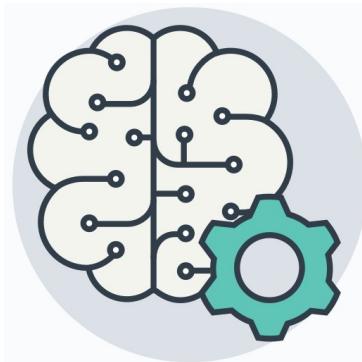


Aprendizado de Máquina

Introdução



Prof. Regis Pires Magalhães

regismagalhaes@ufc.br - <http://bit.ly/ufcregis>

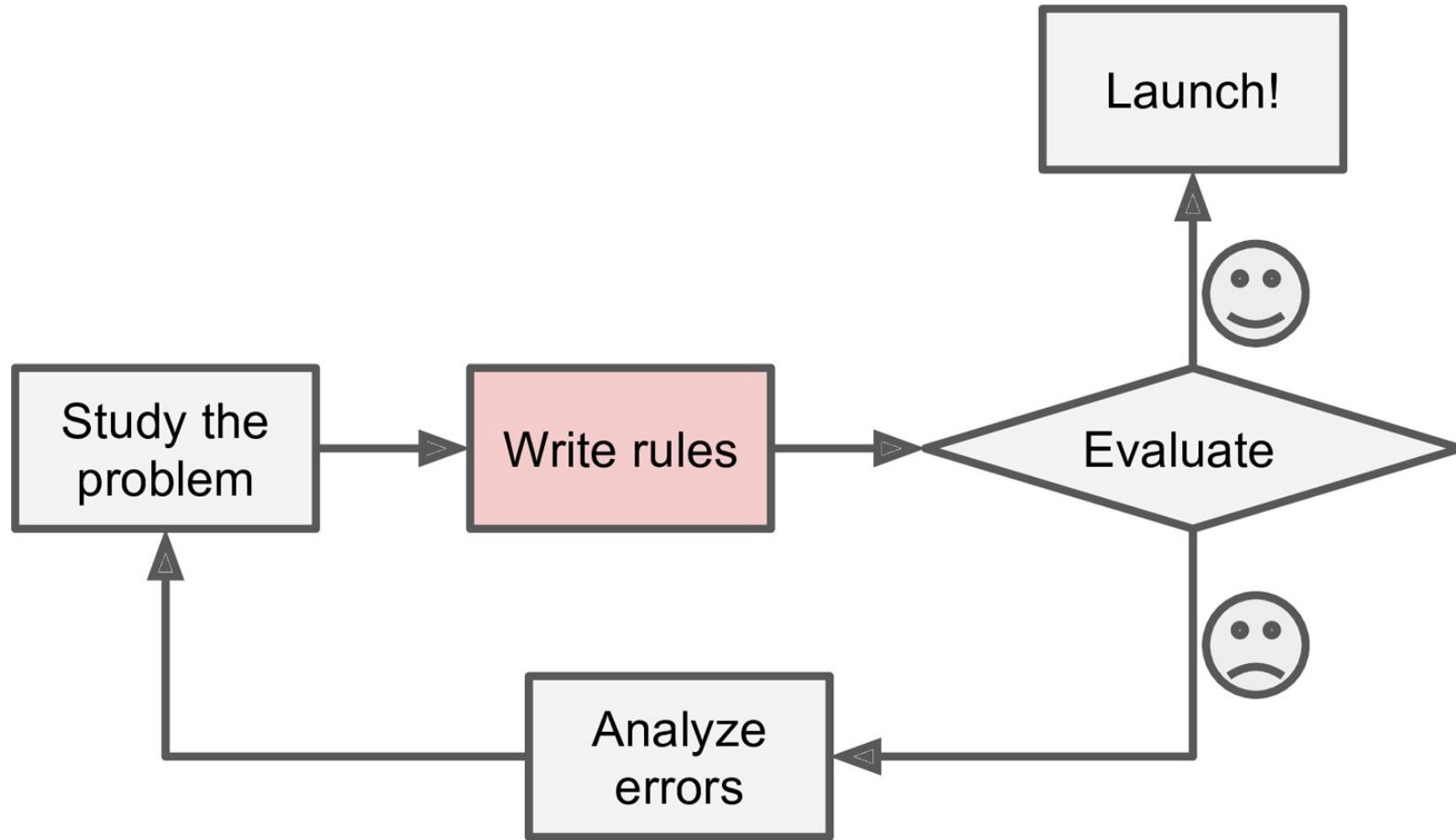
What is Machine Learning?

