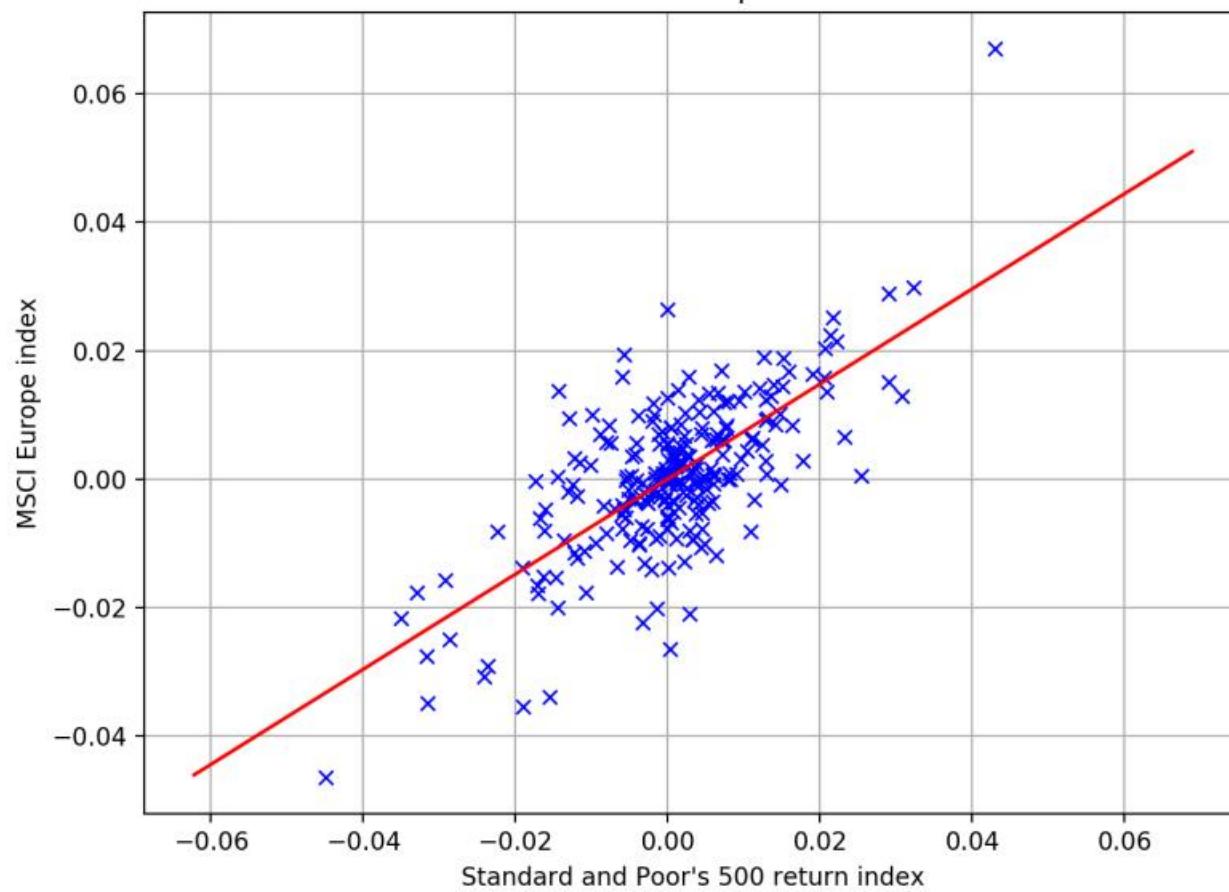
Exemple 1 One dimensional Linear regression



249 observations in year 2010. (Source: UCI)

We want to predict the variation of the MSCI European index by observing Standard and Poor 500 return index. quad

Indexes: The least squares solution



- Observation: x is the value of the variation of Standard and Poor's (SP) 500 return index on a given day.
- Target: t is the value of the variation of the MSCI European index (MSCI) on the same day.

There is clearly some relationship between the two values.

But not one-to-one

A model for approximating the data

A linear model $y(x)$ that predicts t given x :

$$t \approx y \quad \text{where} \quad y = wx$$

For instance

$$t_1 \approx y_1 \quad \text{where} \quad y_1 = wx_1$$

$$\text{or } t_2 \approx y_2 \quad \text{where} \quad y_2 = wx_2$$

We want $y(x)$ to be similar to $t(x)$ for any x .

A solution based on optimization

Idea:

- quantify **how wrong** is each estimate using some **measure**,
- make this measure as small as possible **on average**.

Measure = **loss function**

Overall evaluation (e.g. average) = **objective function**

Properties of the square error loss

- is even: $(t - y)^2 = (y - t)^2$
- grows more than linearly, giving heavier weight to larger errors
- is differentiable with respect to the model output:

$$\frac{d}{dy} \lambda_{\text{SE}}(y, t) = 2(y - t)$$

Square error loss

$$\lambda_{\text{SE}}(y, t) = (y - t)^2$$

The objective function

Generic goal: **minimize the mean value of the loss over the whole data set**

$$J = \frac{1}{N} \sum_{l=1}^N \lambda(y_l, t_l).$$

J = Objective function or cost function

In the specific case of the square error loss:

$$J_{\text{MSE}} = \frac{1}{N} \sum_{l=1}^N (y_l - t_l)^2.$$

J_{MSE} = mean square error objective.

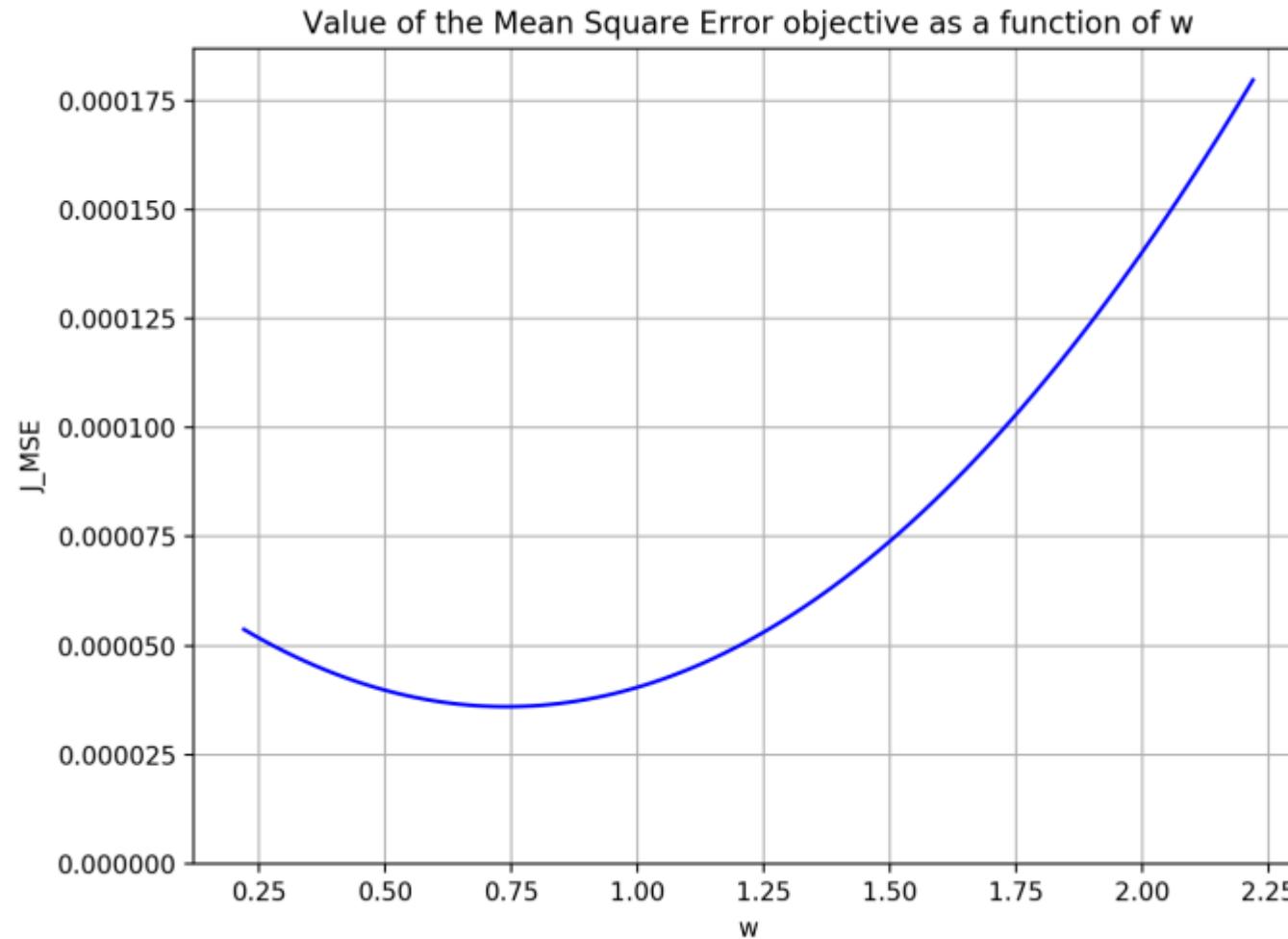
The objective function

A loss function $\lambda(y, t)$ is a function of the two arguments y and t .

We can look at the objective in two complementary situations:

- Given the data set, the targets t_1, \dots, t_N are fixed
→ **The objective depends only on the model parameter(s)**
- Given a model, the parameters are fixed.
However we can apply this model to various data sets.
→ **The objective depends only on the data.**

The objective as a function of the model parameter(s)



- When **building a model (= training)**, the objective is a function of the parameters in the model and the data are fixed
- When **using a model (= inference)**, the objective is a function of the data, which uses the (now fixed) model parameters.

The least squares method

Problem:

minimize J_{MSE} with respect to the parameters with fixed data

we know that J_{MSE} is a parabola

Unique solution: w such that

$$\frac{d}{dw} J_{\text{MSE}} = 0$$

Computing $\frac{d}{dw}J_{\text{MSE}}$

$$\frac{d}{dw}\lambda_{\text{SE}}(y, t) = \frac{d}{dy}\lambda_{\text{SE}}(y, t) \times \frac{d}{dw}y$$

Here we have used the chain rule of differentiation to write the derivative:

$$\frac{df(g(x))}{dx} = \frac{df(y)}{dy} \Big|_{y=g(x)} \quad \frac{dg(x)}{dx}$$

$$\downarrow \qquad \qquad \downarrow$$
$$\frac{d}{dy}(y - t)^2 \Big|_{y=wx} \quad \frac{d}{dw}(wx)$$

$$\downarrow \qquad \qquad \downarrow$$
$$2(xw - t) \qquad \quad x$$

Derivative of the loss

This is for one observation

$$\begin{aligned}\frac{d}{dw} \lambda_{\text{SE}}(t, y) &= 2(xw - t)(x) \\ &= 2x^2 w - 2xt.\end{aligned}$$

Derivative of the objective

This is for the average over all observations

$$\frac{d}{dw} J_{\text{MSE}} = \frac{d}{dw} \frac{1}{N} \sum_{l=1}^N \lambda_{\text{SE}}(y_l, t_l)$$

Exchange sum and derivative:

$$\begin{aligned} &= \frac{1}{N} \sum_{l=1}^N \frac{d}{dw} \lambda_{\text{SE}}(y_l, t_l) \\ &= \frac{1}{N} \sum_{l=1}^N 2(x_l^2 w - x_l t_l), \end{aligned}$$

where the constant coefficient ($2/N$) is irrelevant and can be disregarded.

This equation is solved by bringing w outside the sum, since it does not depend on l :

$$w \sum_{l=1}^N x_l^2 - \sum_{l=1}^N x_l t_l = 0$$

$$w = \frac{\sum_{l=1}^N x_l t_l}{\sum_{l=1}^N x_l^2}$$

This is the least squares solution to the linear regression problem.

Metrics to evaluate the regression

Regression

'explained_variance'	<code>metrics.explained_variance_score</code>
'max_error'	<code>metrics.max_error</code>
'neg_mean_absolute_error'	<code>metrics.mean_absolute_error</code>
'neg_mean_squared_error'	<code>metrics.mean_squared_error</code>
'neg_root_mean_squared_error'	<code>metrics.mean_squared_error</code>
'neg_mean_squared_log_error'	<code>metrics.mean_squared_log_error</code>
'neg_median_absolute_error'	<code>metrics.median_absolute_error</code>
'r2'	<code>metrics.r2_score</code>
'neg_mean_poisson_deviance'	<code>metrics.mean_poisson_deviance</code>
'neg_mean_gamma_deviance'	<code>metrics.mean_gamma_deviance</code>
'neg_mean_absolute_percentage_error'	<code>metrics.mean_absolute_percentage_error</code>
'd2_absolute_error_score'	<code>metrics.d2_absolute_error_score</code>
'd2_pinball_score'	<code>metrics.d2_pinball_score</code>
'd2_tweedie_score'	<code>metrics.d2_tweedie_score</code>

R2-score

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Exemple 2

1974 “Motor Trends” car data (four columns)

Model	mpg	disp	hp	weight
Mazda RX4	21	160	110	2.62
Mazda RX4 Wag	21	160	110	2.875
Datsun 710	22.8	108	93	2.32
Hornet 4 Drive	21.4	258	110	3.215
Hornet Sportabout	18.7	360	175	3.44
Valiant	18.1	225	105	3.46
Duster 360	14.3	360	245	3.57
Merc 240D	24.4	146.7	62	3.19
Merc 230	22.8	140.8	95	3.15
Merc 280	19.2	167.6	123	3.44
Merc 280C	17.8	167.6	123	3.44
Merc 450SE	16.4	275.8	180	4.07
Merc 450SL	17.3	275.8	180	3.73
Merc 450SLC	15.2	275.8	180	3.78
Cadillac Fleetwood	10.4	472	205	5.25
Lincoln Continental	10.4	460	215	5.424
Chrysler Imperial	14.7	440	230	5.345
Fiat 128	32.4	78.7	66	2.2
Honda Civic	30.4	75.7	52	1.615
Toyota Corolla	33.9	71.1	65	1.835
Toyota Corona	21.5	120.1	97	2.465
Dodge Challenger	15.5	318	150	3.52
AMC Javelin	15.2	304	150	3.435
Camaro Z28	13.3	350	245	3.84
Pontiac Firebird	19.2	400	175	3.845
Fiat X1-9	27.3	79	66	1.935
Porsche 914-2	26	120.3	91	2.14
Lotus Europa	30.4	95.1	113	1.513
Ford Pantera L	15.8	351	264	3.17
Ferrari Dino	19.7	145	175	2.77
Maserati Bora	15	301	335	3.57
Volvo 142E	21.4	121	109	2.78

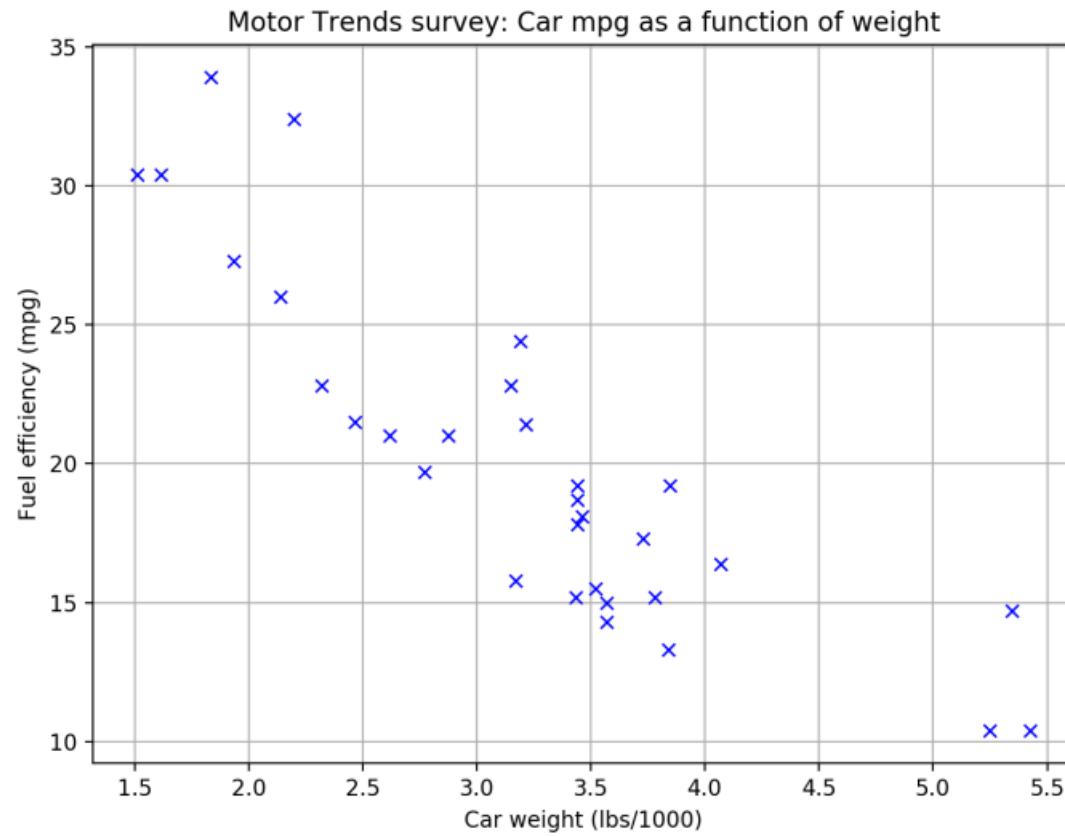
1974 survey on some car models.

Among other variables:

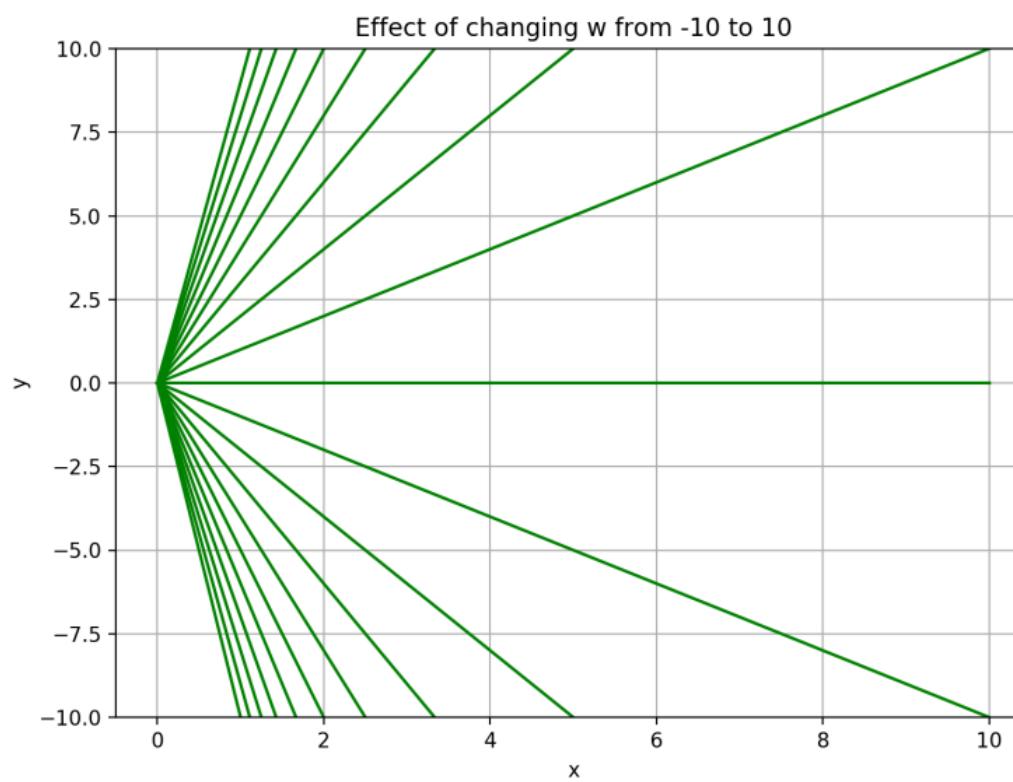
- mpg or miles-per-gallon
- disp or displacement (in cu.in)
- hp or horse-power
- weight, total (in lbs/1000)

(Note that these are USA units.)