- Machine Learning is the science (and art) of programming computers so they can learn from data.
- [Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed.
—Arthur Samuel, 1959
- A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.
—Tom Mitchell, 1997

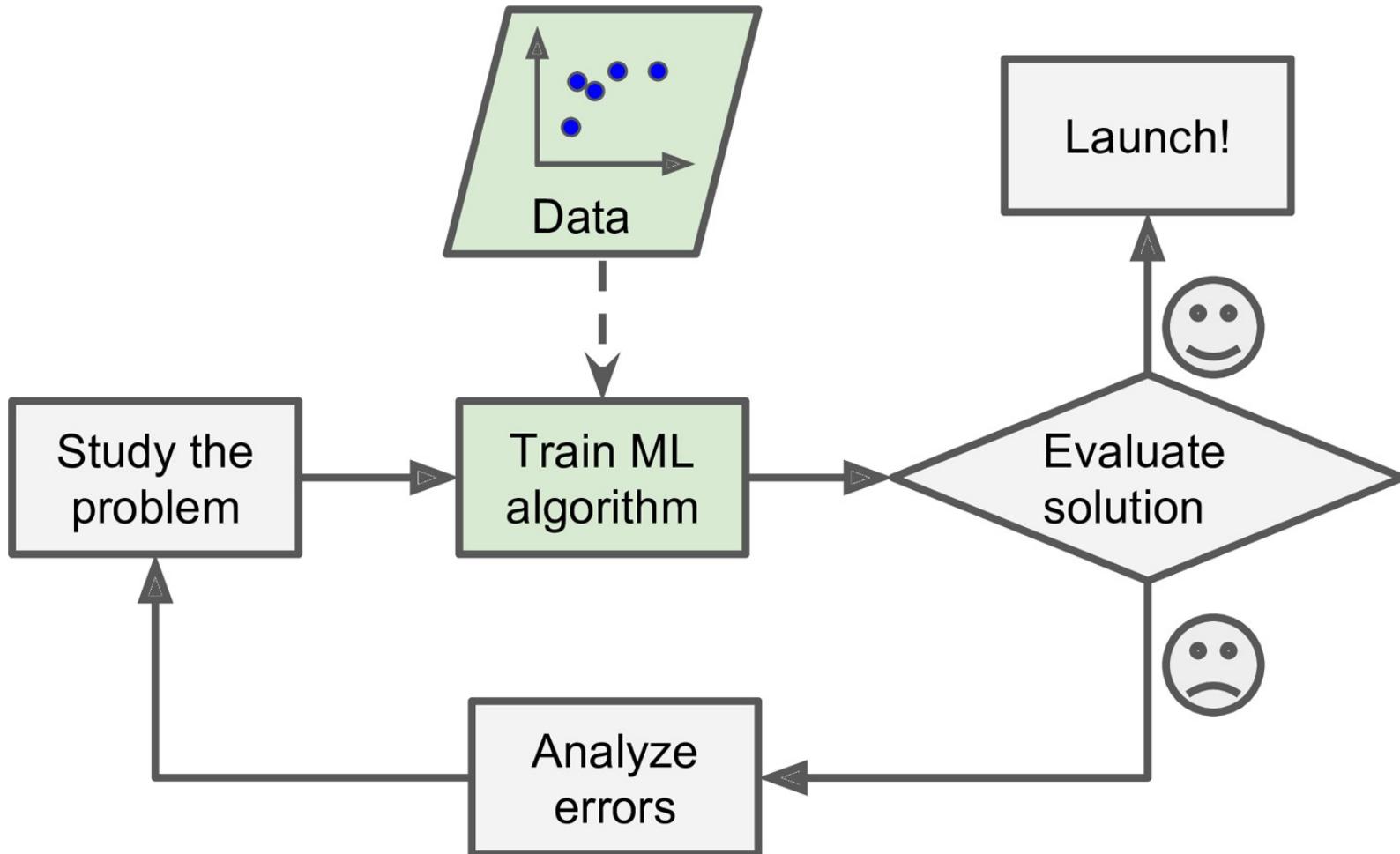
Giving Computers the Ability to Learn from Data