Forecasting mpg with weight

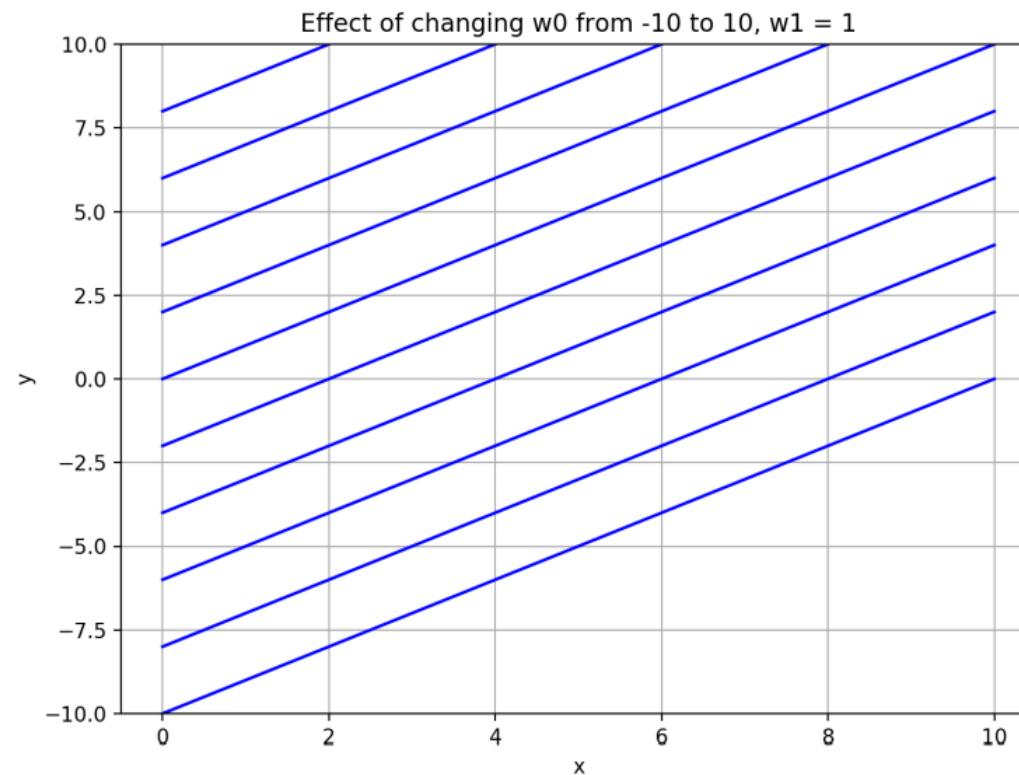


$$y = wx$$



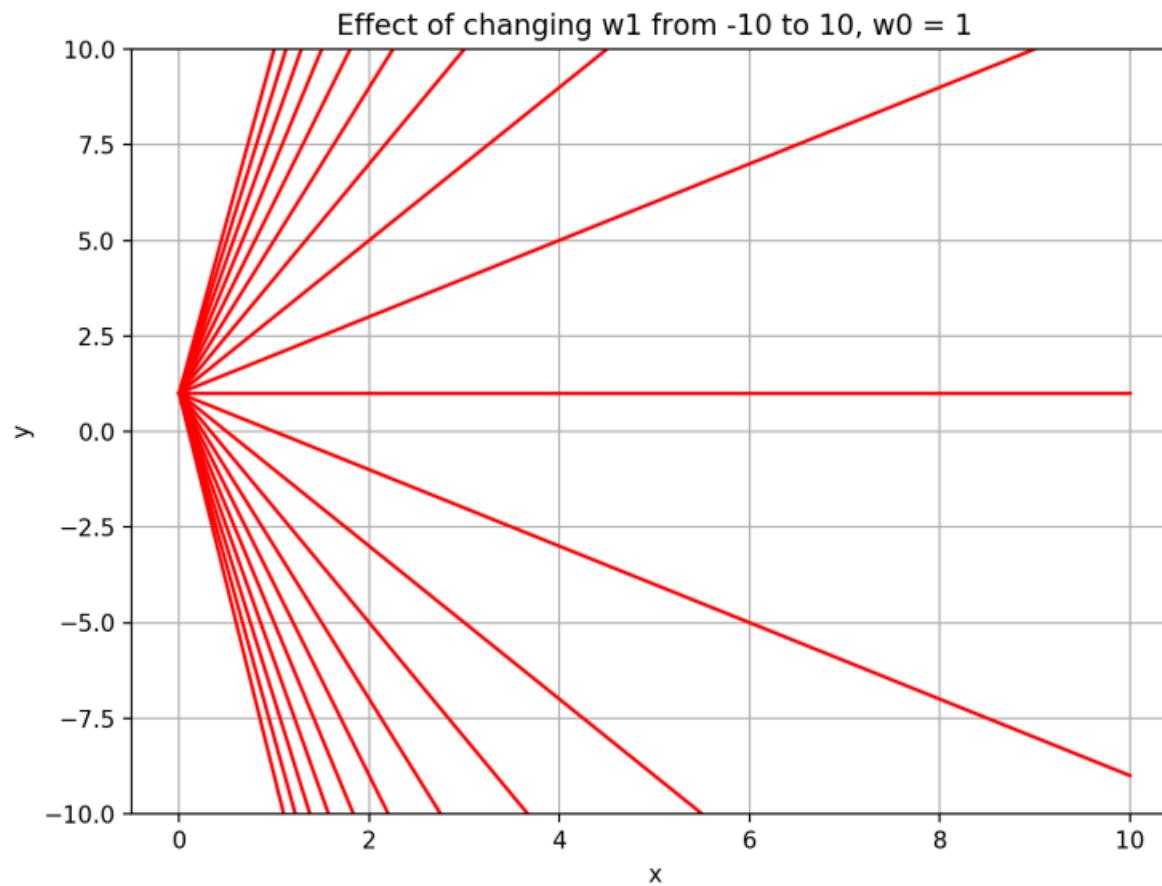
A more flexible model

$$y = w_1 x + w_0$$



A more flexible model

$$y = w_1 x + w_0$$



Solving

The solution in this case can be found by **centering** around the mean \bar{x} of x and \bar{t} of t

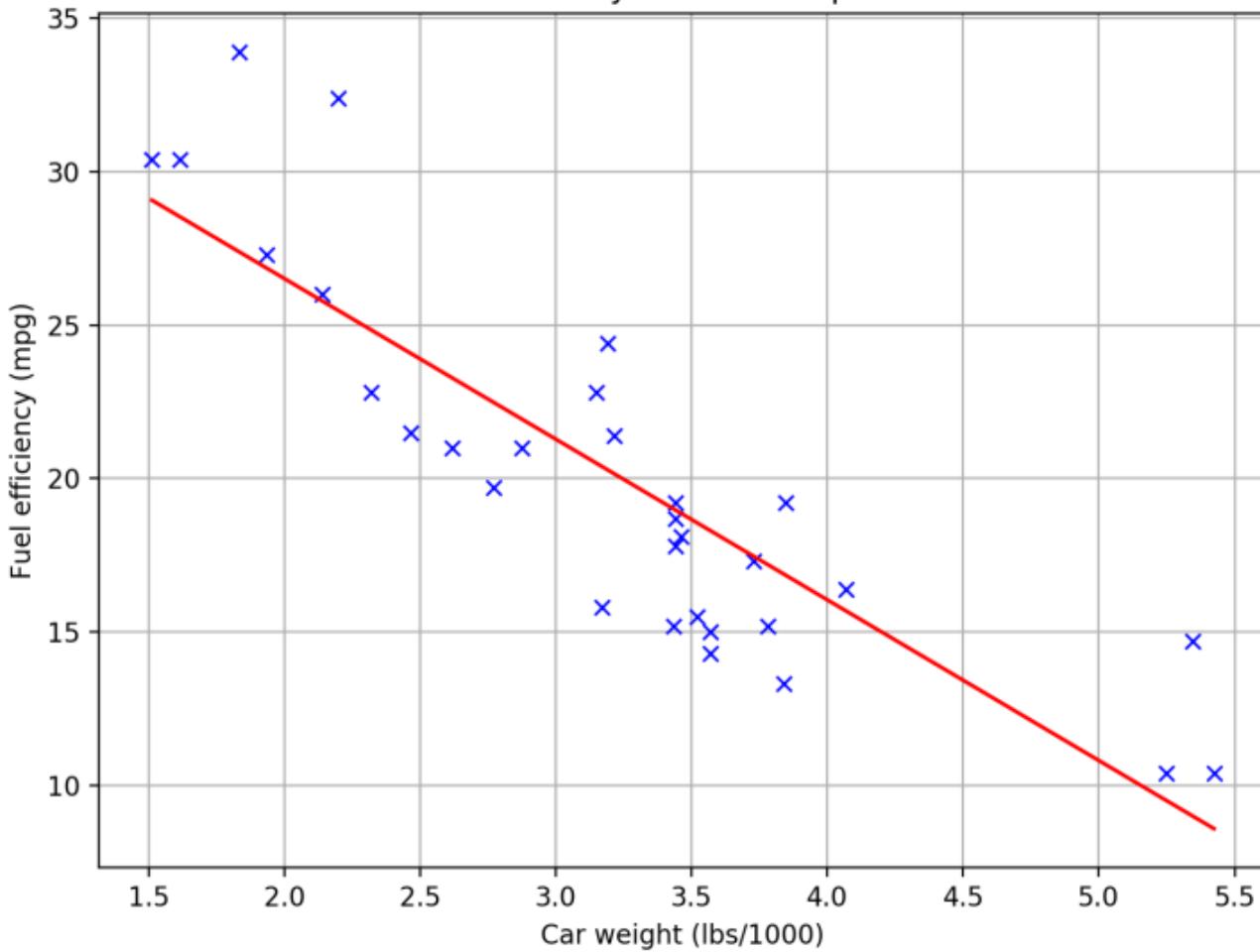
$$\bar{x} = \frac{1}{N} \sum_{l=1}^N x_l \quad \bar{t} = \frac{1}{N} \sum_{l=1}^N t_l$$

$$w_1 = \frac{\sum_{l=1}^N (x_l - \bar{x})(t_l - \bar{t})}{\sum_{l=1}^N (x_l - \bar{x})^2}$$

$$w_0 = \bar{t} - w_1 \bar{x}$$

- w_1 = slope; gain
 w_0 = intercept, offset; bias

Motor Trends survey: The least squares solution



Multidimensional linear regression

The data is now composed of d -dimensional vectors:

$$\mathbf{x}_1 = [x_{1,1}, x_{1,2}, \dots, x_{1,d}]$$

$$\mathbf{x}_2 = [x_{2,1}, x_{2,2}, \dots, x_{2,d}]$$

...

$$\mathbf{x}_N = [x_{N,1}, x_{N,2}, \dots, x_{N,d}]$$

so we can organize them into a $N \times d$ matrix:

$$X = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ x_{3,1} & x_{3,2} & \dots & x_{3,d} \\ \vdots & & & \\ x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix}$$

Since the data are now d -dimensional, we have d parameters in a d -dimensional vector:

$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix}$$

This can be expressed as a matrix-vector multiplication between the data matrix X and the parameter vector \mathbf{w} :

$$\begin{aligned}\mathbf{y} &= \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ x_{3,1} & x_{3,2} & \dots & x_{3,d} \\ & & \vdots & \\ x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} \\ &= \begin{pmatrix} x_{1,1}w_1 + x_{1,2}w_2 + & \dots & +x_{1,d}w_d \\ x_{2,1}w_1 + x_{2,2}w_2 + & \dots & +x_{2,d}w_d \\ x_{3,1}w_1 + x_{3,2}w_2 + & \dots & +x_{3,d}w_d \\ & & \vdots \\ x_{N,1}w_1 + x_{N,2}w_2 + & \dots & +x_{N,d}w_d \end{pmatrix} = X\mathbf{w}\end{aligned}$$

Finally, our goal is to make this model's prediction \mathbf{y} as similar as possible to the measured outputs for each observation, which again can be organized as a vector, this time N -dimensional:

$$\mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_N \end{pmatrix}$$

The square error objective in matrix-vector form

$$\begin{aligned} J_{\text{MSE}} &= \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|^2 \\ &= \frac{1}{2} \|X\mathbf{w} - \mathbf{t}\|^2 \\ &= \frac{1}{2} (X\mathbf{w} - \mathbf{t})^T (X\mathbf{w} - \mathbf{t}) \\ &= \frac{1}{2} (\mathbf{w}^T X^T - \mathbf{t}^T)(X\mathbf{w} - \mathbf{t}) \\ &= \frac{1}{2} (\mathbf{w}^T X^T X \mathbf{w} - \mathbf{w}^T X^T \mathbf{t} - \mathbf{t}^T X \mathbf{w} + \mathbf{t}^T \mathbf{t}) \\ &= \frac{1}{2} \|X\mathbf{w}\|^2 - \mathbf{w}^T X^T \mathbf{t} - \frac{1}{2} \|\mathbf{t}\|^2 \end{aligned}$$

Closed-form solution

It can be proven that we can write:

$$\begin{aligned}\nabla J_{\text{MSE}} &= \frac{\partial}{\partial \mathbf{w}} J_{\text{MSE}} \\ &= \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t}\end{aligned}$$

Setting $\nabla J_{\text{MSE}} = \mathbf{0}$ we get

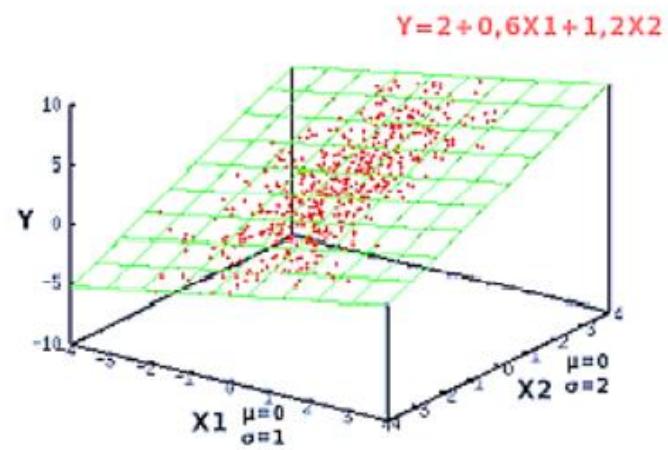
$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{t}$$

By premultiplying both sides by $(\mathbf{X}^T \mathbf{X})^{-1}$, we obtain the closed-form solution:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

The matrix form of the **normal equations** for the least squares problem.

Example



Problem!!

Invertibility and stability

- $X^T X$ might not be invertible. This is when the data have exactly linearly dependent components
- What if the variables are **correlated**?
- Even if $X^T X$ has full rank, in the case of correlated variables it will have a high **condition number** $\lambda_{\text{first}} / \lambda_{\text{last}}$.
- DIFFICULT TO INVERT (numerical precision must be too high)
- This is a problem of **numeric instability**: even very small numeric errors are amplified by the condition number and cause large errors on the result

SOLUTION: Iterative computation by successive approximations

Problem!!



To solve the problem of overfitting in linear regression, the technique of **regularization** is used in such cases.

Polynomial Regression

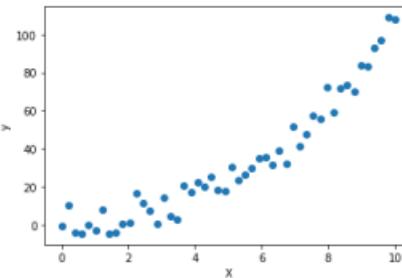
$$y = h(x) = b_0 + b_1x + b_2x^2 + \dots + b_nx^n$$

Python

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import mean_squared_error, r2_score
```

Maintenant, en utilisant les lignes de code suivantes, nous pouvons générer des nuages de points ayant un rapport linéaire :

```
# Étape 1: génération de données
m = 50      #création de 50 échantillons
X = np.linspace(0,10, m).reshape(m, 1) # Vecteur X
y = X**2 + 6*np.random.randn(m, 1)      # Y = X2 + la déviation par
#rapport à x
plt.ylabel('y')
plt.xlabel('X')
plt.scatter(X,y)
```



```
# Étape 2: la préparation des données
nb_degree = 2    # degré de polynome (1 -> y=b1X+b0  2-> y=b2X^2+b1X+b0)
polynomial_features = PolynomialFeatures(degree = nb_degree)
X_TRANSF = polynomial_features.fit_transform(X)

# Étape 3: définir et entraîner le modèle
model = LinearRegression() #L'utilisation de la régression linéaire sur les #données transformées
nous permet d'estimer la fonction de régression des #polynômes.

model.fit(X_TRANSF, y)  # Entraîner le modèle sur les données X,y

# Étape 4: calculer le biais et la variance
Y_NEW = model.predict(X_TRANSF)
rmse = np.sqrt(mean_squared_error(y,Y_NEW))
r2 = r2_score(y,Y_NEW)

print('RMSE: ', rmse)
print('R2: ', r2)
```

```

model.score(X_TRANSF)      # Evaluer le modèle sur les mêmes données
predictions = model.predict(X_TRANSF)
plt.scatter(X,y)
title = 'Degree = {} ; RMSE = {} ; R2 = {}'.format(nb_degree,
round(rmse,2), round(r2,2))

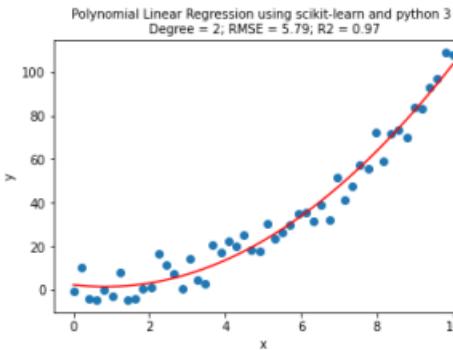
plt.plot(X,predictions,c='r')
plt.title("Polynomial Linear Regression using scikit-learn and python 3
\n " + title,
          fontsize=10)
plt.xlabel('x')
plt.ylabel('y')

```

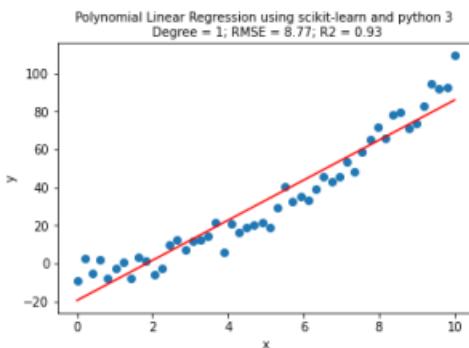
Sortie: 0.965828030085421

l'exécution de ses lignes de code, on obtient un score de 96% qui est le coefficient de détermination R² de moindre carré.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$



Maintenant si on change nb_degree pour quelle soit =1, notre modèle n'est autre que la régression linéaire univariée. Et le score sera dans ce cas de 82%.



Logistic Regression

Example 1: Malignant or Benign

- Consider the example of recognizing whether a tumor in an input image is malignant or benign

Input:



Output:

Benign

Malignant

0

1

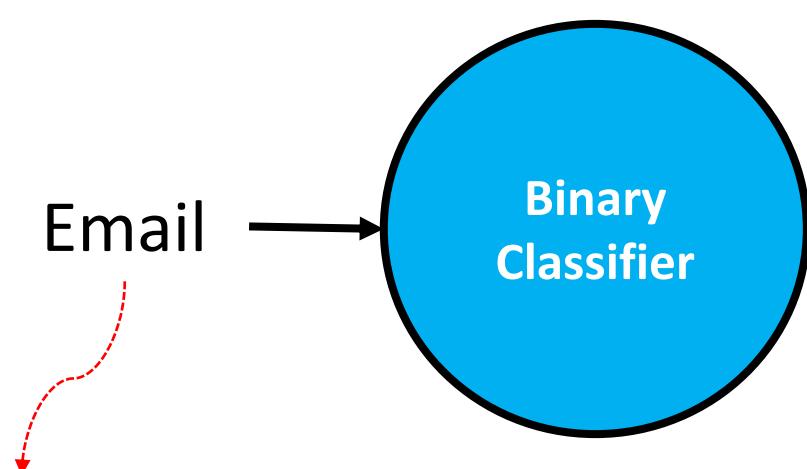
Can be represented as integers!

Can be represented as a **matrix** of pixels

Example 2: Spam or Not Spam

- As another example, consider the problem of detecting whether an email is a spam or not a spam

Input:



Output:

Not Spam

0

Spam

1

Can be represented as integers!

Can be represented as a vector $\mathbf{x} = [x_1, x_2, \dots, x_d]$, with each component x_i corresponding to the presence ($x_i = 1$) or absence ($x_i = 0$) of a particular word (or *feature*) in the email

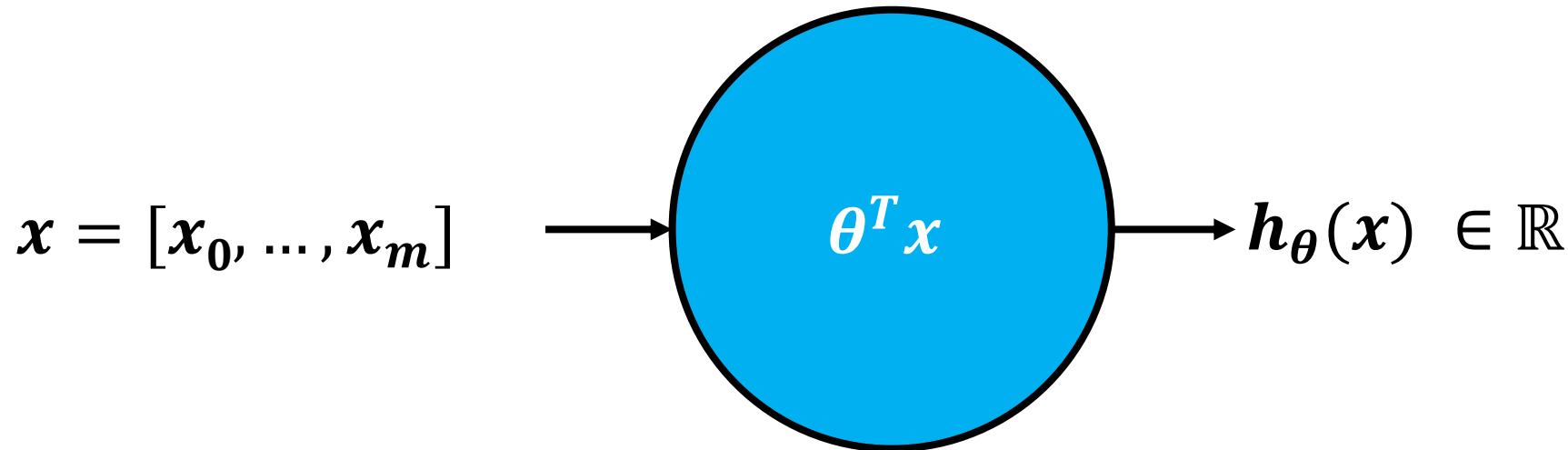
Regression vs. Classification

- What are the possible output values of the *linear regression model*
$$h_{\theta}(x) = \theta^T x$$
?

θ is the *parameter vector*, that is, $\theta = [\theta_0, \theta_1, \dots, \theta_m]$ (assuming $m + 1$ parameters) and x is the *feature vector*, that is, $x = [x_0, x_1, \dots, x_m]$ (assuming $m + 1$ features and $x_0 = 1$ to account for the intercept term, namely, θ_0)

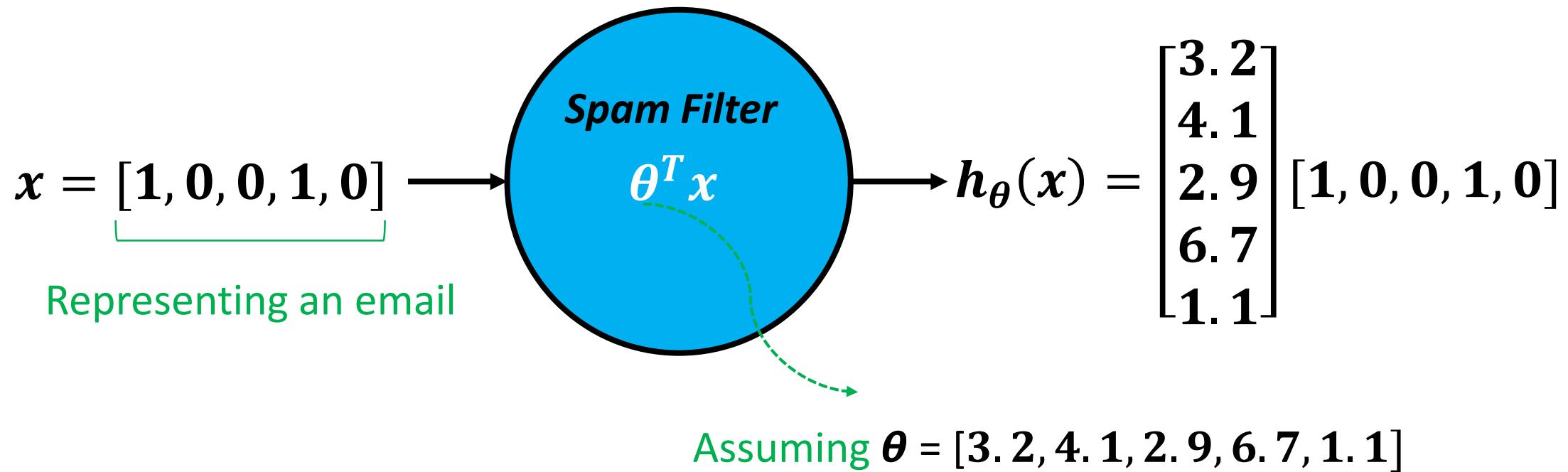
Regression vs. Classification

- What are the possible output values of the *linear regression model* $h_{\theta}(x) = \theta^T x$?
 - Real-valued outputs



Regression vs. Classification

- What are the possible output values of the *linear regression model* $h_{\theta}(x) = \theta^T x$?
 - Real-valued outputs

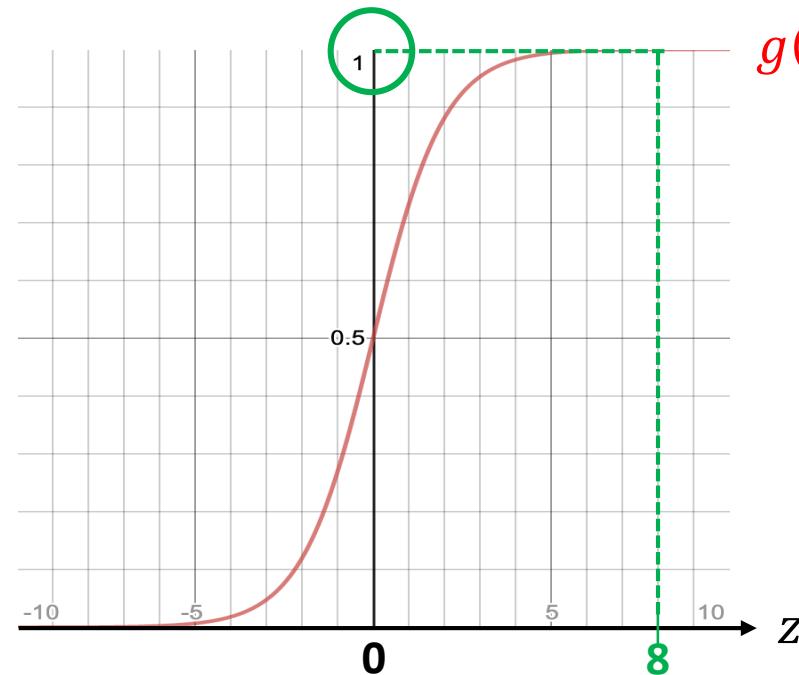


Regression vs. Classification

- How can we make the possible outputs of $h_{\theta}(x) = \theta^T x$ discrete-valued (as opposed to real-valued)?
 - By using an ***activation function*** (e.g., ***sigmoid or logistic function***)

$$g(z) = \frac{1}{1 + e^{-z}}$$

$z \in \mathbb{R}$, but
 $g(z) \in [0,1]$



$g(z)$ Assume a labeled example (x, y) :

If $y = 1$, we want $g(z) \approx 1$ (i.e., we want a correct prediction)

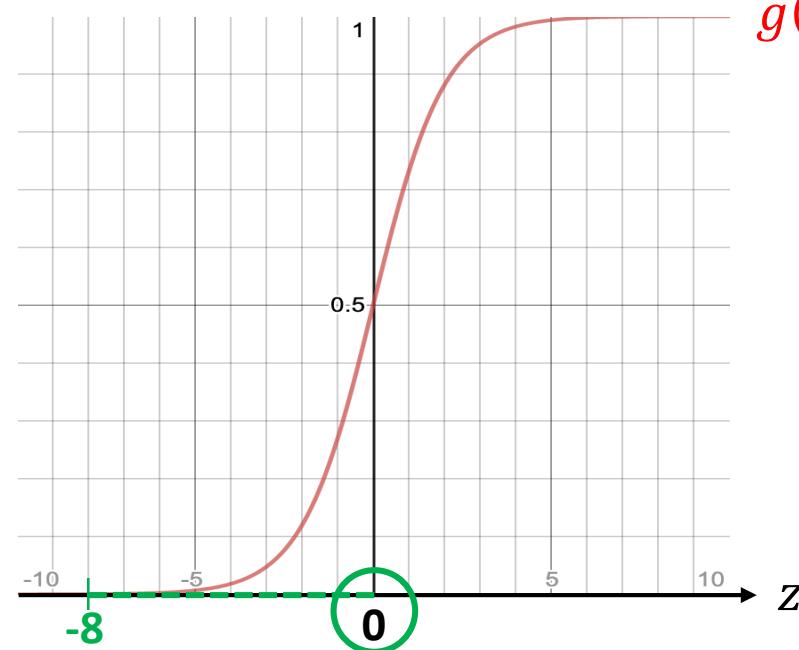
For this to happen, $z \gg 0$

Regression vs. Classification

- How can we make the possible outputs of $h_{\theta}(x) = \theta^T x$ discrete-valued (as opposed to real-valued)?
 - By using an ***activation function*** (e.g., ***sigmoid or logistic function***)

$$g(z) = \frac{1}{1 + e^{-z}}$$

$z \in \mathbb{R}$, but
 $g(z) \in [0,1]$



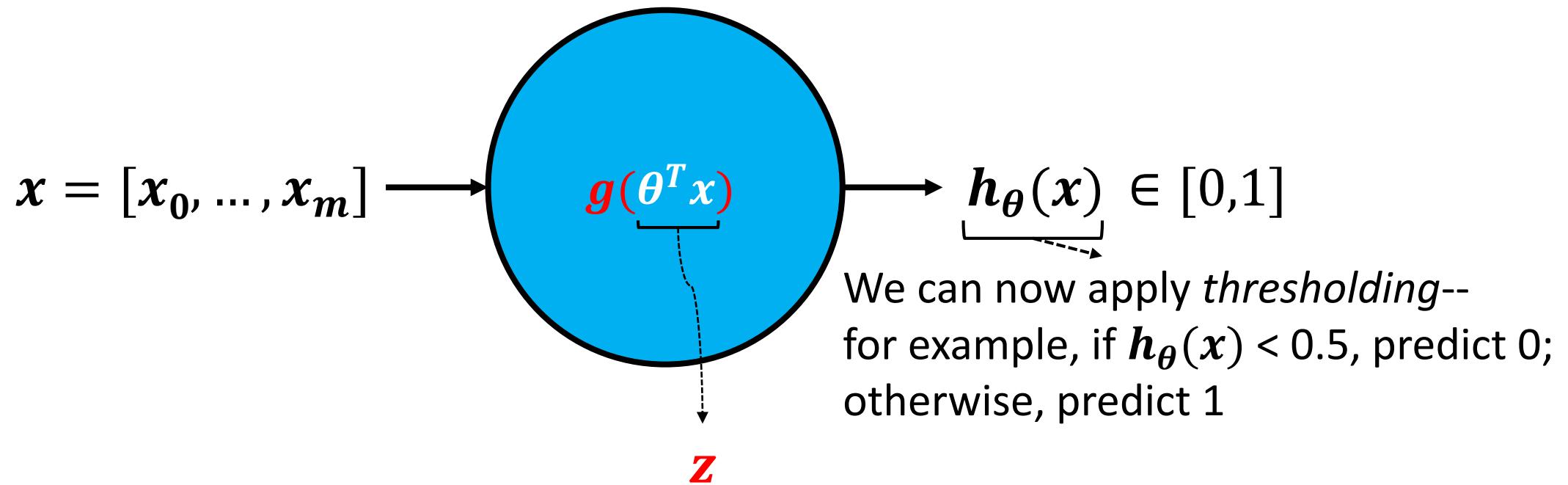
$g(z)$ Assume a labeled example (x, y) :

If $y = 0$, we want $g(z) \approx 0$ (i.e., we want a correct prediction)

For this to happen, $z \ll 0$

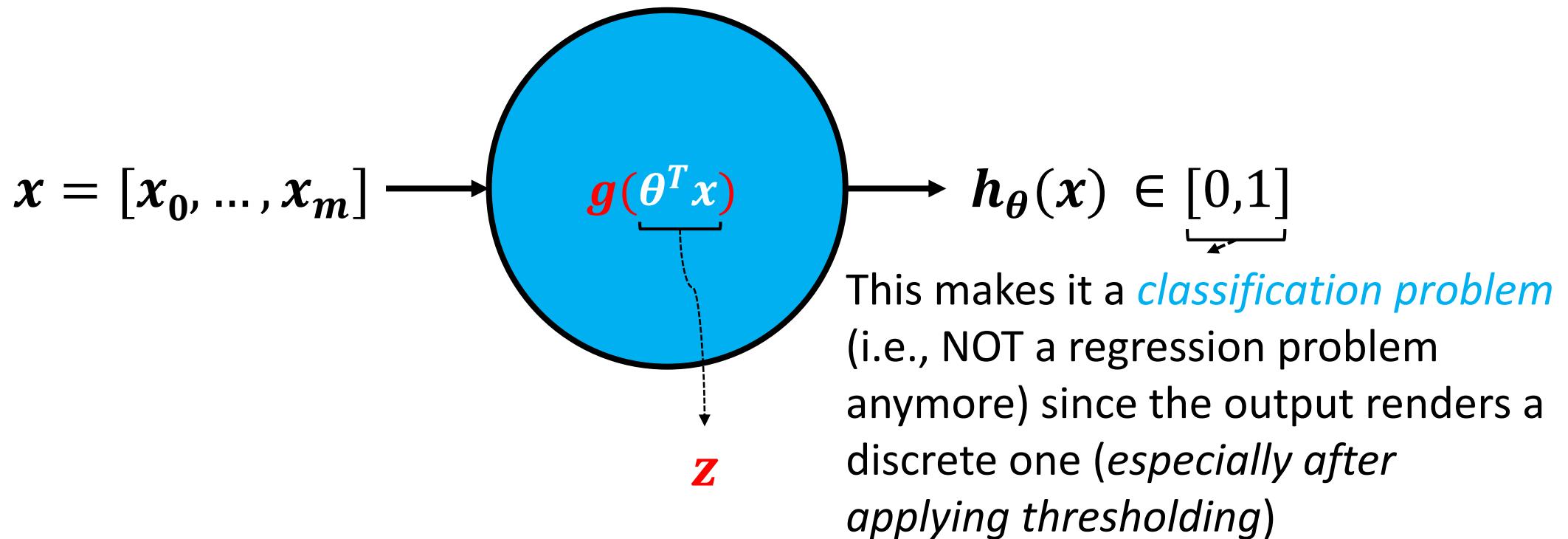
Regression vs. Classification

- How can we make the possible outputs of $\mathbf{h}_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$ discrete-valued (as opposed to real-valued)?
 - By using an *activation function* (e.g., *sigmoid or logistic function*)



Regression vs. Classification

- How can we make the possible outputs of $h_{\theta}(x) = \theta^T x$ discrete-valued (as opposed to real-valued)?
 - By using an *activation function* (e.g., *sigmoid or logistic function*)



So, What is Logistic Regression?

- Logistic regression is a machine learning algorithm that can be used to *classify* input data into discrete output (e.g., *input emails into spam or non-spam and tumour images into benign or malignant*)
 - **Note:** The word “regression” in the name does not mean that the algorithm is a regression algorithm (rather it is a classification algorithm)
- Major questions about logistic regression:
 - What is the ***hypothesis function*** (or ***model***)?
 - What is the ***cost function***?
 - How can we learn the parameters of the model?

The Logistic Regression Model

- What will be the output of the model $h_{\theta}(x) = \theta^T x$, where $\theta = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - Real-valued
- How can we make the output of $h_{\theta}(x)$ discrete?
 - By using the logistic function as follows:
$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

This is the logistic regression model or hypothesis function
 - And then applying thresholding *after* learning the model to predict the output as follows:

$$\begin{cases} \text{if } h_{\theta}(x) < 0.5 \text{ predict 0} \\ \text{if } h_{\theta}(x) \geq 0.5 \text{ predict 1} \end{cases}$$

Towards Identifying the Logistic Regression Cost Function

- How to learn a *logistic regression model* $\mathbf{h}_{\theta}(x) = \mathbf{g}(\boldsymbol{\theta}^T x)$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - Perhaps, by minimizing *Mean Squared Error (MSE)*. That is:

$$\text{minimize} \frac{1}{2n} \sum_{i=1}^n (y'^{(i)} - y^{(i)})^2$$

≡

$$\text{minimize}_{\theta} \frac{1}{2n} \sum_{i=1}^n ((g(\boldsymbol{\theta}^T x))^{(i)} - y^{(i)})^2$$

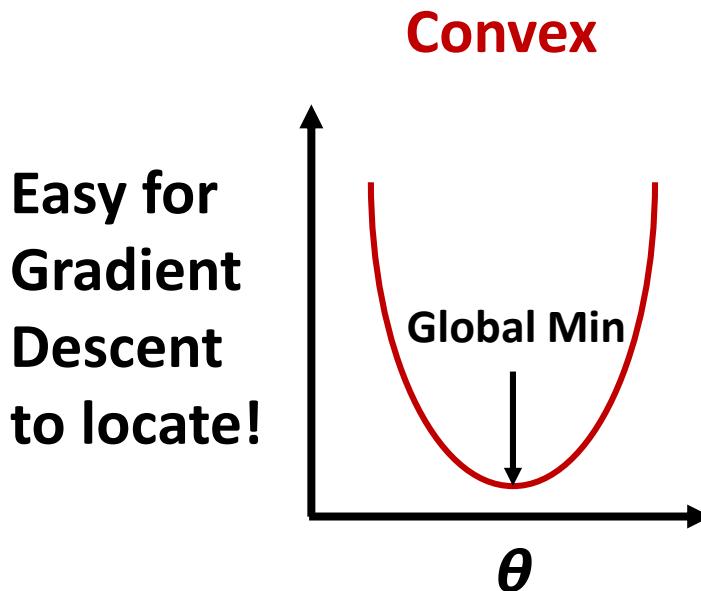
≡

$$\text{minimize}_{\theta} J(\boldsymbol{\theta})$$

objective function
 $J(\boldsymbol{\theta})$

Towards Identifying the Logistic Regression Cost Function

- How to learn a *logistic regression model* $\mathbf{h}_{\theta}(x) = \mathbf{g}(\boldsymbol{\theta}^T x)$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - Perhaps, by minimizing **Mean Squared Error (MSE)**. That is:



$$\text{minimize}_{\theta} \frac{1}{2n} \sum_{i=1}^n (y'^{(i)} - y^{(i)})^2$$

$$\equiv$$
$$\text{minimize}_{\theta} \frac{1}{2n} \sum_{i=1}^n ((g(\boldsymbol{\theta}^T x))^{(i)} - y^{(i)})^2$$

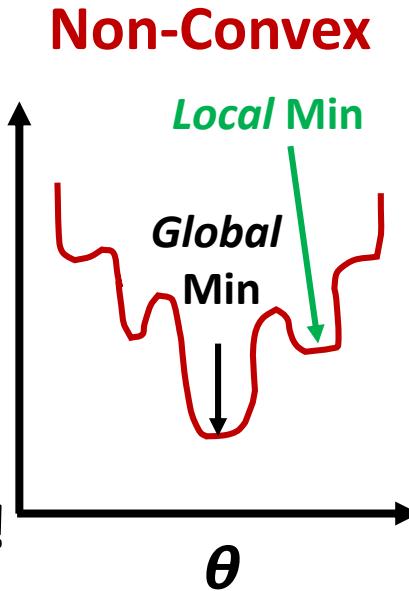
$$\equiv$$
$$\text{minimize}_{\theta} J(\boldsymbol{\theta})$$

Unfortunately, if we plot this cost function, it will turn out to be “**non-convex**”

Towards Identifying the Logistic Regression Cost Function

- How to learn a *logistic regression model* $\mathbf{h}_\theta(\mathbf{x}) = \mathbf{g}(\boldsymbol{\theta}^T \mathbf{x})$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $\mathbf{x} = [x_0, \dots, x_m]$?
 - Perhaps, by minimizing **Mean Squared Error (MSE)**. That is:

Gradient Descent might get stuck at a *local* min and fail to locate the *global* min!



$$\text{minimize}_{\theta} \frac{1}{2n} \sum_{i=1}^n (y'^{(i)} - y^{(i)})^2$$

≡

$$\text{minimize}_{\theta} \frac{1}{2n} \sum_{i=1}^n ((g(\boldsymbol{\theta}^T \mathbf{x}))^{(i)} - y^{(i)})^2$$

≡

$$\text{minimize}_{\theta} J(\boldsymbol{\theta})$$

Unfortunately, if we plot this cost function, it will turn out to be "**non-convex**"

The Logistic Regression Cost Function

- How to learn a *logistic regression model* $h_{\theta}(x) = g(\theta^T x)$, where $\theta = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - Let us try a different cost function. That is:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

Equivalent To

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

The Logistic Regression Cost Function

- How to learn a *logistic regression model* $\mathbf{h}_{\theta}(\mathbf{x}) = \mathbf{g}(\boldsymbol{\theta}^T \mathbf{x})$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $\mathbf{x} = [x_0, \dots, x_m]$?
 - Let us try a different cost function. That is:

$$\text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}), y) = -y \log(\mathbf{h}_{\theta}(\mathbf{x})) - (1 - y) \log(1 - \mathbf{h}_{\theta}(\mathbf{x}))$$

This function still assumes real-valued outputs for $\mathbf{h}_{\theta}(\mathbf{x})$ (i.e., still entails a *regression problem*), while logistic regression should predict discrete values (i.e., logistic regression is a *classification problem*)

The Logistic Regression Cost Function

- How to learn a *logistic regression model* $h_{\theta}(x) = g(\theta^T x)$, where $\theta = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - Let us try a different cost function. That is:

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

We still need to apply to it the logistic function:

$$g(z) = \frac{1}{1 + e^{-z}}$$

The Logistic Regression Cost Function

- How to learn a *logistic regression model* $\mathbf{h}_{\theta}(x) = \mathbf{g}(\boldsymbol{\theta}^T x)$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - By minimizing the following cost function:

$$\text{Cost}(\mathbf{h}_{\theta}(x), y) = -y \log(\mathbf{h}_{\theta}(x)) - (1 - y) \log(1 - \mathbf{h}_{\theta}(x))$$

$$= -y \log(\mathbf{g}(\boldsymbol{\theta}^T x)) - (1 - y) \log(1 - \mathbf{g}(\boldsymbol{\theta}^T x))$$

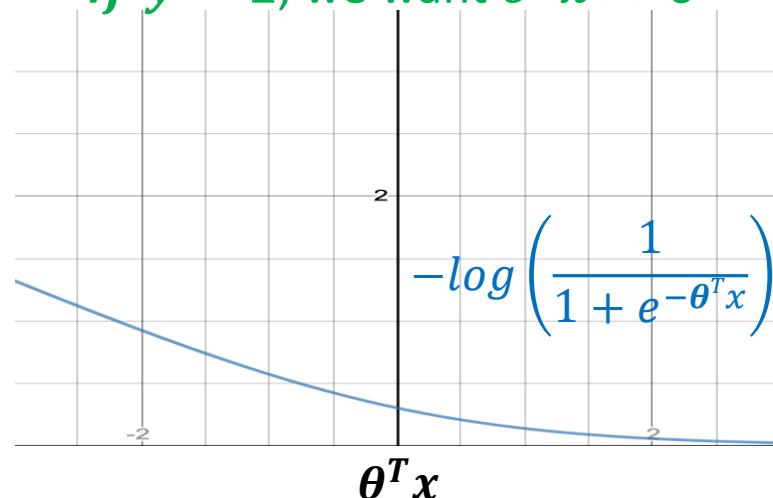
$$= -y \log\left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T x}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\boldsymbol{\theta}^T x}}\right)$$

The Logistic Regression Cost Function

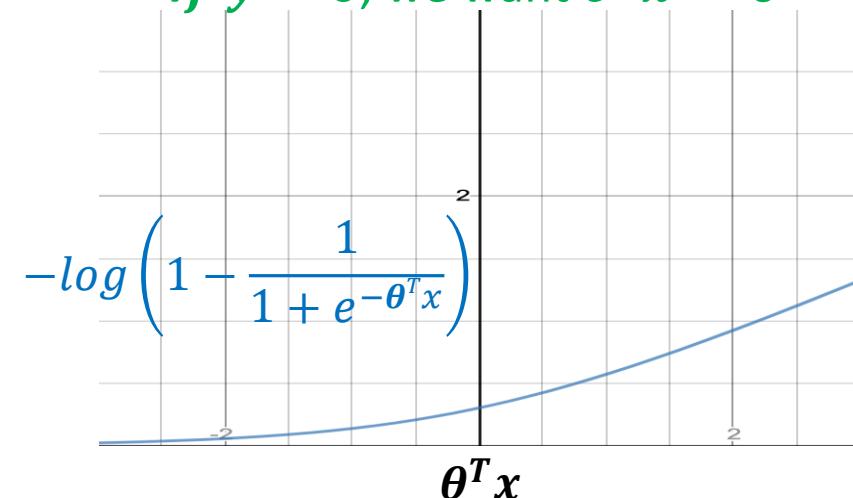
- How to learn a *logistic regression model* $h_{\theta}(x) = g(\theta^T x)$, where $\theta = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - By minimizing the following cost function:

$$\text{Cost}(h_{\theta}(x), y) = -y \log\left(\frac{1}{1 + e^{-\theta^T x}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

If $y = 1$, we want $\theta^T x >> 0$



If $y = 0$, we want $\theta^T x << 0$



Learning a Logistic Regression Model

- How to learn a *logistic regression model* $\mathbf{h}_{\theta}(x) = \mathbf{g}(\boldsymbol{\theta}^T x)$, where $\boldsymbol{\theta} = [\theta_0, \dots, \theta_m]$ and $x = [x_0, \dots, x_m]$?
 - By minimizing the following cost function:

$$\text{Cost}(\mathbf{h}_{\theta}(x), y) = -y \log\left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T x}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\boldsymbol{\theta}^T x}}\right)$$

- That is:

$$\underset{\theta}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \text{Cost}(\mathbf{h}_{\theta}(x^{(i)}), y^{(i)})$$

≡

$$\underset{\theta}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n -y^{(i)} \log\left(\frac{1}{1 + e^{-\boldsymbol{\theta}^T x^{(i)}}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\boldsymbol{\theta}^T x^{(i)}}}\right)$$

Cost function
 $J(\boldsymbol{\theta})$

Gradient Descent For Logistic Regression

- **Outline:**
 - Have the objective function $J(\theta)$, where $\theta = [\theta_0, \dots, \theta_m]$
 - Start off with some guesses for $\theta_0, \dots, \theta_m$
 - It does not really matter what values you start off with, but a common choice is to set them all initially to zero
 - Repeat until convergence{

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

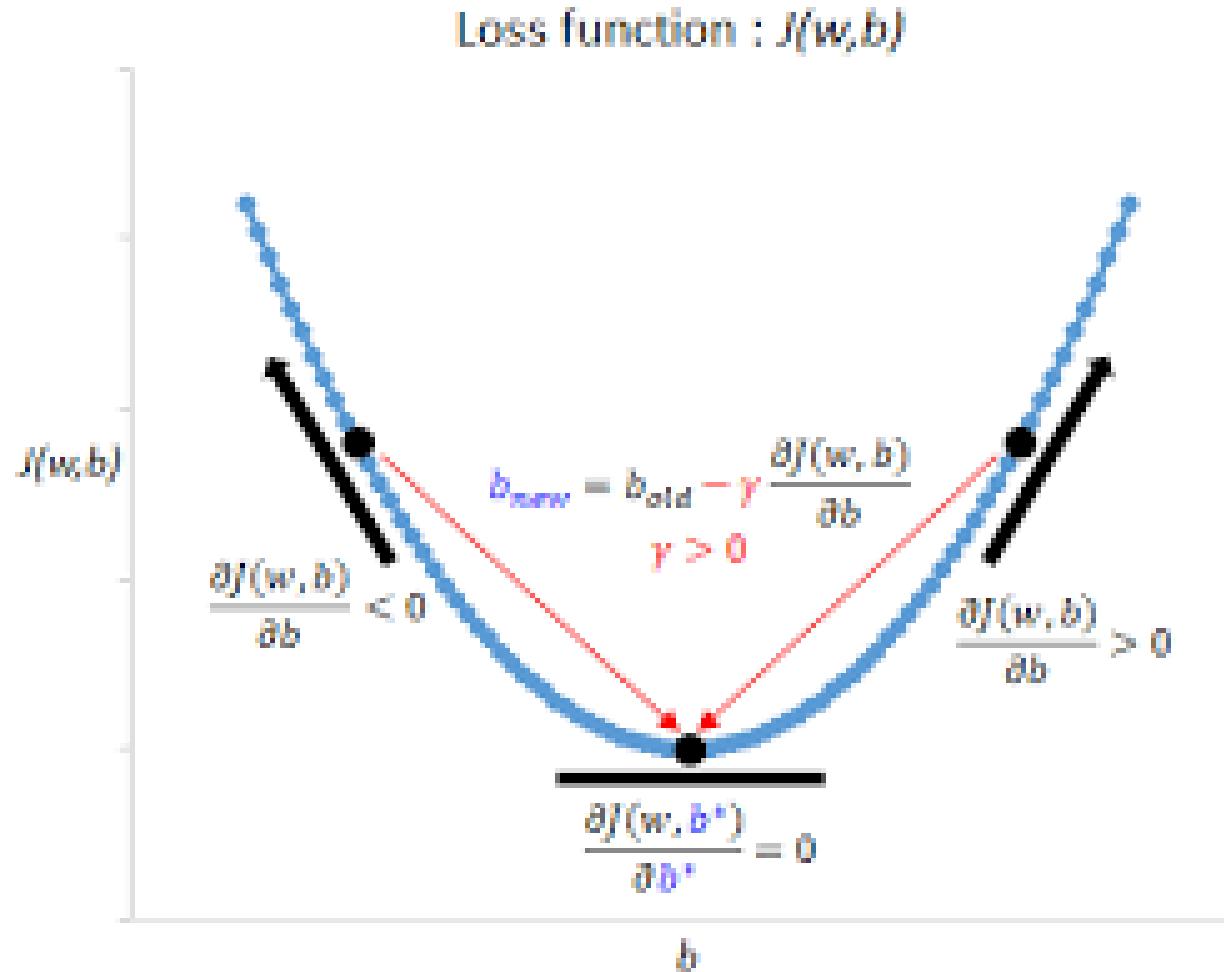
Learning rate, which controls how big a step we take when we update θ_j

Partial derivative

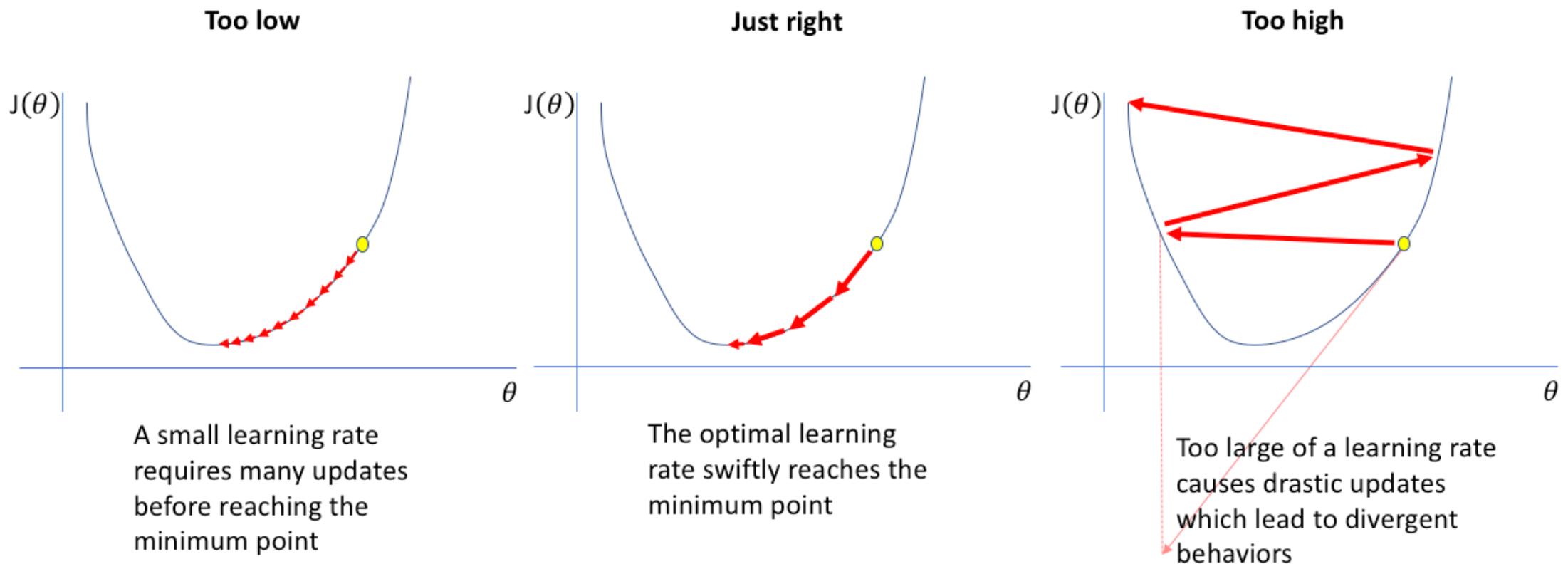
Note: Update all θ_j simultaneously

}

Gradient Descent



Gradient Descent



Gradient Descent For Logistic Regression

- **Outline:**
 - Have cost function $J(\theta)$, where $\theta = [\theta_0, \dots, \theta_m]$
 - Start off with some guesses for $\theta_0, \dots, \theta_m$
 - It does not really matter what values you start off with, but a common choice is to set them all initially to zero
 - Repeat until convergence{

$$\theta_j = \theta_j - \alpha \left(\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_j^{(i)} \right)$$

*The final formula
after applying
partial derivatives*

Inference After Learning

- After learning the parameters $\theta = [\theta_0, \dots, \theta_m]$, we can predict the output of any new unseen $x = [x_0, \dots, x_m]$ as follows:

$$\begin{cases} \text{if } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} < 0.5 \text{ predict 0} \\ \text{Else if } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \geq 0.5 \text{ predict 1} \end{cases}$$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

	and	vaccine	the	of	nigeria	y
Email a	1	1	0	1	1	1
Email b	0	0	1	1	0	0
Email c	0	1	1	0	0	1
Email d	1	0	0	1	0	0
Email e	1	0	1	0	1	1
Email f	1	0	1	1	0	0

A Training Dataset

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

	and	vaccine	the	of	nigeria	y
Email a	1	1	0	1	1	1
Email b	0	0	1	1	0	0
Email c	0	1	1	0	0	1
Email d	1	0	0	1	0	0
Email e	1	0	1	0	1	1
Email f	1	0	1	1	0	0



1 entails that a word (i.e., "and") is *present* in an email (i.e., "Email a")

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

	and	vaccine	the	of	nigeria	y
Email a	1	1	0	1	1	1
Email b	0	0	1	1	0	0
Email c	0	1	1	0	0	1
Email d	1	0	0	1	0	0
Email e	1	0	1	0	1	1
Email f	1	0	1	1	0	0



0 entails that a word (i.e., "and") is *absent* in an email (i.e., "Email b")

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$
- 5 words (or *features*) = $[x_1, x_2, x_3, x_4, x_5]$
- We define 6 parameters
(the first one, i.e., θ_0 ,
is the intercept)

	$x_1 = \text{and}$	$x_2 = \text{vaccine}$	$x_3 = \text{the}$	$x_4 = \text{of}$	$x_5 = \text{nigeria}$	y
Email a	1	1	0	1	1	1
Email b	0	0	1	1	0	0
Email c	0	1	1	0	0	1
Email d	1	0	0	1	0	0
Email e	1	0	1	0	1	1
Email f	1	0	1	1	0	0

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$ The parameter vector:
 $\theta = [\theta_0, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5]$
- $x = [x_0, x_1, x_2, x_3, x_4, x_5] \rightarrow$ The feature vector

	$x_0 = 1$	$x_1 = \text{and}$	$x_2 = \text{vaccine}$	$x_3 = \text{the}$	$x_4 = \text{of}$	$x_5 = \text{nigeria}$	y
Email a	1	1	1	0	1	1	1
Email b	1	0	0	1	1	0	0
Email c	1	0	1	1	0	0	1
Email d	1	1	0	0	1	0	0
Email e	1	1	0	1	0	1	1
Email f	1	1	0	1	1	0	0

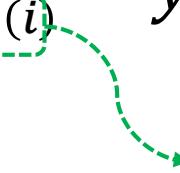
To account for the intercept

Recap: Gradient Descent For Logistic Regression

- **Outline:**
 - Have cost function $J(\theta)$, where $\theta = [\theta_0, \dots, \theta_m]$
 - Start off with some guesses for $\theta_0, \dots, \theta_m$
 - It does not really matter what values you start off with, but a common choice is to set them all initially to zero
 - Repeat until convergence{

$$\theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_j^{(i)}$$

}



First, let us calculate this factor
for every example in our
training dataset

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$
[1,1,1,0,1,1]	1	$[0,0,0,0,0,0] \times [1,1,1,0,1,1] = 0$
[1,0,0,1,1,0]	0	$[0,0,0,0,0,0] \times [1,0,0,1,1,0] = 0$
[1,0,1,1,0,0]	1	$[0,0,0,0,0,0] \times [1,0,1,1,0,0] = 0$
[1,1,0,0,1,0]	0	$[0,0,0,0,0,0] \times [1,1,0,0,1,0] = 0$
[1,1,0,1,0,1]	1	$[0,0,0,0,0,0] \times [1,1,0,1,0,1] = 0$
[1,1,0,1,1,0]	0	$[0,0,0,0,0,0] \times [1,1,0,1,1,0] = 0$

Recap: Gradient Descent For Logistic Regression

- **Outline:**

- Have cost function $J(\theta)$, where $\theta = [\theta_0, \dots, \theta_m]$
- Start off with some guesses for $\theta_0, \dots, \theta_m$
 - It does not really matter what values you start off with, but a common choice is to set them all initially to zero
- Repeat until convergence{

$$\theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_j^{(i)}$$

Second, let us calculate this equation for every example in our training dataset and for every θ_j , where j is between 0 and m

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_0$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_0$
[1,1,1,0,1,1]	1	[0,0,0,0,0,0]×[1,1,1,0,1,1]=0	$(\frac{1}{1+e^{-0}} - 1) \times 1 = -0.5$
[1,0,0,1,1,0]	0	[0,0,0,0,0,0]×[1,0,0,1,1,0]=0	$(\frac{1}{1+1} - 0) \times 1 = 0.5$
[1,0,1,1,0,0]	1	[0,0,0,0,0,0]×[1,0,1,1,0,0]=0	$(\frac{1}{1+1} - 1) \times 1 = -0.5$
[1,1,0,0,1,0]	0	[0,0,0,0,0,0]×[1,1,0,0,1,0]=0	$(\frac{1}{1+1} - 0) \times 1 = 0.5$
[1,1,0,1,0,1]	1	[0,0,0,0,0,0]×[1,1,0,1,0,1]=0	$(\frac{1}{1+1} - 1) \times 1 = -0.5$
[1,1,0,1,1,0]	0	[0,0,0,0,0,0]×[1,1,0,1,1,0]=0	$(\frac{1}{1+1} - 0) \times 1 = 0.5$

Recap: Gradient Descent For Logistic Regression

- **Outline:**
 - Have cost function $J(\theta)$, where $\theta = [\theta_0, \dots, \theta_m]$
 - Start off with some guesses for $\theta_0, \dots, \theta_m$
 - It does not really matter what values you start off with, but a common choice is to set them all initially to zero
 - Repeat until convergence{

$$\theta_j = \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_j^{(i)}$$

}

Third, let us compute every θ_j

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_0$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0.5
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	-0.5
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0.5
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	-0.5
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_0^{(i)} = 0$$

Then,

$$\theta_0 = \theta_0 - \alpha \times 0$$

New θ_0

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_0$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0.5
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	-0.5
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0.5
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	-0.5
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_0^{(i)} = 0$$

Then,

$$\theta_0 = \theta_0 - \alpha \times 0$$

Old θ_0

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_0$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0.5
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	-0.5
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0.5
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	-0.5
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_0^{(i)} = 0$$

Then,

$$\begin{aligned} \theta_0 &= \theta_0 - \alpha \times 0 \\ &= 0 - 0.5 \times 0 = 0 \end{aligned}$$

New Parameter Vector:
 $\theta = [0, \theta_1, \theta_2, \theta_3, \theta_4, \theta_5]$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_1$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_1$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	0
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0.5
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	-0.5
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_1^{(i)} = 0$$

Then,

$$\begin{aligned} \theta_1 &= \theta_1 - \alpha \times 0 \\ &= 0 - 0.5 \times 0 = 0 \end{aligned}$$

New Parameter Vector:
 $\theta = [0, 0, \theta_2, \theta_3, \theta_4, \theta_5]$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_2$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_2$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	-0.5
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	0
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_2^{(i)} = -1$$

Then,

$$\begin{aligned} \theta_2 &= \theta_2 - \alpha \times (-1) \\ &= 0 - 0.5 \times (-1)/6 = 0.5/6 \end{aligned}$$

New Parameter Vector:
 $\theta = [0, 0, 0.5/6, \theta_3, \theta_4, \theta_5]$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x^T	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_3$
[1,1,1,0,1,1]	1	[0,0,0,0,0]×[1,1,1,0,1,1]=0	
[1,0,0,1,1,0]	0	[0,0,0,0,0]×[1,0,0,1,1,0]=0	
[1,0,1,1,0,0]	1	[0,0,0,0,0]×[1,0,1,1,0,0]=0	
[1,1,0,0,1,0]	0	[0,0,0,0,0]×[1,1,0,0,1,0]=0	
[1,1,0,1,0,1]	1	[0,0,0,0,0]×[1,1,0,1,0,1]=0	
[1,1,0,1,1,0]	0	[0,0,0,0,0]×[1,1,0,1,1,0]=0	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_3$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	0
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0.5
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	-0.5
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	-0.5
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_3^{(i)} = 0$$

Then,

$$\theta_3 = \theta_3 - \alpha \times 0$$

$$= 0 - 0.5 \times 0 = 0$$

New Parameter Vector:
 $\theta = [0, 0, 0.1, 0, \theta_4, \theta_5]$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_4$
[1,1,1,0,1,1]	1	[0,0,0,0,0]×[1,1,1,0,1,1]=0	
[1,0,0,1,1,0]	0	[0,0,0,0,0]×[1,0,0,1,1,0]=0	
[1,0,1,1,0,0]	1	[0,0,0,0,0]×[1,0,1,1,0,0]=0	
[1,1,0,0,1,0]	0	[0,0,0,0,0]×[1,1,0,0,1,0]=0	
[1,1,0,1,0,1]	1	[0,0,0,0,0]×[1,1,0,1,0,1]=0	
[1,1,0,1,1,0]	0	[0,0,0,0,0]×[1,1,0,1,1,0]=0	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_4$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	-0.5
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	0.5
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	0
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	0.5
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	0
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	0.5

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_4^{(i)} = 1$$

Then,

$$\theta_4 = \theta_4 - \alpha \times 1$$

$$= 0 - 0.5 \times 1/6 = -0.5/6$$

New Parameter Vector:

$$\theta = [0, 0, 0.5/6, 0, -0.5/6, \theta_5]$$

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_5$
[1,1,1,0,1,1]	1	$[0,0,0,0,0] \times [1,1,1,0,1,1] = 0$	
[1,0,0,1,1,0]	0	$[0,0,0,0,0] \times [1,0,0,1,1,0] = 0$	
[1,0,1,1,0,0]	1	$[0,0,0,0,0] \times [1,0,1,1,0,0] = 0$	
[1,1,0,0,1,0]	0	$[0,0,0,0,0] \times [1,1,0,0,1,0] = 0$	
[1,1,0,1,0,1]	1	$[0,0,0,0,0] \times [1,1,0,1,0,1] = 0$	
[1,1,0,1,1,0]	0	$[0,0,0,0,0] \times [1,1,0,1,1,0] = 0$	

A Concrete Example: The Training Phase

- Let us apply logistic regression on the spam email recognition problem, assuming $\alpha = 0.5$ and starting with $\theta = [0, 0, 0, 0, 0, 0]$

x	y	$\theta^T x$	$(\frac{1}{1+e^{-\theta^T x}} - y)x_5$
[1,1,1,0,1,1]	1	[0,0,0,0,0]×[1,1,1,0,1,1]=0	-0.5
[1,0,0,1,1,0]	0	[0,0,0,0,0]×[1,0,0,1,1,0]=0	0
[1,0,1,1,0,0]	1	[0,0,0,0,0]×[1,0,1,1,0,0]=0	0
[1,1,0,0,1,0]	0	[0,0,0,0,0]×[1,1,0,0,1,0]=0	0
[1,1,0,1,0,1]	1	[0,0,0,0,0]×[1,1,0,1,0,1]=0	-0.5
[1,1,0,1,1,0]	0	[0,0,0,0,0]×[1,1,0,1,1,0]=0	0

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{1 + e^{-\theta^T x^{(i)}}} - y^{(i)} \right) x_5^{(i)} = -1$$

Then,

$$\begin{aligned} \theta_5 &= \theta_5 - \alpha \times (-1) \\ &= 0 - 0.5 \times (-1)/6 = 0.5/6 \end{aligned}$$

New Parameter Vector:

$$\theta = [0, 0, 0.5/6, 0, -0.5/6, 0.5/6]$$

A Concrete Example: Inference

- Let us infer whether a given new email, say, $k = [1, 0, 1, 0, 0, 1]$ is a spam or not, using logistic regression with the just learnt parameter vector $\theta = [0, 0, 0.5/6, 0, -0.5/6, 0.5/6]$

	$x_0 = 1$	$x_1 = \text{and}$	$x_2 = \text{vaccine}$	$x_3 = \text{the}$	$x_4 = \text{of}$	$x_5 = \text{nigeria}$	y
Email a	1	1	1	0	1	1	1
Email b	1	0	0	1	1	0	0
Email c	1	0	1	1	0	0	1
Email d	1	1	0	0	1	0	0
Email e	1	1	0	1	0	1	1
Email f	1	1	0	1	1	0	0

Our Training Dataset

A Concrete Example: Inference

- Let us infer whether a given new email, say, $k = [1, 0, 1, 0, 0, 1]$ is a spam or not, using logistic regression with the just learnt parameter vector $\theta = [0, 0, 0.5/6, 0, -0.5/6, 0.5/6]$

	$x_0 = 1$	$x_1 = \text{and}$	$x_2 = \text{vaccine}$	$x_3 = \text{the}$	$x_4 = \text{of}$	$x_5 = \text{nigeria}$	y
Email a	1	1	1	0	1	1	1
Email b	1	0	0	1	1	0	0
Email c	1	0	1	1	0	0	1
Email d	1	1	0	0	1	0	0
Email e	1	1	0	1	0	1	1
Email f	1	1	0	1	1	0	0
Email k	1	0	1	0	0	1	?

A Concrete Example: Inference

- Let us infer whether a given new email, say, $k = [1, 0, 1, 0, 0, 1]$ is a spam or not, using logistic regression with the just learnt parameter vector $\theta = [0, 0, 0.5/6, 0, -0.5/6, 0.5/6]$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad \begin{bmatrix} 0 \\ 0 \\ 0.5/6 \\ 0 \\ -0.5/6 \\ 0.5/6 \end{bmatrix} [1, 0, 1, 0, 0, 1] = (0.5/6 \times 1) + (0.5/6 \times 1) = 1/6$$
$$= \frac{1}{1 + e^{-1/6}}$$

$\geq 0.5 \rightarrow$ Class 1 (i.e., Spam)

TF-IDF (Term frequency-Inverse dense or document frequency)

TF-IDF is scored between 0 and 1. The **higher the numerical weight value, the rarer the term**.
The **smaller the weight, the more common the term**.

For a term i in document j :

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

tf_{ij} = number of occurrences of i in j

df_i = number of documents containing i

N = total number of documents

Example

Let's cover an example of 3 documents -

Document 1 It is going to rain today.

Document 2 Today I am not going outside.

Document 3 I am going to watch the season premiere.

Word	Count
going	3
to	2
today	2
i	2
am	2
.	1
is	1
rain	1

Vocab of document

Words/ Documents	Document 1	Document 2	Document 3
going	0.16	0.16	0.12
to	0.16	0	0.12
today	0.16	0.16	0
i	0	0.16	0.12
am	0	0.16	0.12
it	0.16	0	0
is	0.16	0	0
rain	0.16	0	0

TF for the document

Words	IDF Value
going	$\log(3/3)$
to	$\log(3/2)$
today	$\log(3/2)$
i	$\log(3/2)$
am	$\log(3/2)$
It	$\log(3/1)$
is	$\log(3/1)$
rain	$\log(3/1)$

IDF for document

Words/ Documents	going	to	today	i	am	it	is	rain
Document 1	0	0.07	0.07	0	0	0.17	0.17	0.17
Document 2	0	0	0.07	0.07	0.07	0	0	0
Document 3	0	0.05	0	0.05	0.05	0	0	0

Remember, the final equation = TF-IDF = TF * IDF

SMSSpamCollection database

jupyter SMSSpamCollection 03/24/2023

Logout

Plain Text