- Intuição:
 - Exemplo do empréstimo pessoal

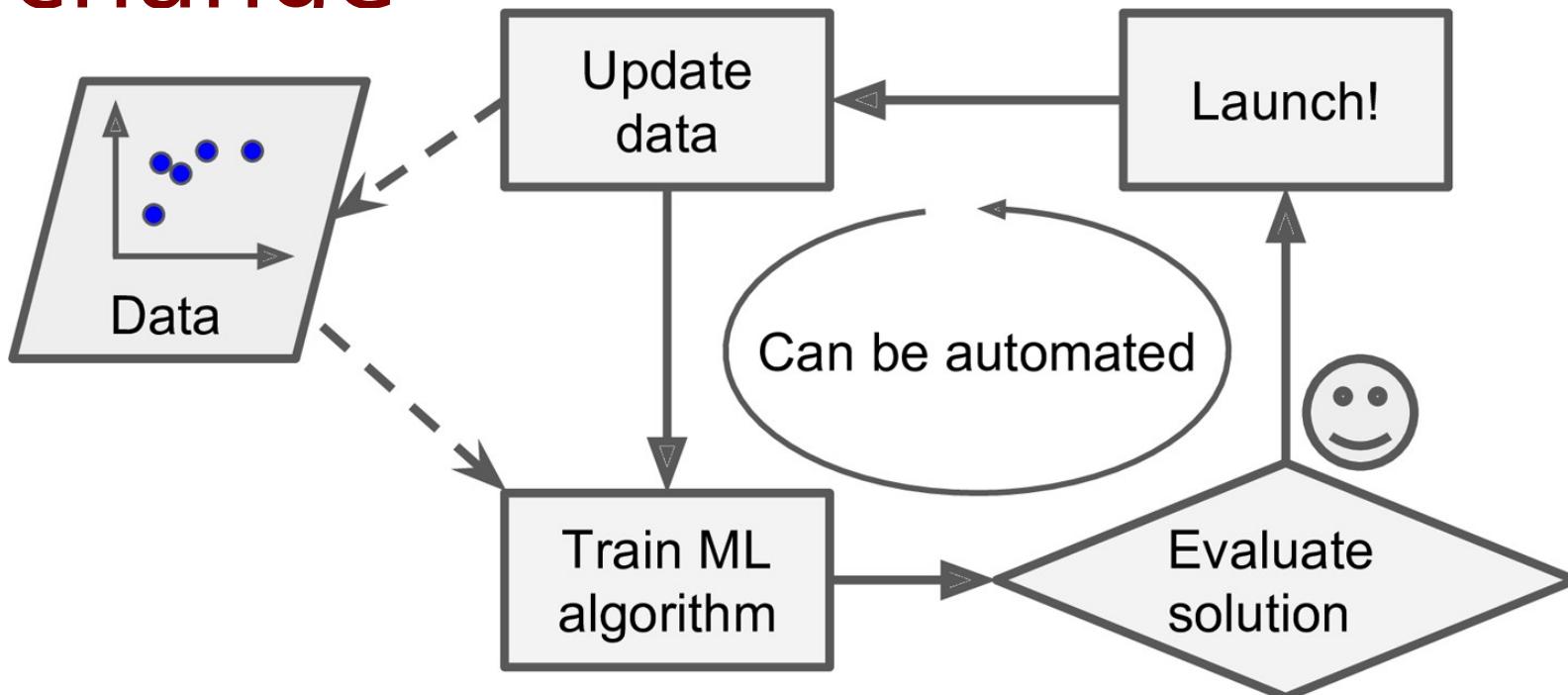
Traditional approach



Machine Learning approach



Automatically adapting to change



ARTIFICIAL INTELLIGENCE

A program that can sense, reason,
act, and adapt

MACHINE LEARNING

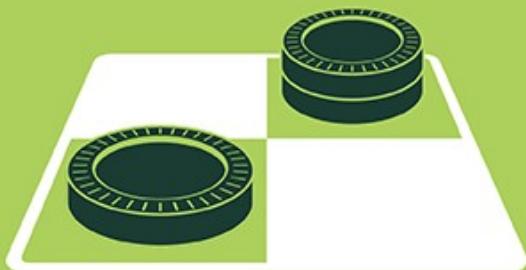
Algorithms whose performance improve
as they are exposed to more data over time