```
1 ham=Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat...
2 ham=Ok lar... Joking wif u oni...
3 spam--Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry question(std txt rate)T&C's
apply 08452810075over18's
4 ham=U dun say so early hor... U c already then say...
5 ham=Nah I don't think he goes to usf, he lives around here though
6 spam--FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb ok! XXX std chgs to
send, £1.50 to rcv
7 ham-Even my brother is not like to speak with me. They treat me like aids patient.
8 ham-As per your request 'Melle Melle (Oru Minnaminunginte Nurungu Vettam)' has been set as your callertune for all callers. Press *9 to
copy your friends Callertune
9 spam--WINNER!! As a valued network customer you have been selected to receivea £900 prize reward! To claim call 09061701461. Claim code
KL341. Valid 12 hours only.
10 spam--Had your mobile 11 months or more? U R entitled to Update to the latest colour mobiles with camera for Free! Call The Mobile
Update Co FREE on 08002986030
11 ham-I've gonna be home soon and i don't want to talk about this stuff anymore tonight, k? I've cried enough today.
12 spam--SIX chances to win CASH! From 100 to 20,000 pounds txt> CSH11 and send to 87575. Cost 150p/day, 6days, 16+ TsandCs apply Reply HL
4 info
13 spam--URGENT! You have won a 1 week FREE membership in our £100,000 Prize Jackpot! Txt the word: CLAIM to No: 81010 T&C www.dbuk.net
LCLTD POBOX 4403LDNW1A7RN18
14 ham-I've been searching for the right words to thank you for this breather. I promise i wont take your help for granted and will fulfil my
promise. You have been wonderful and a blessing at all times.
15 ham-I HAVE A DATE ON SUNDAY WITH WILL!!
16 spam--XXXMobileMovieClub: To use your credit, click the WAP link in the next txt message or click here>> http://wap.
xxxmobilemovieclub.com?n=Q1KGIGHJJGCB
17 ham=Oh k...i'm watching here:)
18 ham=Eh u remember how 2 spell his name... Yes i did. He v naughty make until i v wet.
19 ham=Fine if that's the way u feel. That's the way its gotta b
20 spam--England v Macedonia - dont miss the goals/team news. Txt ur national team to 87077 eg ENGLAND to 87077 Try:WALES, SCOTLAND 4txt/
Ú1.20 POBOXox36504W45WQ 16+
21 ham=Is that seriously how you spell his name?
22 ham=I'm going to try for 2 months ha ha only joking
23 ham=So ü pay first lar... Then when is da stock comin...
24 ham=Aft i finish my lunch then i go str down lor. Ard 3 smth lor. U finish ur lunch already?
25 ham=FFFFFFFFFF. Alright no way I can meet up with you sooner?
26 ham=Just forced myself to eat a slice. I'm really not hungry tho. This sucks. Mark is getting worried. He knows I'm sick when I turn down
pizza. Lol
```

```
In [2]: import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split, cross_val_score

df = pd.read_csv('SMSSpamCollection', delimiter='\t', header=None)
```

```
In [3]: X_train_raw, X_test_raw, y_train, y_test = train_test_split(df[1], df[0])

vectorizer = TfidfVectorizer()
X_train = vectorizer.fit_transform(X_train_raw)
```

```
In [4]: classifier = LogisticRegression()
classifier.fit(X_train, y_train)

X_test = vectorizer.transform(['URGENT! Your Mobile No 1234 was awarded a Prize', 'Hey honey, whats up?', 'how are you'])
print(X_test)
```

```
In [3]: X_train_raw, X_test_raw, y_train, y_test = train_test_split(df[1],df[0])  
vectorizer = TfidfVectorizer()  
X_train = vectorizer.fit_transform(X_train_raw)
```

```
In [4]: classifier = LogisticRegression()  
classifier.fit(X_train, y_train)  
  
X_test = vectorizer.transform( ['URGENT! Your Mobile No 1234 was awarded a Prize', 'Hey honey, whats up?', 'how are you'] )  
print(X_test)
```

```
(0, 7511)    0.25466740289936407  
(0, 7192)    0.3461278313956526  
(0, 7014)    0.4264937405685157  
(0, 5300)    0.4190022423989018  
(0, 4679)    0.2984334164110299  
(0, 4430)    0.3820593750636906  
(0, 1142)    0.47206819237223135  
(1, 7276)    0.5599731789780006  
(1, 6995)    0.33337693489630504  
(1, 3403)    0.6304810448926348  
(1, 3346)    0.42164381904103376  
(2, 7505)    0.3925217496458997  
(2, 3439)    0.679214073335713  
(2, 1021)    0.6201571725277635
```

```
In [5]: predictions = classifier.predict(X_test)  
print(predictions)  
  