DEEP LEARNING

Subset of machine learning in
which multilayered neural
networks learn from
vast amounts of data

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

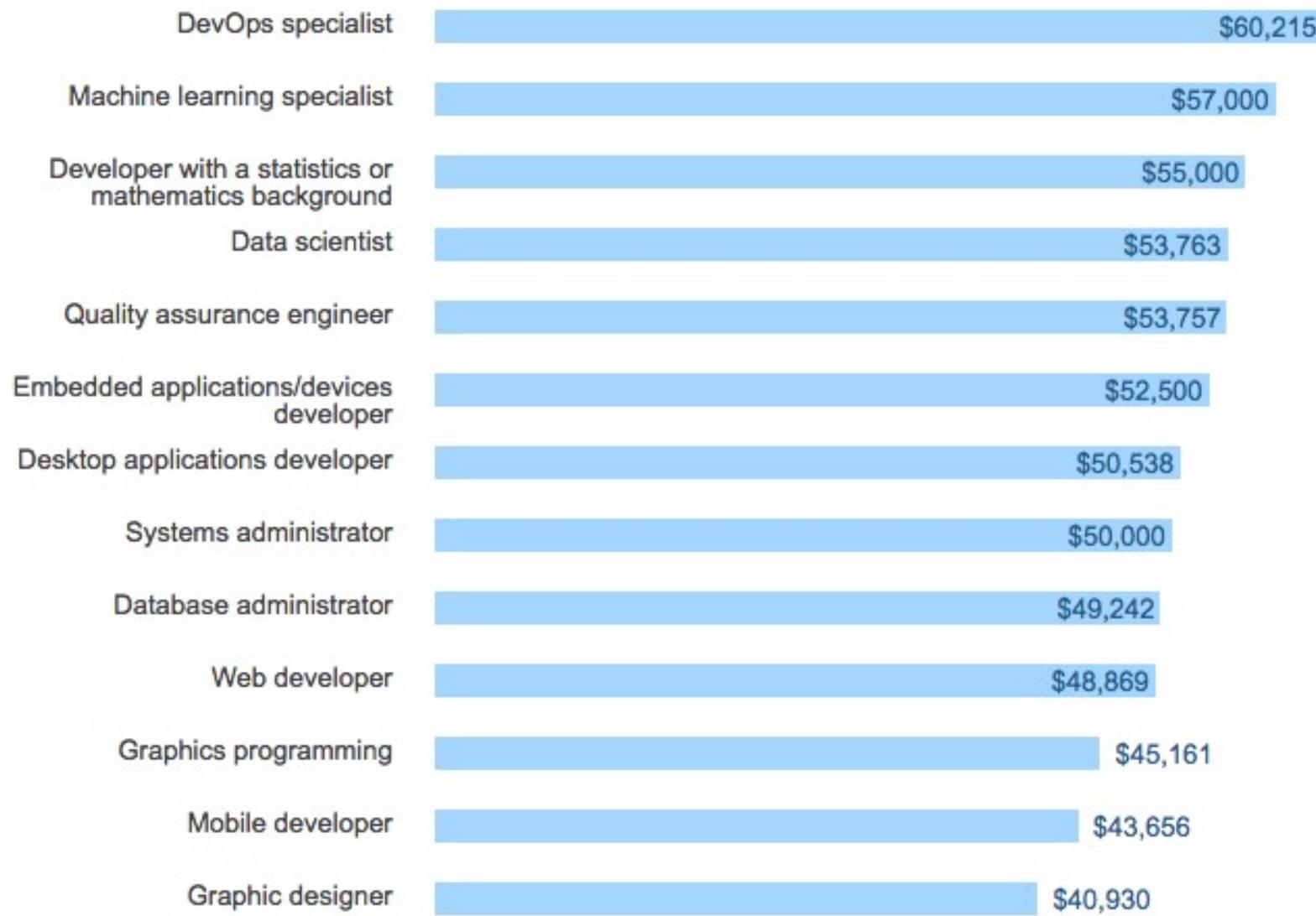
Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Motivações

- Barateamento e popularização dos dispositivos móveis e dos mais diversos tipos de sensores para coleta de dados.
- Aumento do poder de processamento e armazenamento dos computadores.
- Otimização de recursos em um mercado global cada vez mais competitivo.
- Amadurecimento, melhoria e disseminação de técnicas de aprendizagem de máquina.