['spam' 'ham' 'ham']
```

Classification metrics

- ◆ **Terms:**

- P : positive examples, involving major interesting classes
- N : negative examples (other examples)
- TP: true positive examples (positive examples that are correctly classified by the classifier)
- TN: true negative examples (negative examples that are correctly classified by the classifier)
- FP: false positive examples (negative examples that are incorrectly marked as positive examples)
- FN: false negative examples (positive examples that are incorrectly marked as negative examples)
- ◆ Confusion matrix: a table (at least $m \times m$). $CM_{i,j}$, element of the first i -th line and j -th column is the number of examples known to be in class i but marked as j by the classifier.
 - Ideally, for a highly accurate classifier, most data samples should be represented by the elements on the diagonal from $CM_{1,1}$ and $CM_{m,m}$. Other elements are 0 or close to 0. In other words, FP and FN are close to 0.

Predicted Observed	Yes	No	Total
Yes	TP	FN	P
No	FP	TN	N
Total	P'	N'	$P + N$

Confusion matrix

Classification metrics

Measure	Formula
Accuracy or correct classification rate	$\frac{TP + TN}{P + N}$
Error rate or false classification rate	$\frac{FP + FN}{P + N}$
Sensitivity, true positive rate, or <i>recall</i>	$\frac{TP}{P}$
Specificity or true negative rate	$\frac{TN}{N}$
<i>Precision</i>	$\frac{TP}{TP + FP}$
<i>F</i> score: harmonic mean of precision or recall	$\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$
F_β (β is a non-negative real number)	$\frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$

Example

- ◆ We have trained a machine learning model to identify that whether the object in an image is a cat. Now we use 200 images to test its performance. Among 200 images, 170 contain cats and 30 do not contain cats. The identification result of the model is that 160 images contain cats and 40 do not contain cats.

$$\text{Precision: } P = \frac{TP}{TP+FP} = \frac{140}{140+20} = 87.5\%$$

$$\text{Recall: } R = \frac{TP}{P} = \frac{140}{170} = 93.3\%.$$

$$\text{Accuracy: } ACC = \frac{TP+TN}{P+N} = \frac{140+10}{170+30} = 85\%$$

Predicted Observed	Yes	No	Total
Yes	140	30	170
No	20	10	30
Total	160	40	200

Debugging a learning algorithm

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

- You have built you awesome linear regression model predicting price
- Work perfectly on you testing data

Debugging a learning algorithm

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

- You have built you awesome linear regression model predicting price
- Work perfectly on your training data
- Then it fails miserably when you test it on the new data you collected
- What to do now?

Things You Can Try

- Get more data
 - Try different features
 - Try tuning your hyperparameter
-
- But which should I try first?

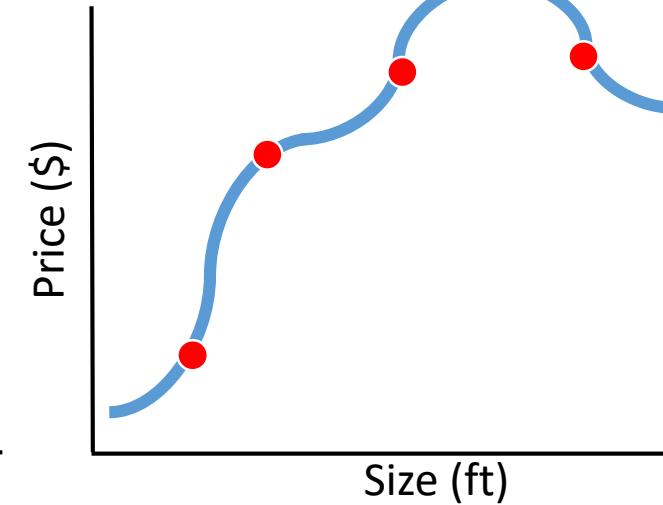
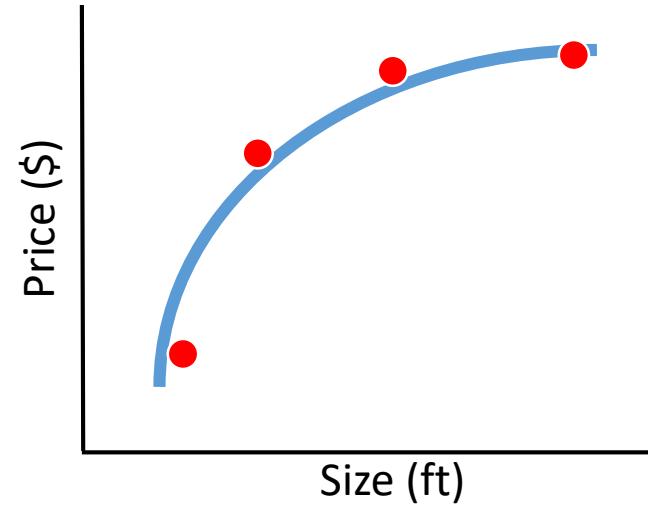
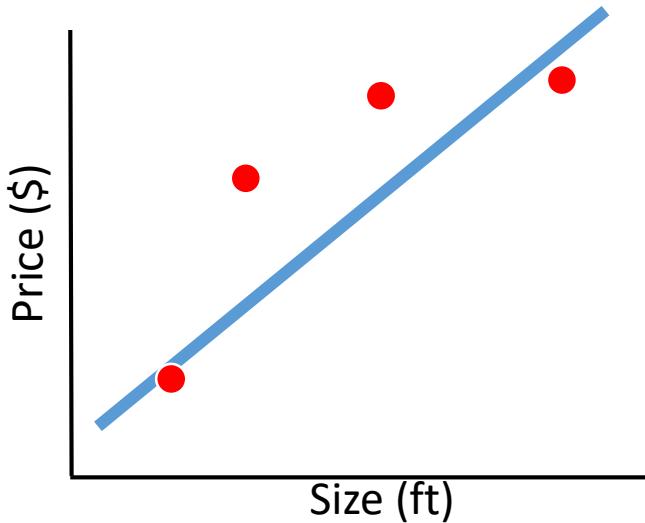
Diagnosing Machine Learning System

- Figure out what is wrong first
- Diagnosing your system takes time, but it can save your time as well
- Ultimate goal: **low generalization error**

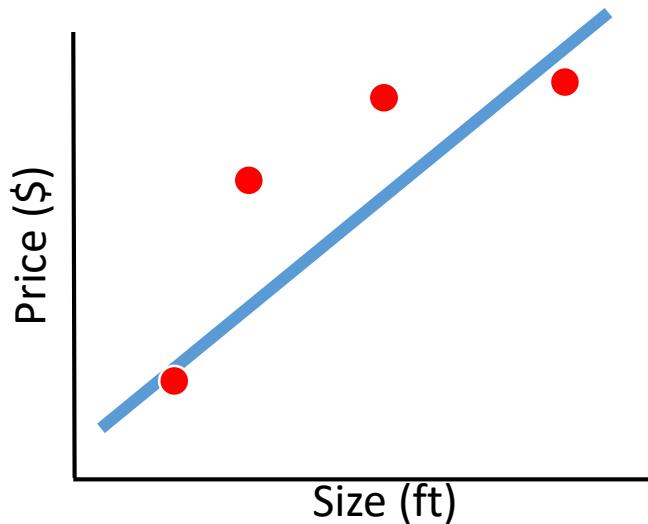
Problem: Fail to Generalize

- Model does not generalize to unseen data
 - Fail to predict things that are not in training sample
 - Pick a model that has **lower generalization error**

Evaluate Your Hypothesis

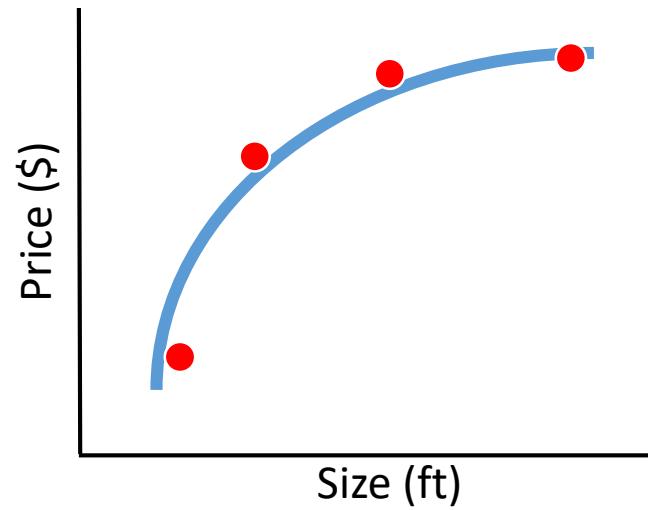


Evaluate Your Hypothesis



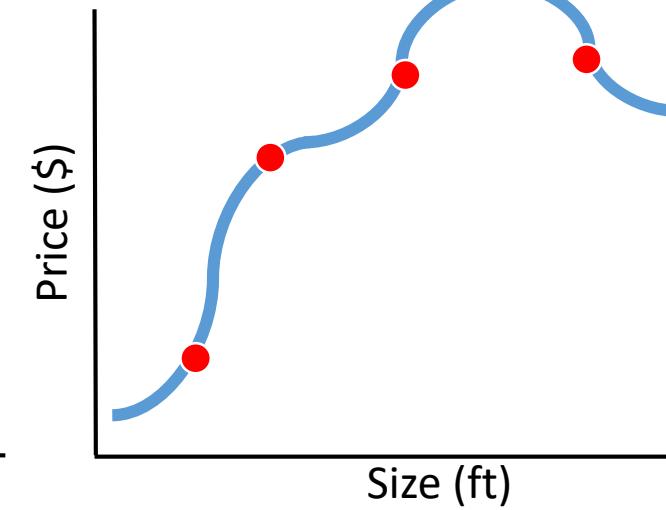
$$\theta_0 + \theta_1 x$$

Underfit



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

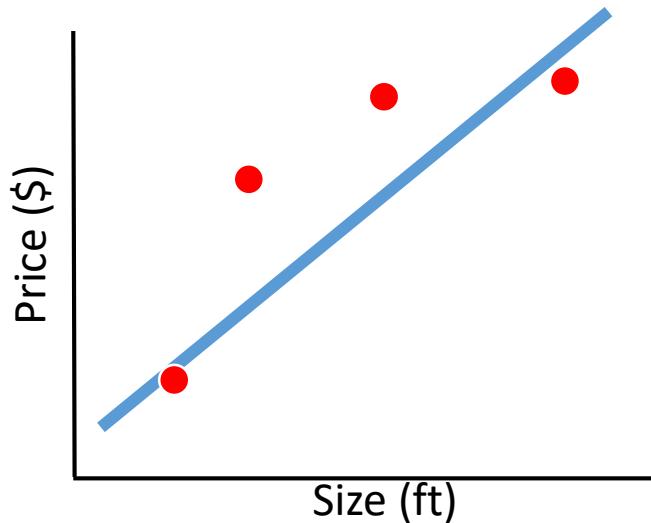
Just right



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

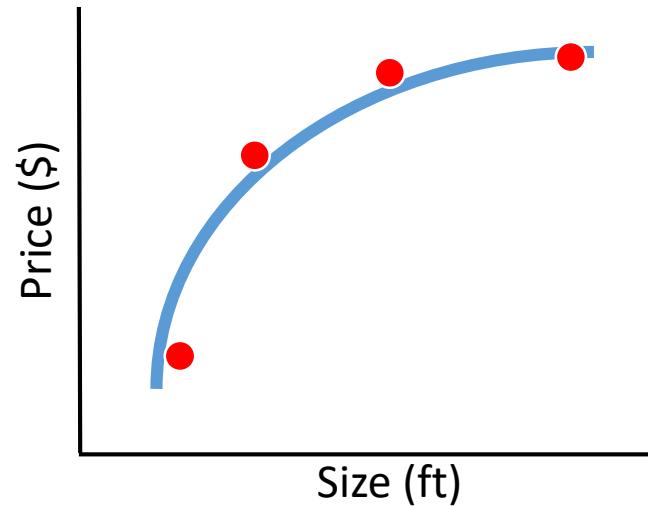
Overfit

Evaluate Your Hypothesis



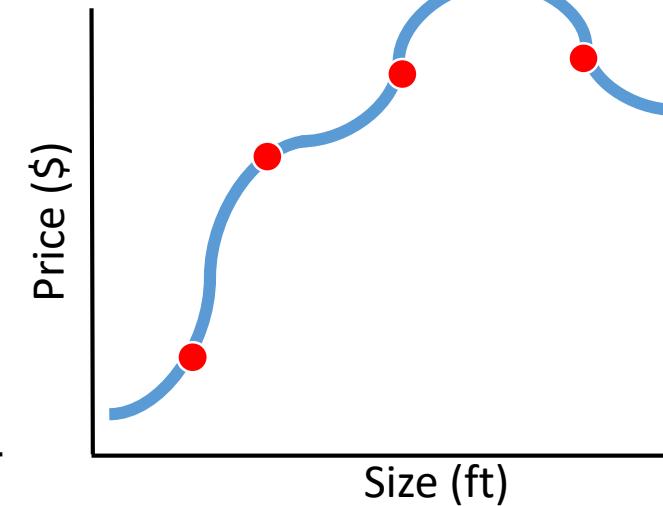
$$\theta_0 + \theta_1 x$$

Underfit



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Just right



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfit

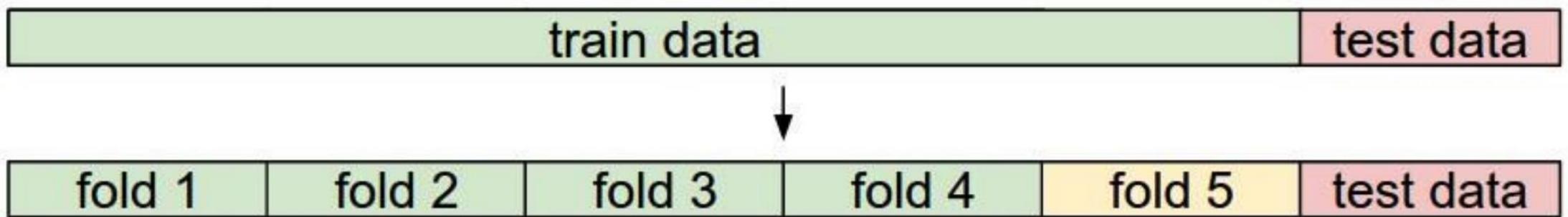
- What if the feature dimension is too high?

Model Selection

- Model does not generalize to unseen data
 - Fail to predict things that are not in training sample
 - Pick a model that has **lower generalization error**

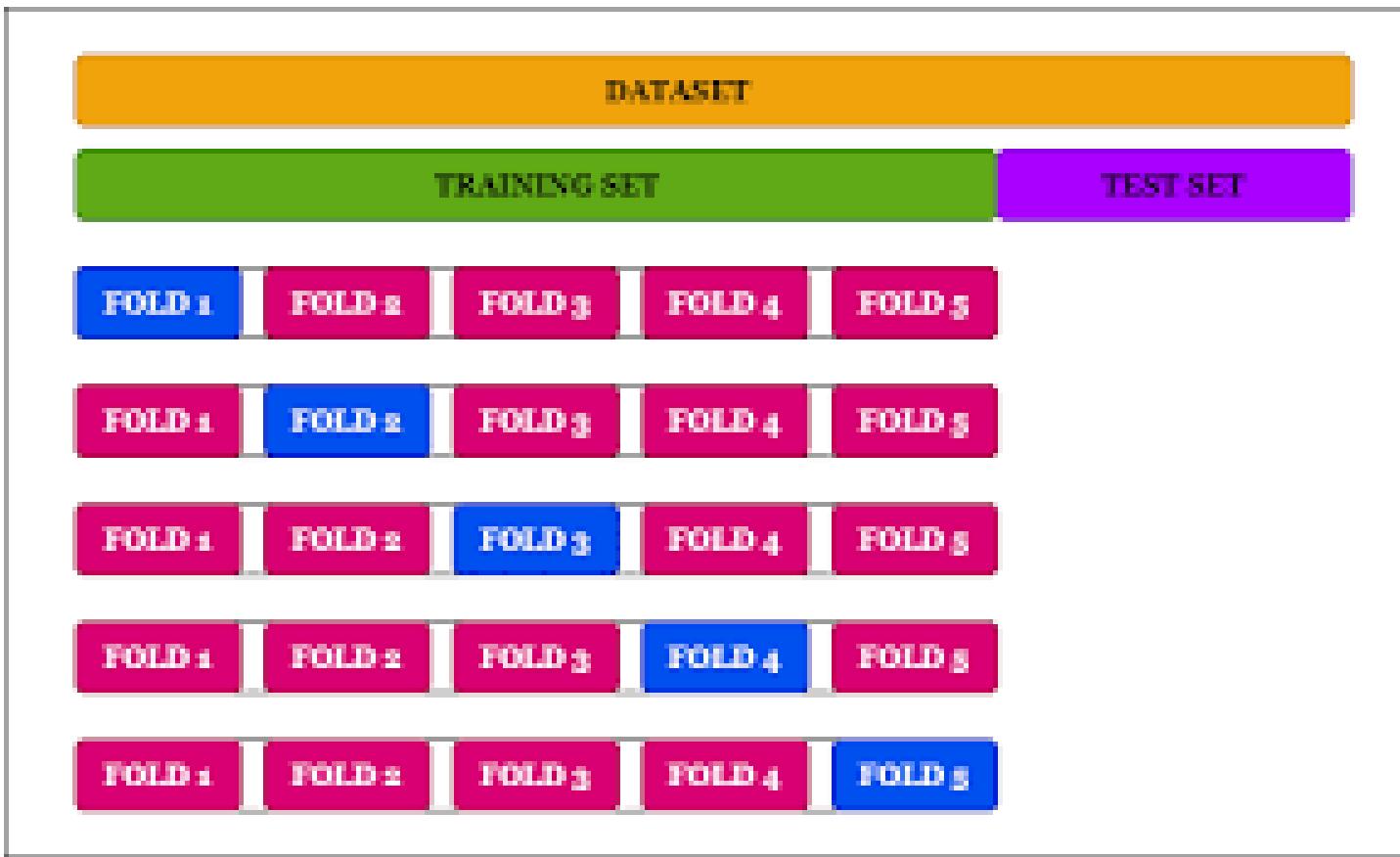
Model Selection

- Model does not generalize to unseen data
 - Fail to predict things that are not in training sample
 - Pick a model that has **lower generalization error**



Model Selection

- Model does not generalize to unseen data
 - Fail to predict things that are not in training sample
 - Pick a model that has **lower generalization error**
- How to evaluate generalization error?
 - Split your data into *train*, *validation*, and *test set*.
 - Use *test set error* as an **estimator** of generalization error
 - Technique K-fold cross validation



Model Selection

- Training error

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

- Validation error

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

- Test error

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Model Selection

- Training error

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

- Validation error

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

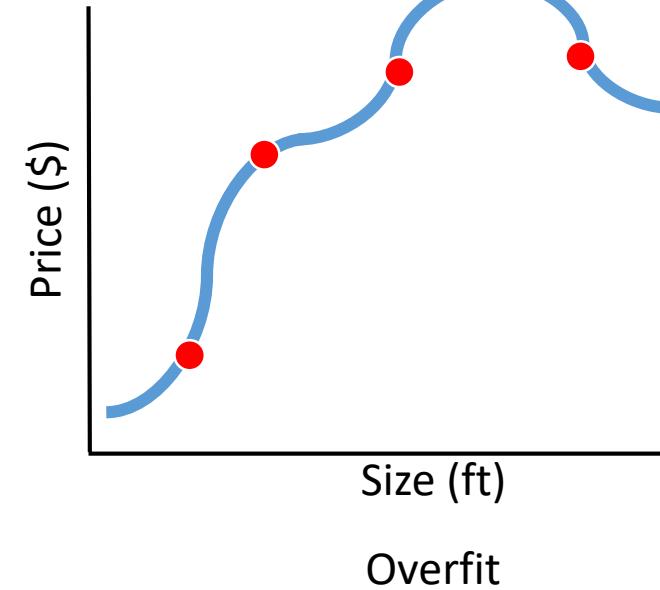
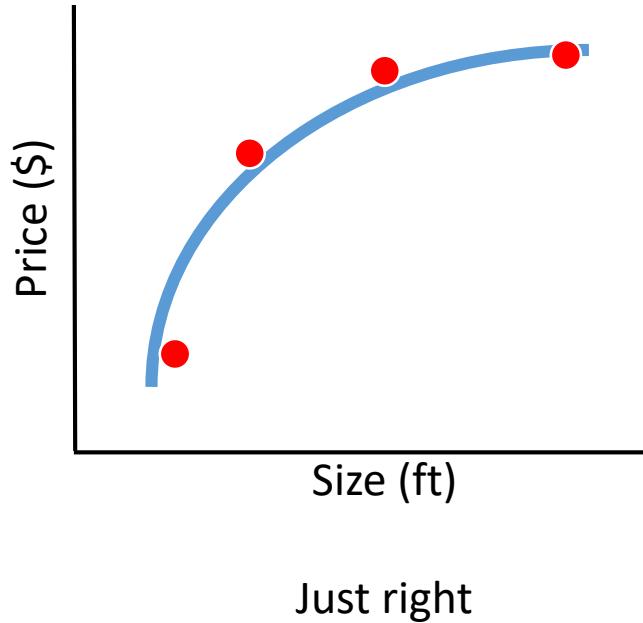
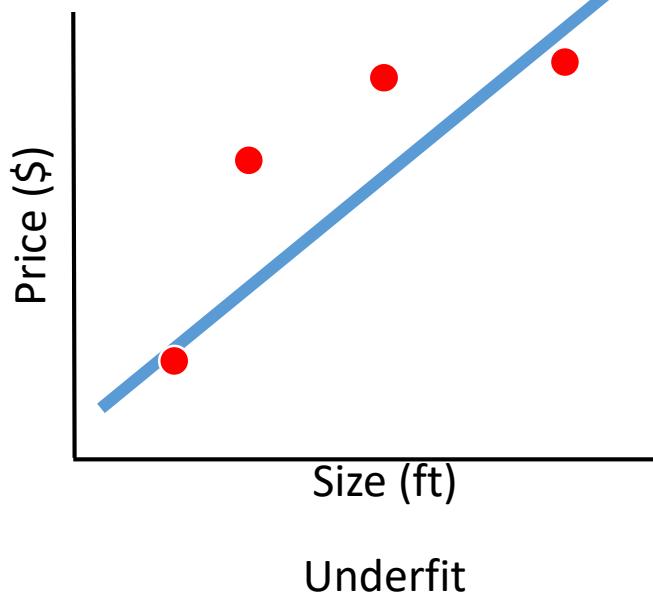
- Test error

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

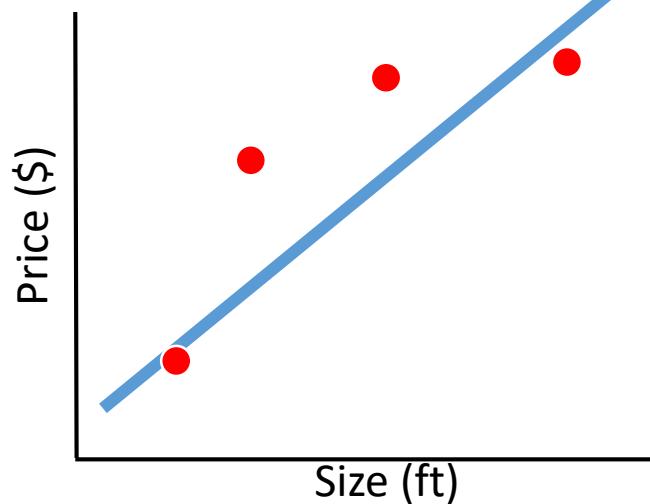
Procedure:

- Step 1. Train on training set
- Step 2. Evaluate validation error
- Step 3. Pick the best model based on Step 2.
- Step 4. Evaluate the test error

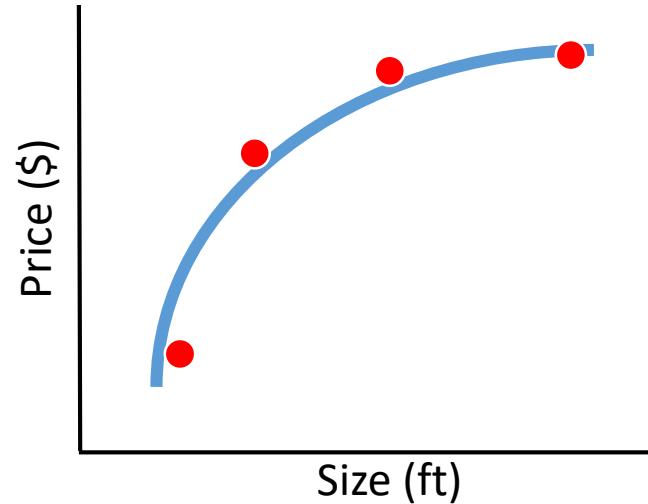
Bias/Variance Trade-off



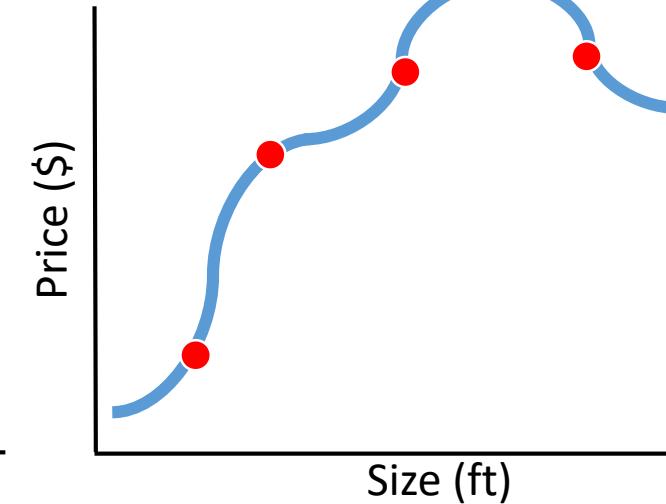
Bias/Variance Trade-off



Underfit
High bias

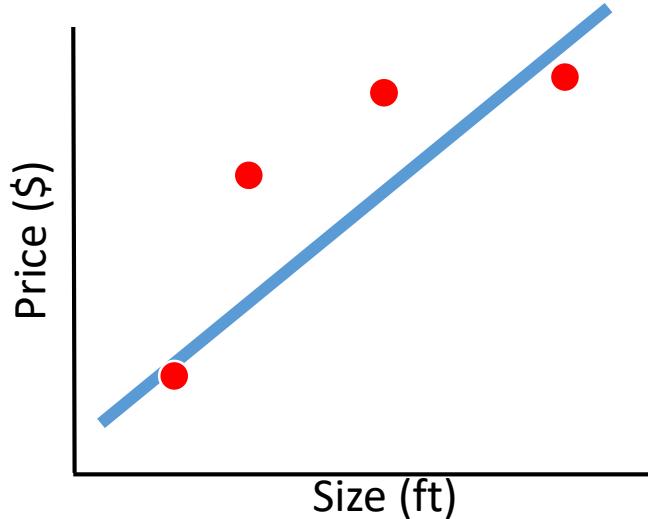


Just right

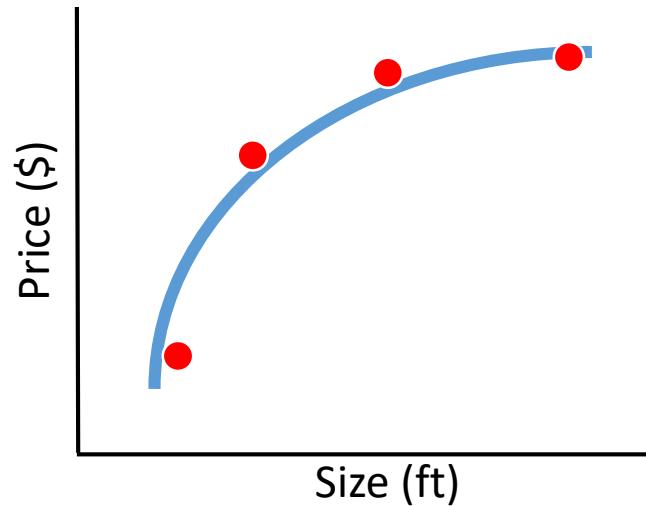


Overfit
High Variance

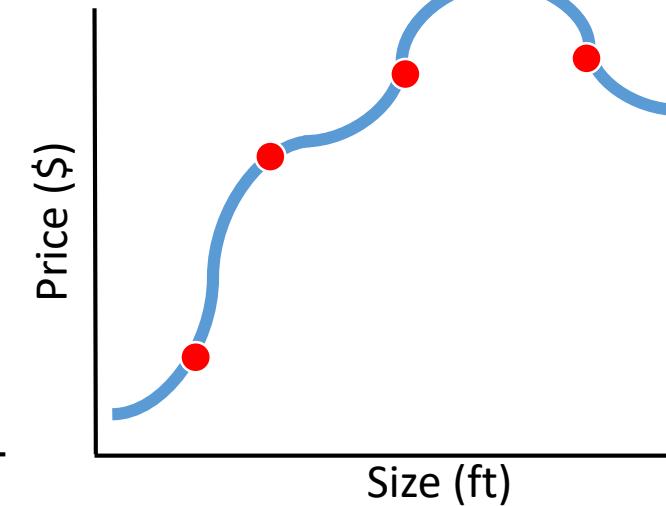
Bias/Variance Trade-off



Underfit
High bias
Too simple



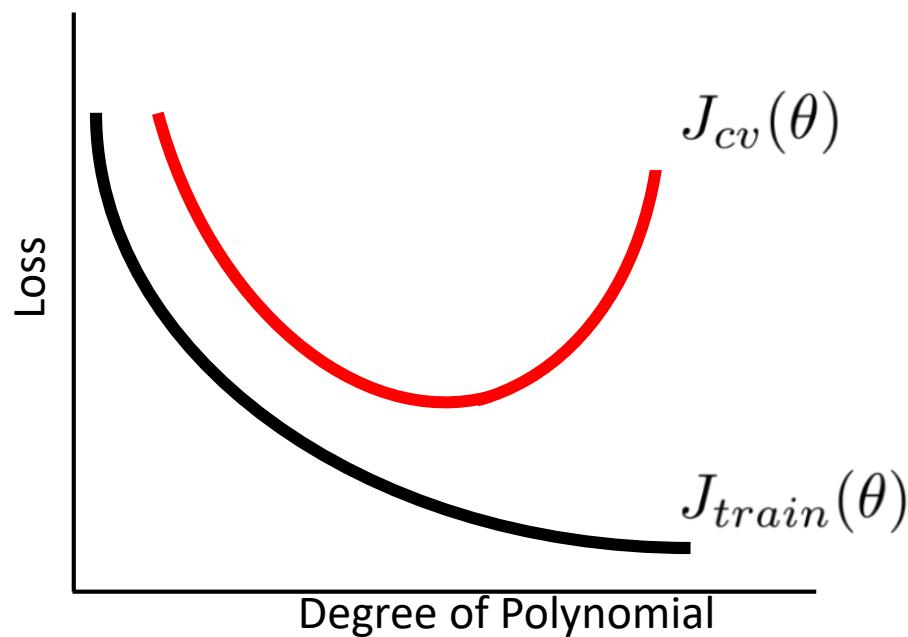
Just right



Overfit
High Variance
Too Complex

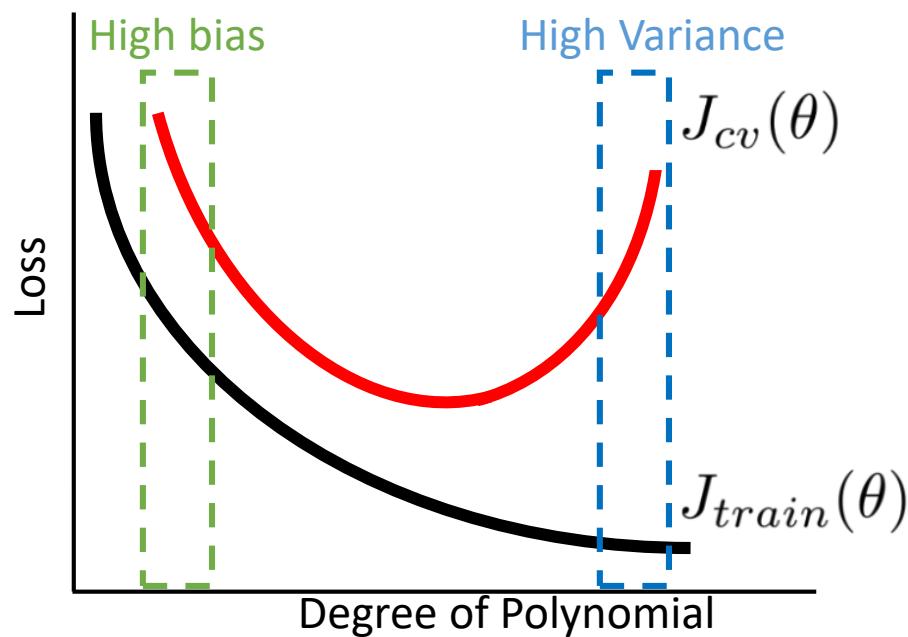
Bias / Variance Trade-off

- Training error $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
- Cross-validation error $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$



Bias / Variance Trade-off

- Training error $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
- Cross-validation error $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$

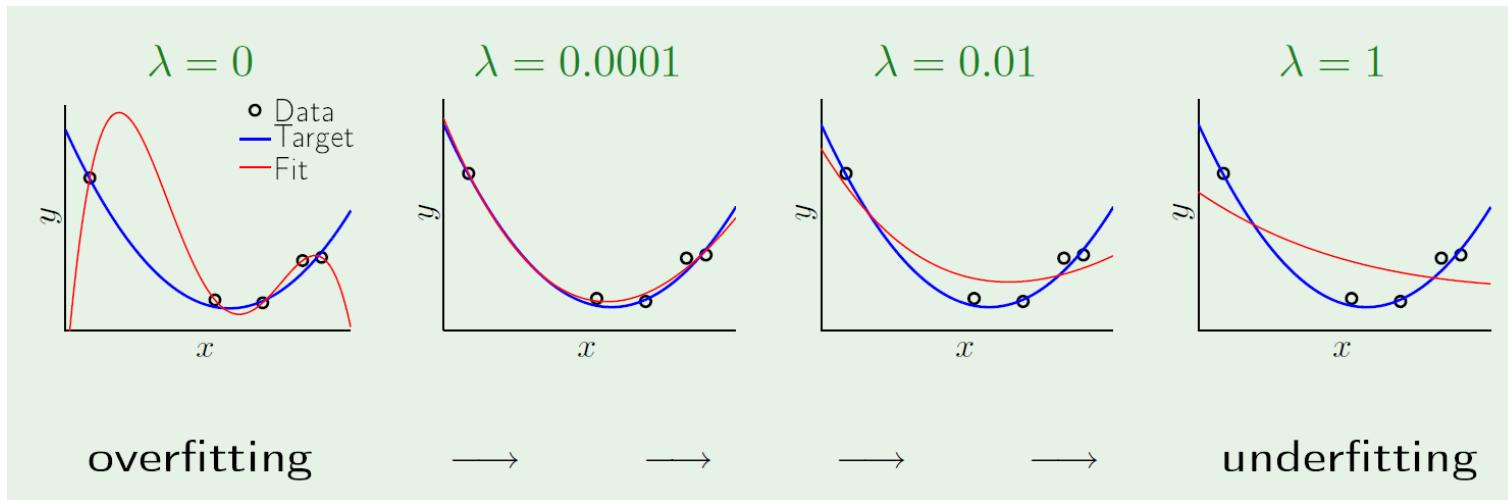


Regularization

- Involves shrinking the estimated coefficients toward zero relative to the OLS estimates; has the effect of reducing variance.
- Adds another term to the loss function to be reduced= λ . Regularization term, where $\lambda \geq 0$ is a *tuning parameter*, to be determined separately.
- The **weight decay coefficient λ** determines how dominant the regularization is during the gradient computation
- It could be
 - L1 regularization : Lasso (order 1)
 - L2 regularization: Ridge (order 2)

Regularization or Shrinkage

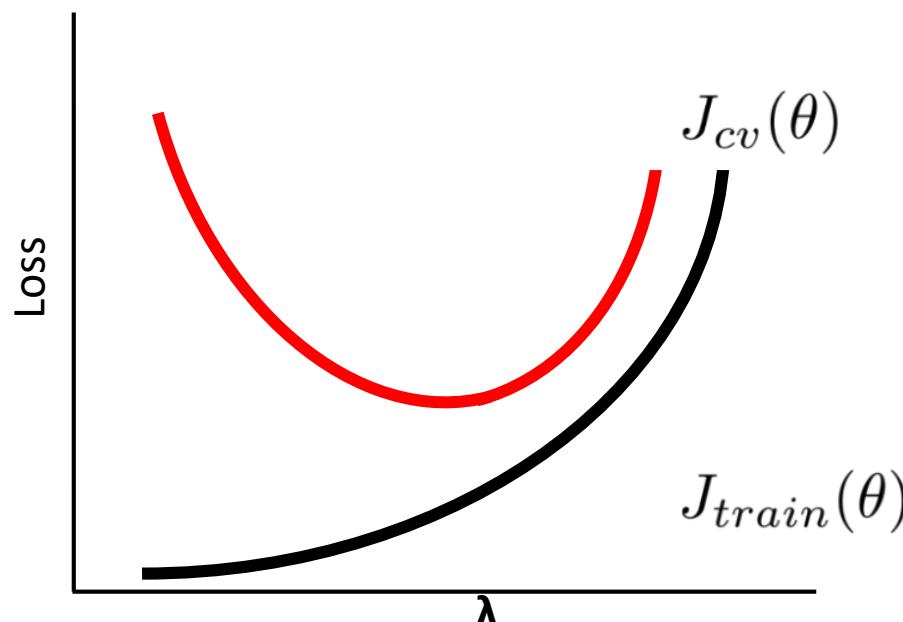
- Look at what a little regularization can do:



Bias / Variance Trade-off with Regularization

L2

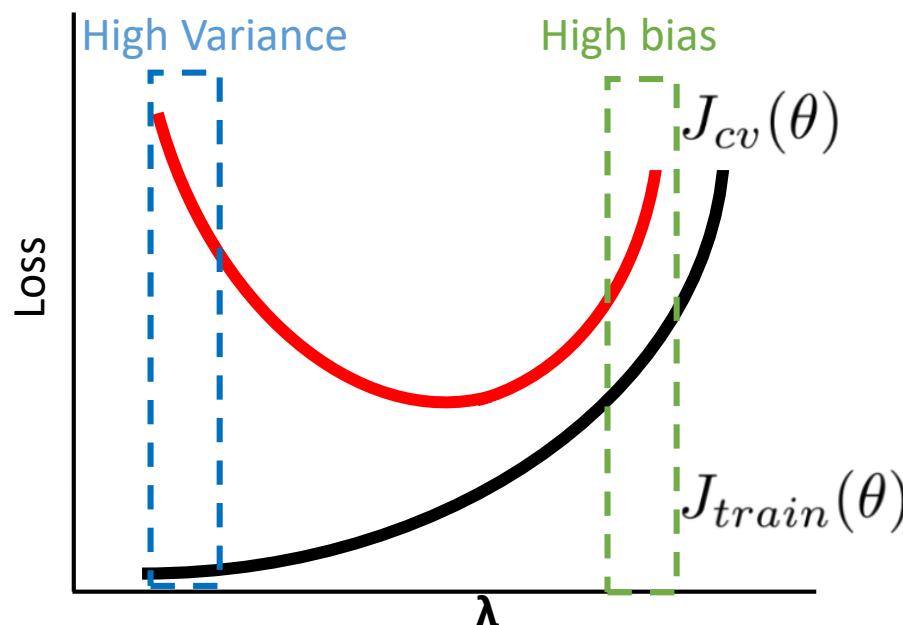
- Training error $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$
- Cross-validation error $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2 + \frac{\lambda}{2m_{cv}} \sum_{j=1}^{m_{cv}} \theta_j^2$



Bias / Variance Trade-off with Regularization

L2

- Training error $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$
- Cross-validation error $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2 + \frac{\lambda}{2m_{cv}} \sum_{j=1}^{m_{cv}} \theta_j^2$



Problem: Fail to Generalize

- Should we get more data?

Problem: Fail to Generalize

- Should we get more data?
- Getting more data does not always help

Problem: Fail to Generalize

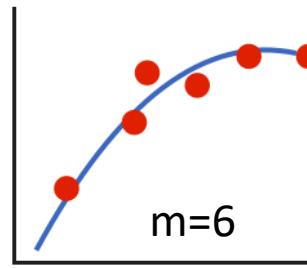
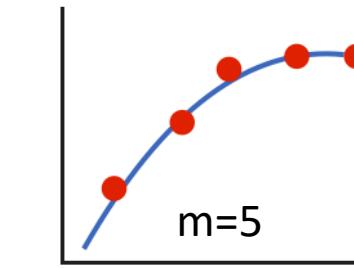
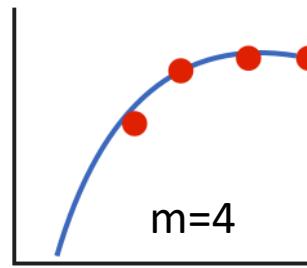
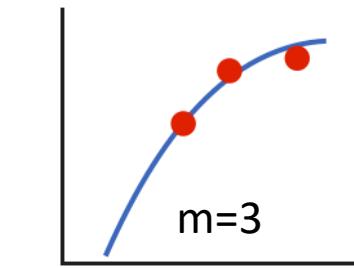
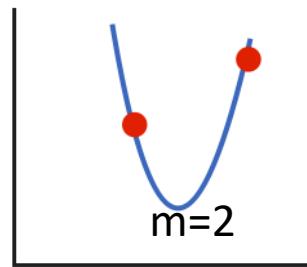
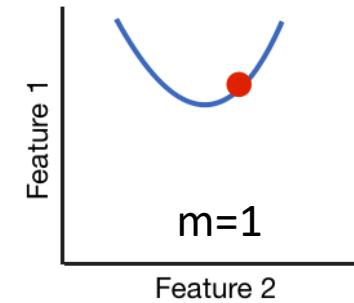
- Should we get more data?
- Getting more data does not always help
- How do we know if we should collect more data?

Learning Curve

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

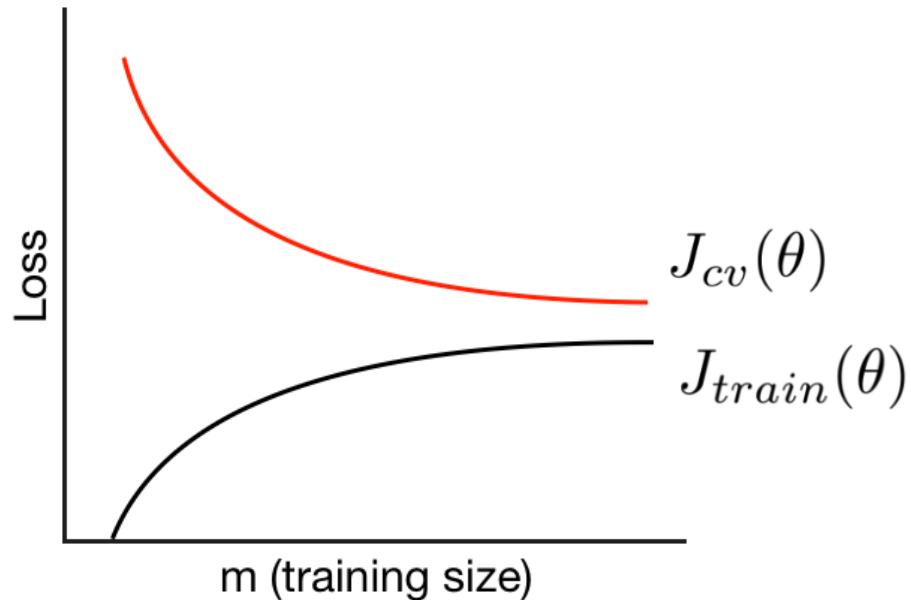
$$\theta_0 + \theta_1 x + \theta_2 x^2$$



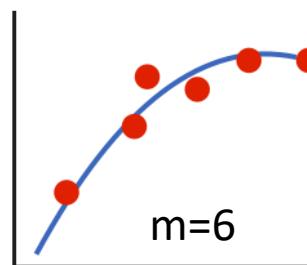
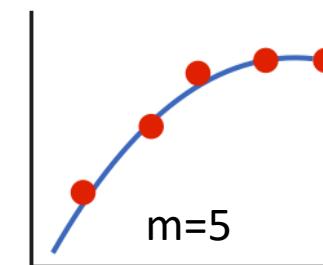
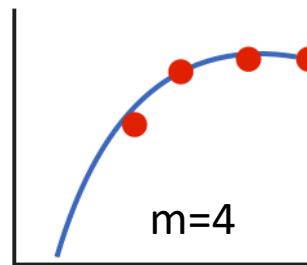
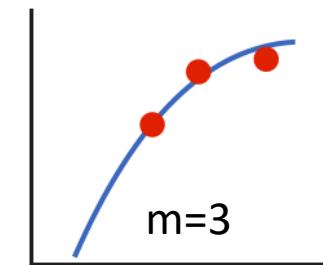
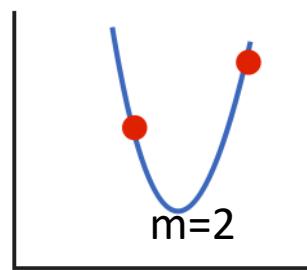
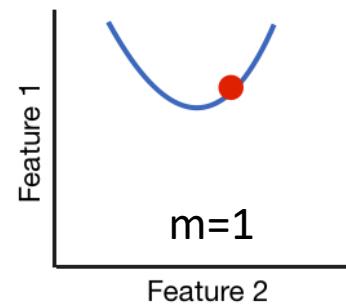
Learning Curve

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

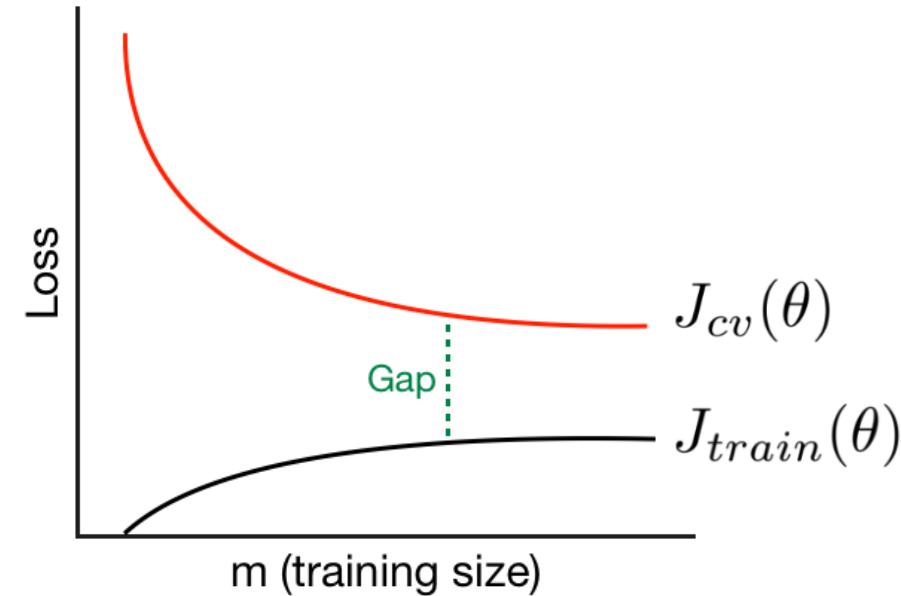
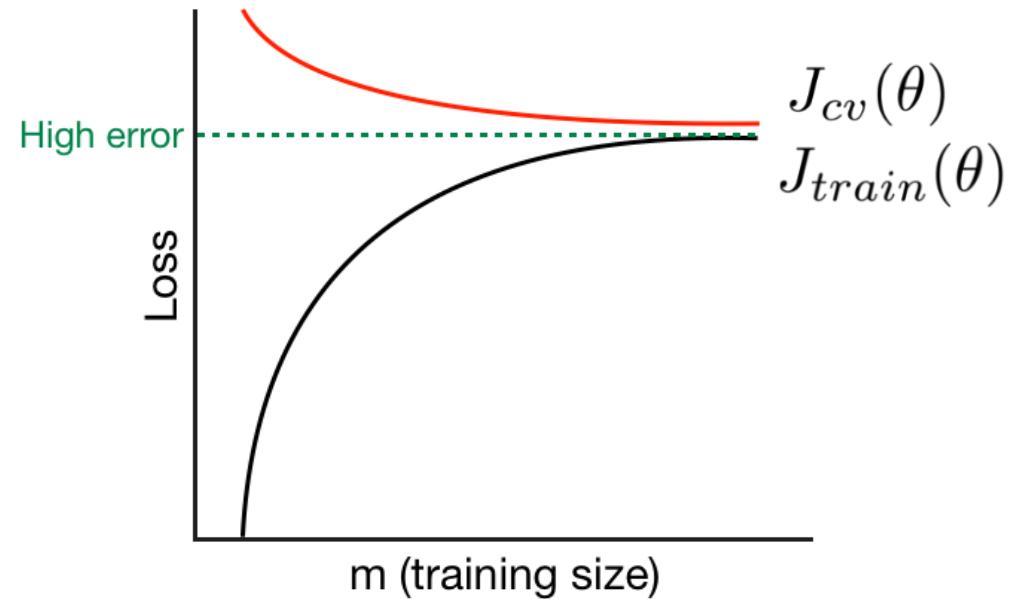
$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



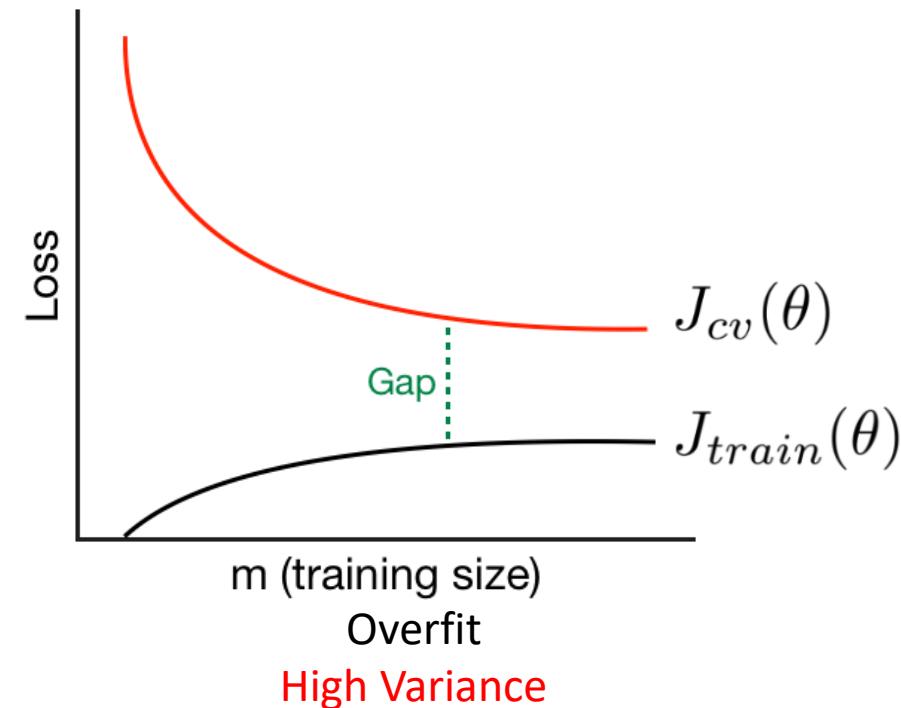
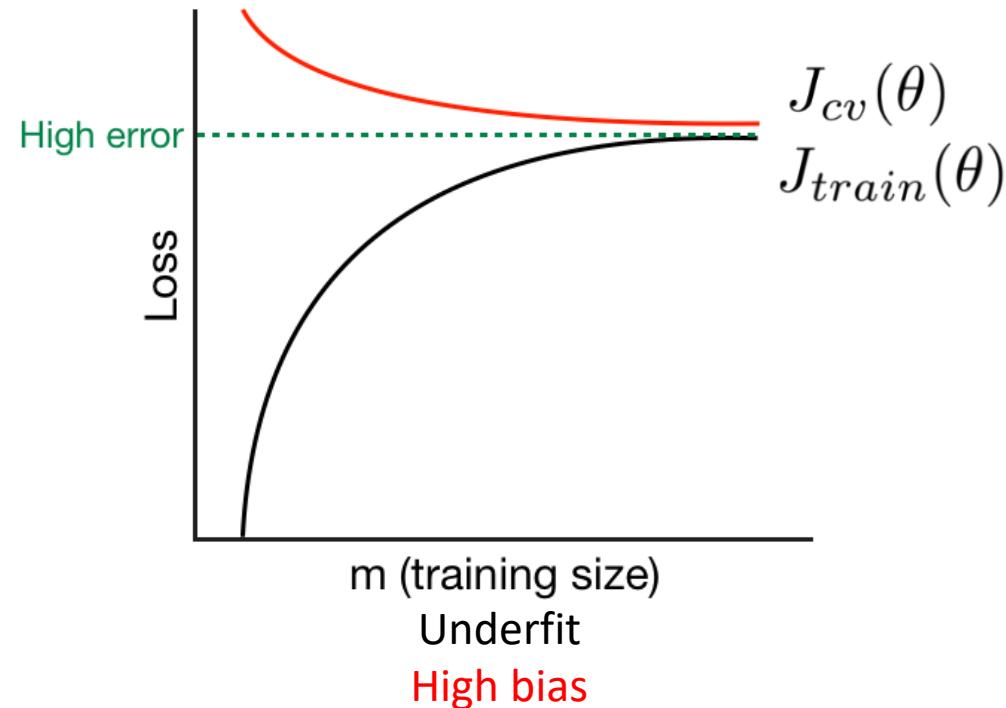
$$\theta_0 + \theta_1 x + \theta_2 x^2$$



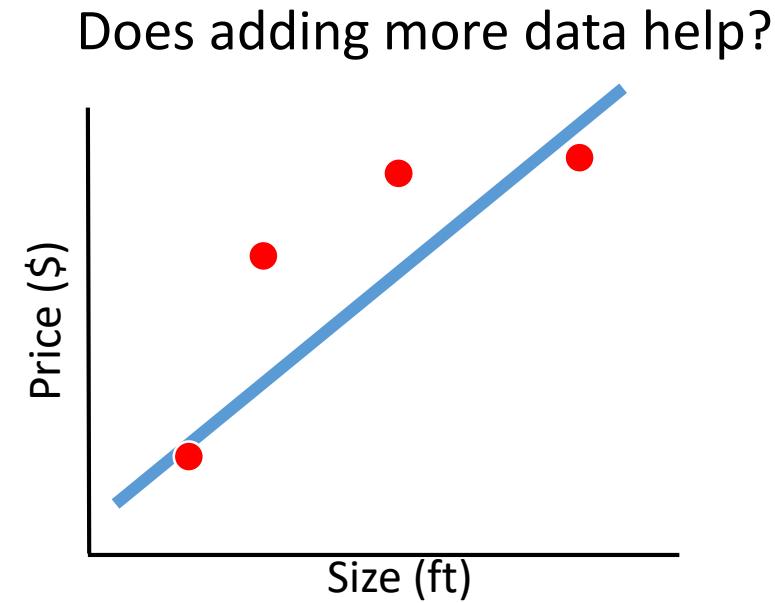
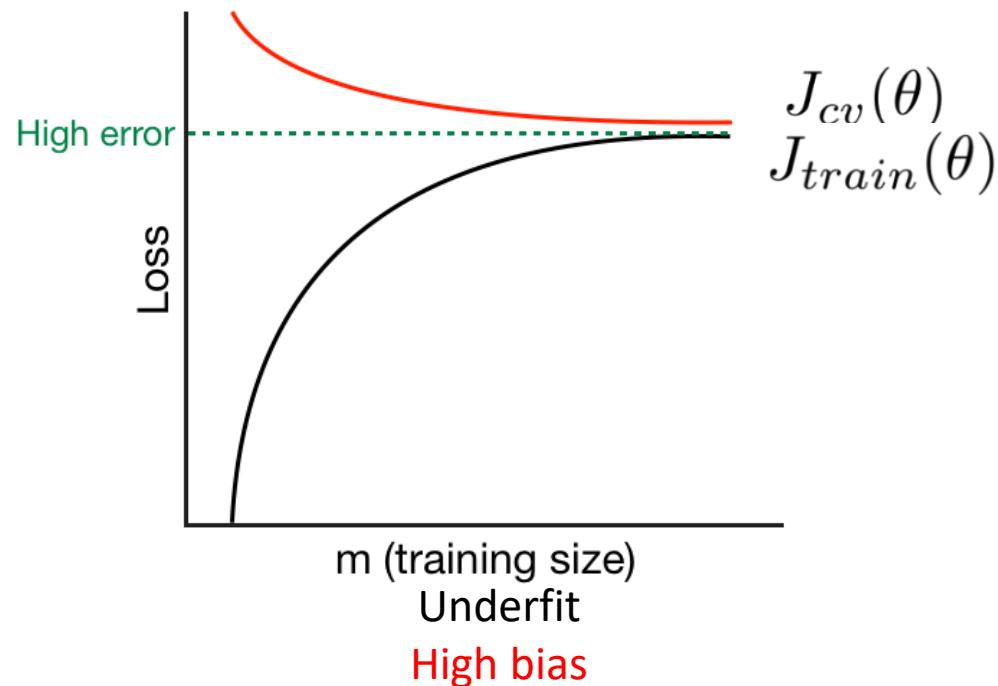
Learning Curve



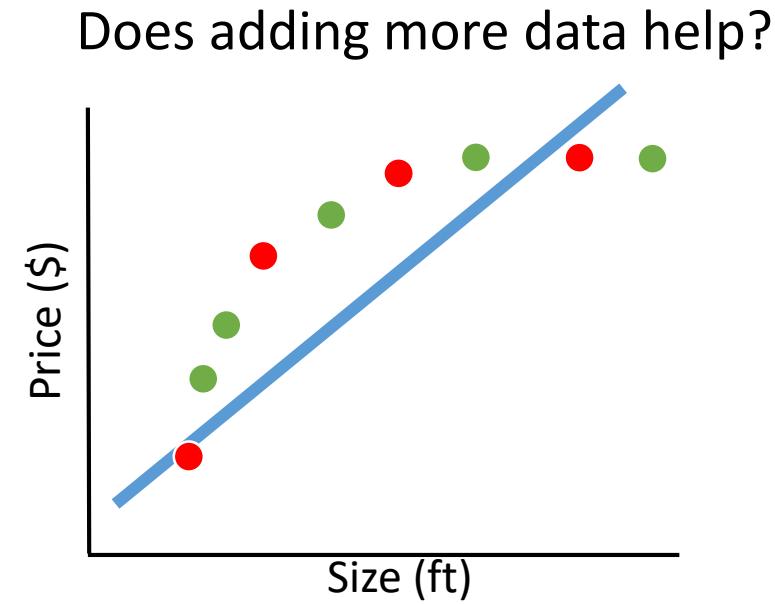
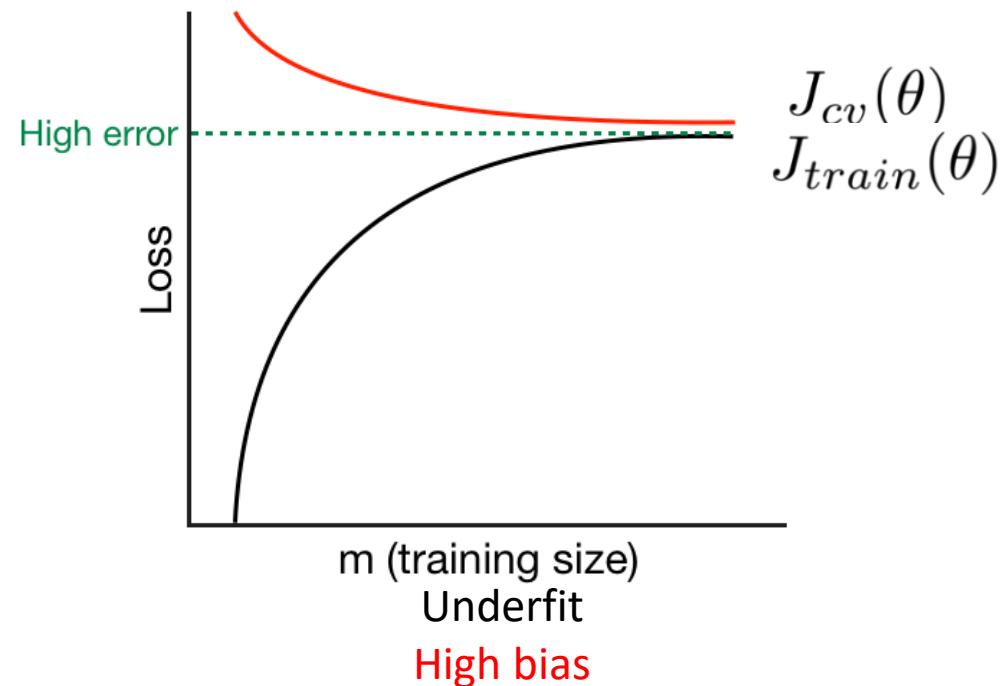
Learning Curve



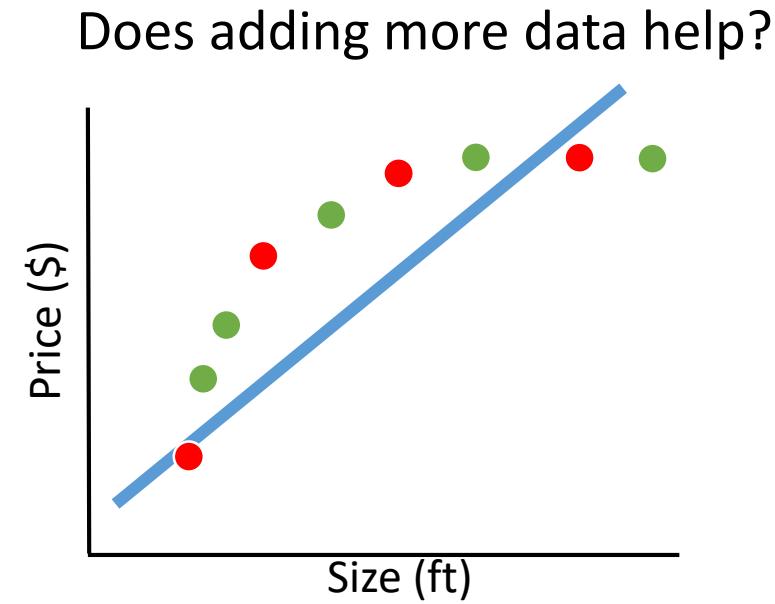
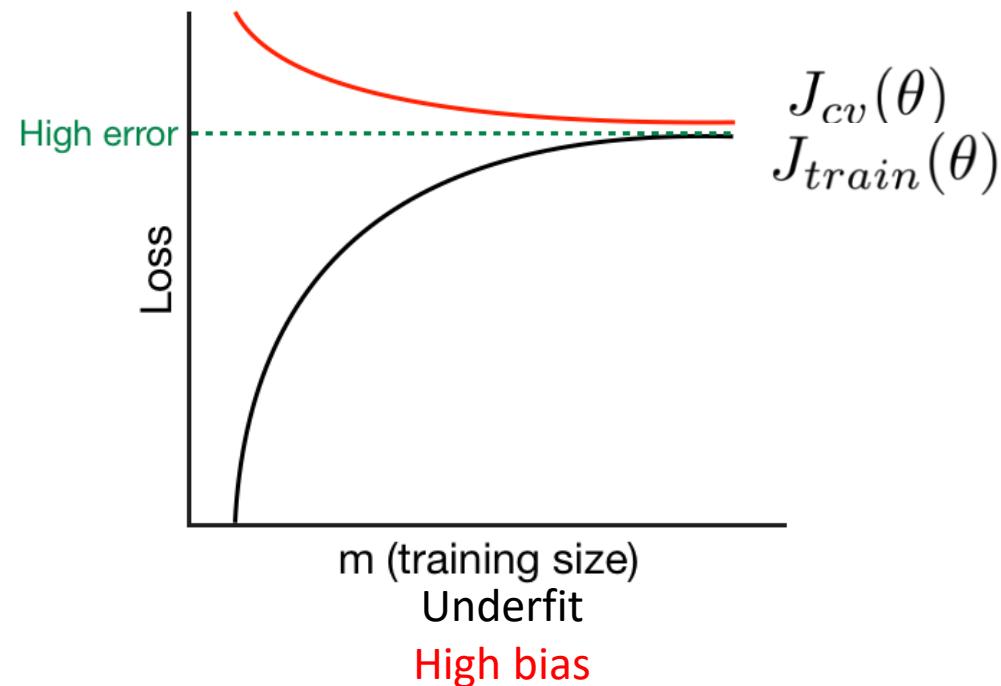
Learning Curve



Learning Curve

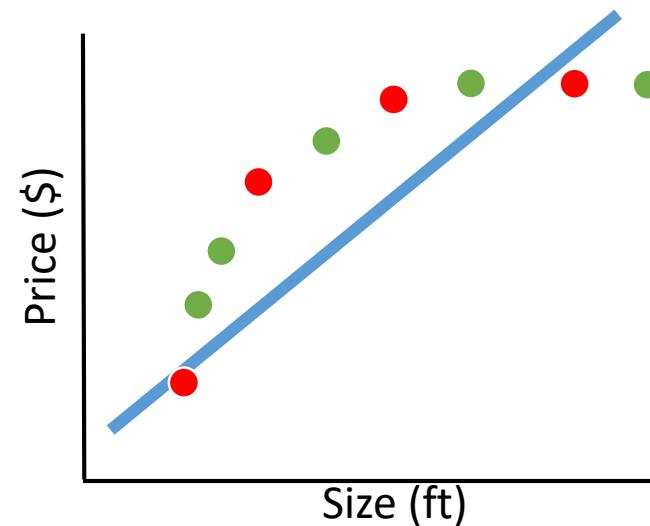
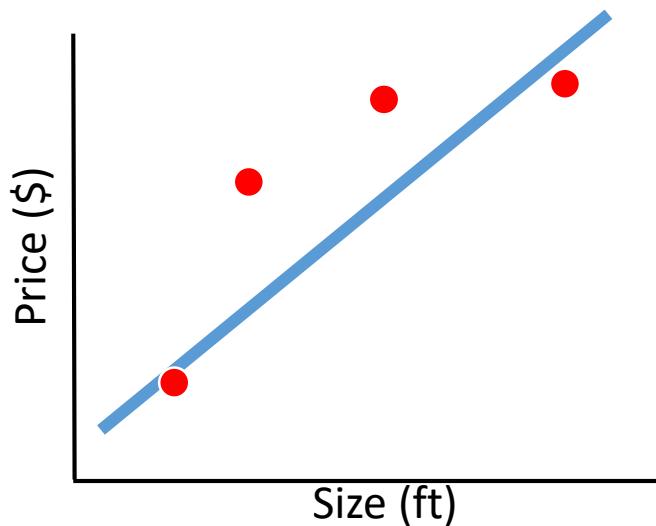


Learning Curve



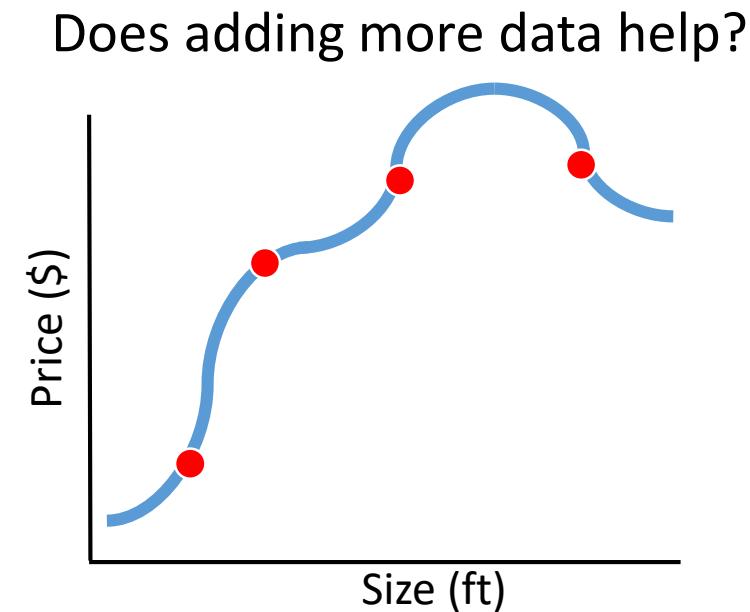
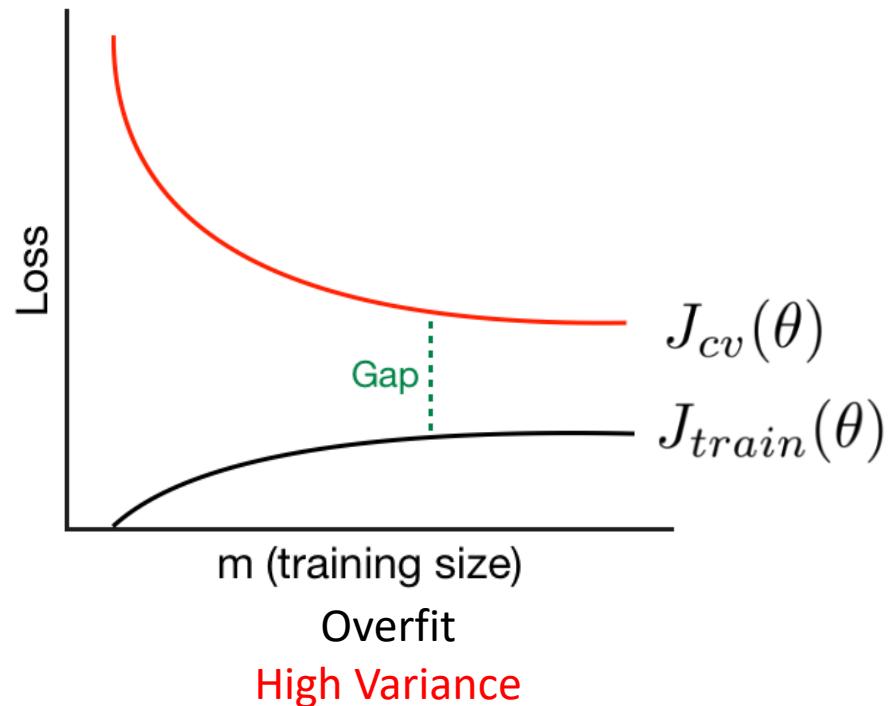
Learning Curve

Does adding more data help?

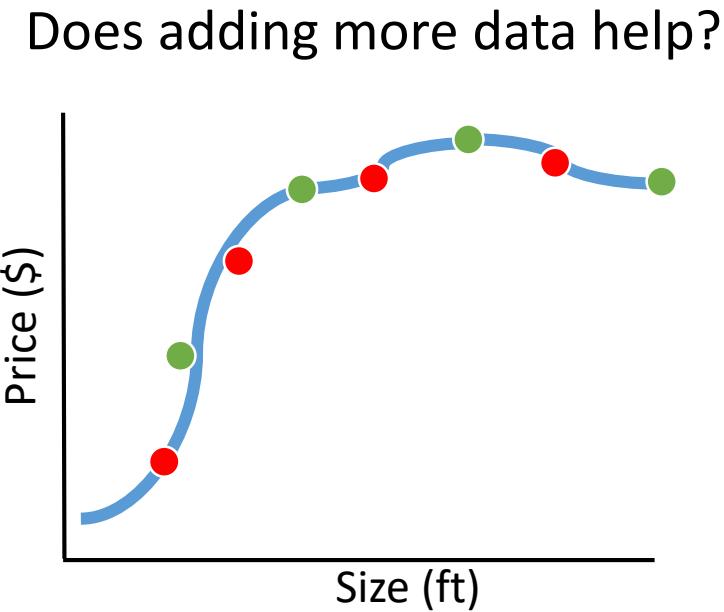
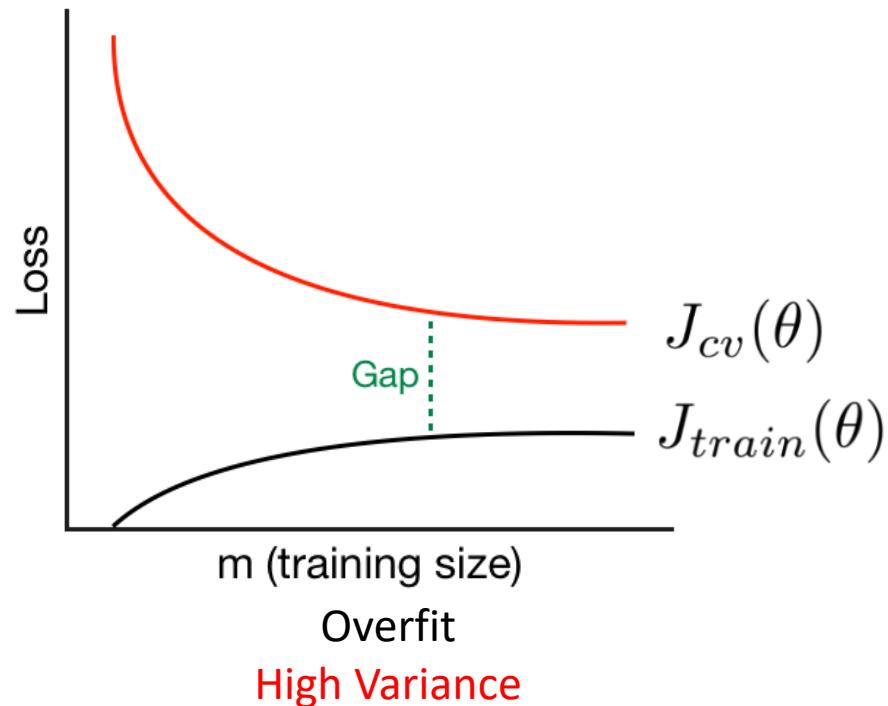


More data doesn't help when your model has **high bias**

Learning Curve

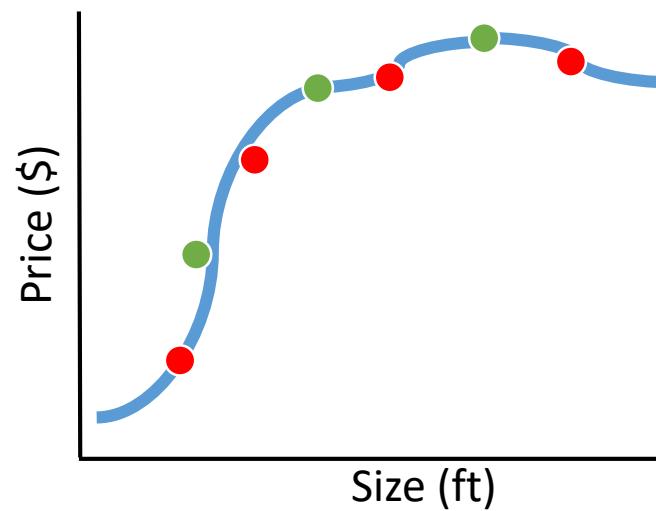
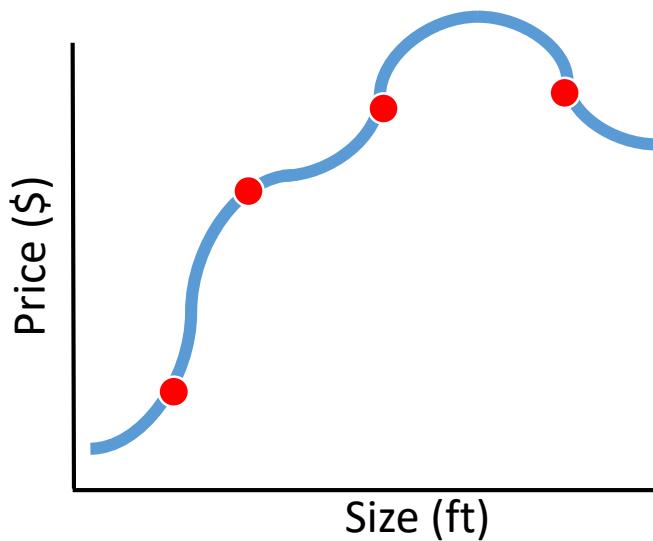


Learning Curve



Learning Curve

Does adding more data help?



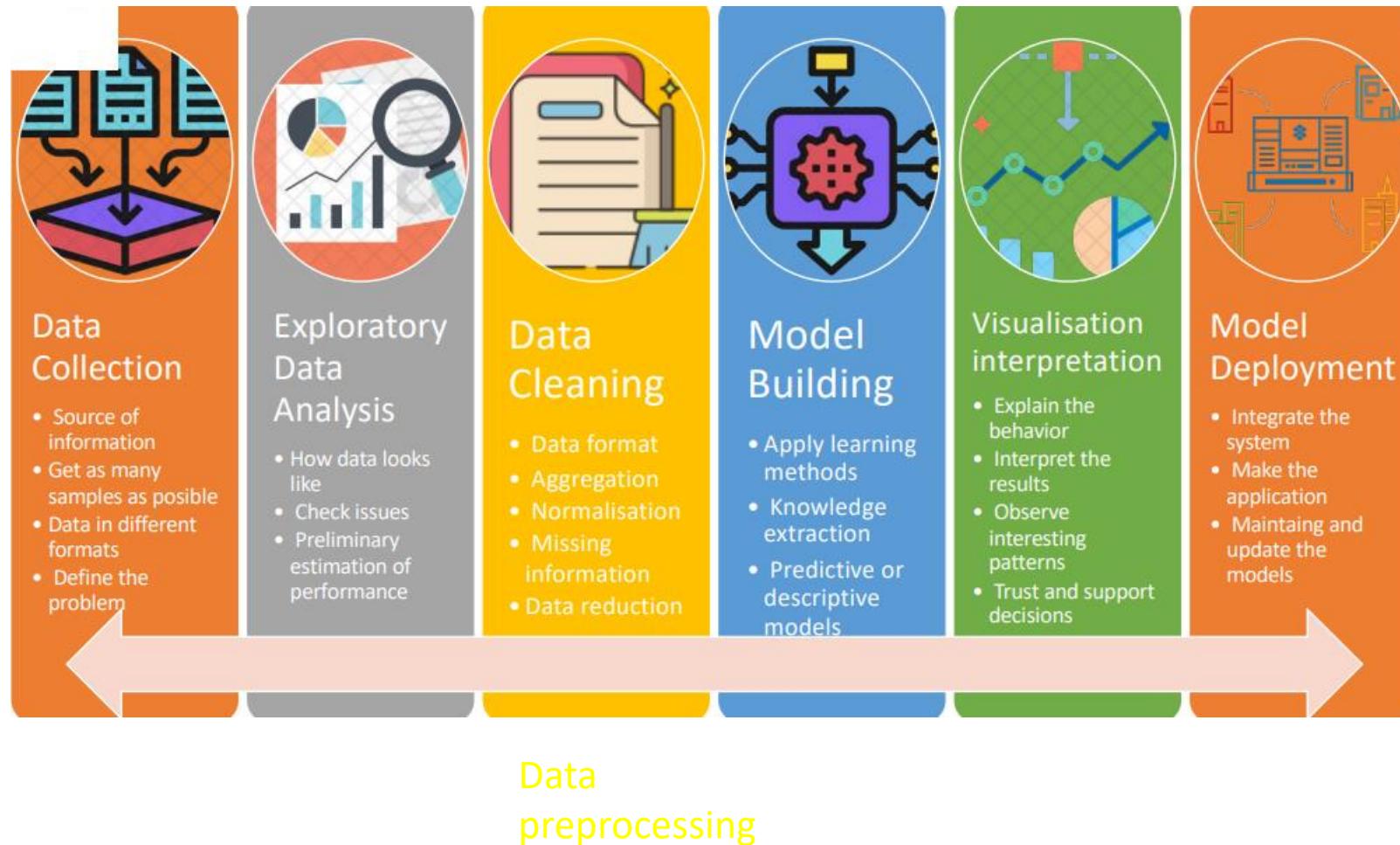
More data is **likely** to help when your model has **high variance**

Things You Can Try

- Get more data
 - When you have **high variance**
- Try different features
 - Adding feature helps fix **high bias**
 - Using smaller sets of feature fix **high variance**
- Try tuning your decay parameter
 - Decrease regularization when **bias is high**
 - Increase regularization when **variance is high**

Analyze your model before you act

ML pipeline



Data Preprocessing

What does Data Preprocessing include?

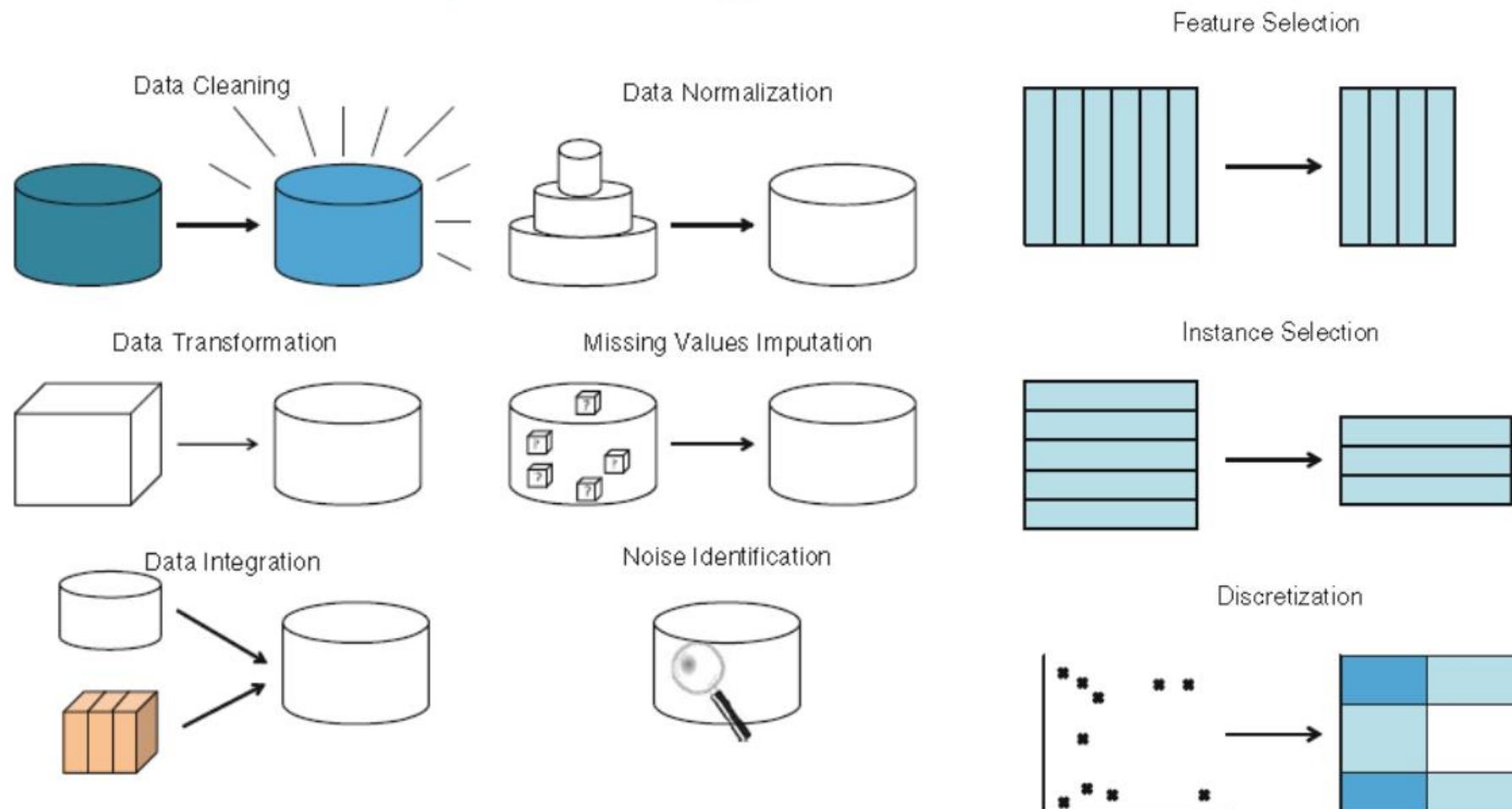
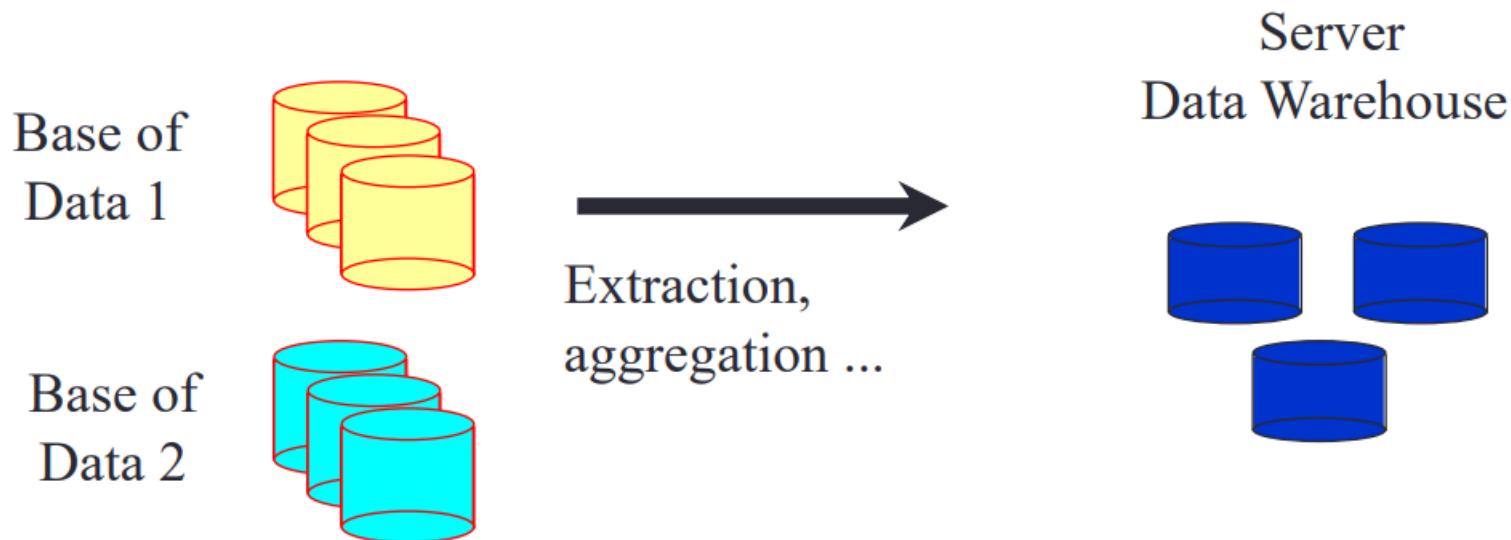


Fig. 1.3 Forms of data preparation

data reduction

Data Integration

- ✿ Obtains data from different sources of information
- ✿ Solves representation and coding problems
- ✿ Integrates data from different tables to create homogeneous information, ...



Data cleansing

- **Objectives:**
 - resolving inconsistencies
 - Fill/impute missing values,
 - smoothing the noise in the data,
 - identify or eliminate *outliers* ...

Data cleansing

- Data need to be formatted for a given software tool
- Data need to be made adequate for a given method
- Data in the real world is dirty
 - incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - e.g., occupation=""
 - noisy: containing errors or outliers
 - e.g., Salary="-10", Age="222"
 - inconsistent: containing discrepancies in codes or names
 - e.g., Age="42" Birthday="03/07/1997"
 - e.g., Was rating "1,2,3", now rating "A, B, C"
 - e.g., discrepancy between duplicate records

Normalisation (MinMaxScaling)

- **Objective:** To move the values of an attribute to a better rank.
- Useful for some techniques such as AANN or distance-based methods (nearest neighbours,...).
- Some standardisation techniques:
 - **min-max normalisation:** Performs a linear transformation of the original data.

$$\begin{aligned} [\min_A, \max_A] &\rightarrow [nuevo_{\min_A}, nuevo_{\max_A}] \\ v' &= \frac{v - \min_A}{\max_A - \min_A} (nuevo_{\max_A} - nuevo_{\min_A}) + nuevo_{\min_A} \end{aligned}$$

Relationships between the original data are preserved.

Standardisation

- **Zero-mean normalisation.** It is normalised according to the mean and the standard deviation.

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

Useful when limits are unknown or when outliers may dominate the min-max normalisation.

- **Normalisation by decimal scale.** Normalises by moving the decimal point of the attribute values. The number of decimal points moved depends on the maximum absolute value of A.

$$v' = \frac{v}{10^j}$$

with j equal to the smallest integer such that $\max(|v'|) < 1$.

E.g. : if the data are in [-986,917], then $j=3$.

Data Transformation

- **Objective:** To transform the data in the best possible way for the application of DM algorithms.
- Some typical operations:
 - Aggregation. E.g. aggregation of monthly sales into a single attribute that is annual sales, ...
 - Data generalisation. The aim is to obtain higher level data from current data, using hierarchies of concepts.
 - City streets
 - Numerical age {young, adult, middle-aged, old}
 - Normalisation: Change the range [-1,1] or [0,1].
 - Linear, quadratic, polynomial transformations, ...

Different examples of variables or features

- ID numbers, Names of people
- eye color, zip codes
- rankings (e.g., taste of potato chips on a scale from 1-10),
grades, height in {tall, medium, short}
- calendar dates, temperatures in Celsius or Fahrenheit, GRE
(Graduate Record Examination) and IQ scores
- temperature in Kelvin, length, time, counts

Different examples of features

- Examples:
 - US State Code (50 values)
 - Profession Code (7,000 values, but only few frequent)
- Ignore ID-like fields whose values are unique for each record
- For other fields, group values “naturally”:
 - e.g. 50 US States = 3 or 5 regions
 - Profession – select most frequent ones, group the rest

Encoding

- the machine learning model can't process categorical variables.
- So when we have categorical variables in our dataset then we need to convert them into numerical variables.
- There are many ways to convert categorical values into numerical values. Most popular
 - One hot encoding
 - Label encoding

One hot encoding

Country	India	Australia	Russia	America
India	1	0	0	0
Australia	0	1	0	0
Russia	0	0	1	0
America	0	0	0	1

Label encoding

State	Confirmed	Recovered	Deaths
Punjab	2000	1500	200
Haryana	2321	1222	345
Kerala	3455	2365	400