Salary by Developer Type



Median of 12,475 responses; USD

<https://stackoverflow.com/insights/survey/2017/>

Salaries by Geography

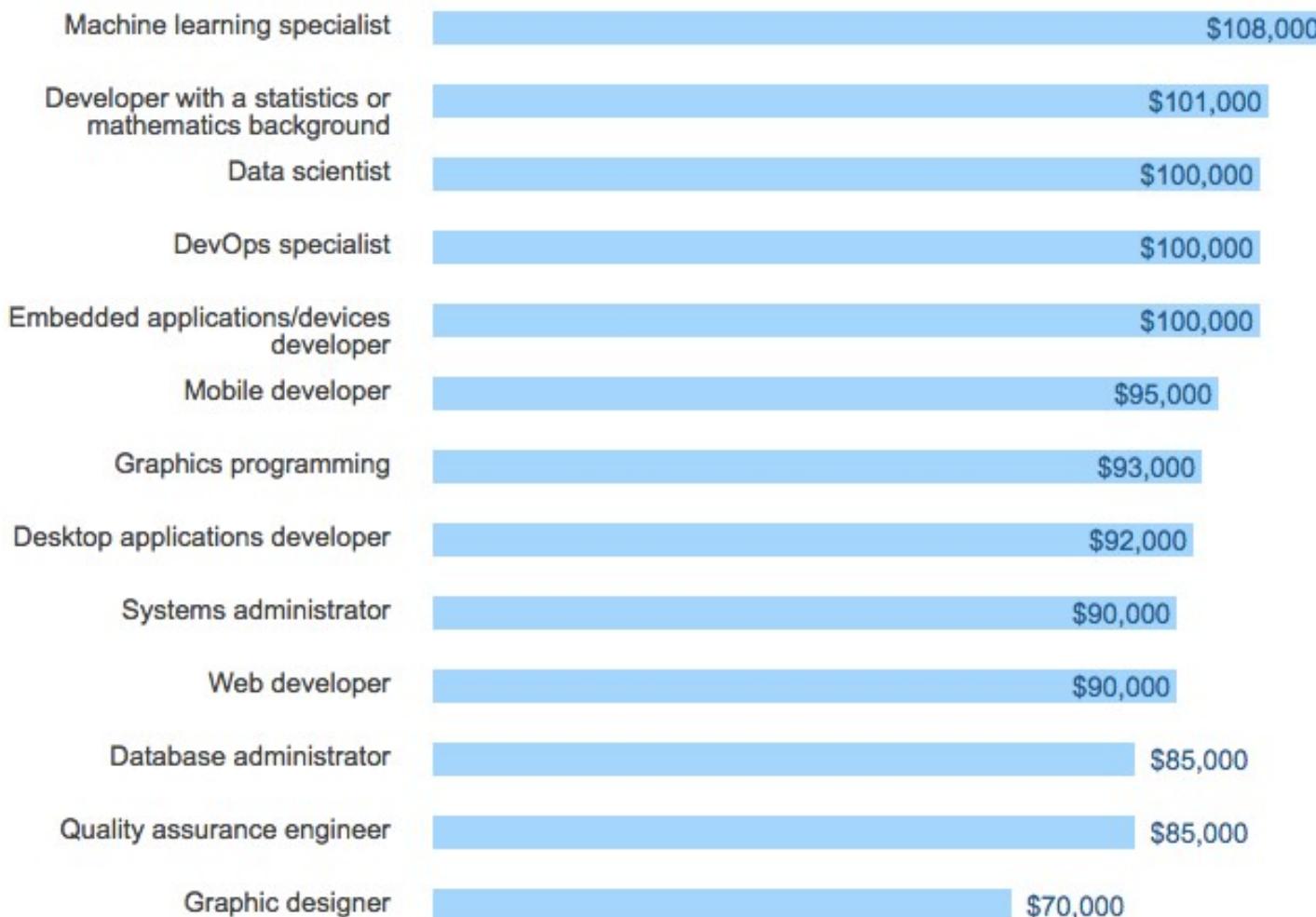
United States

Canada

United Kingdom

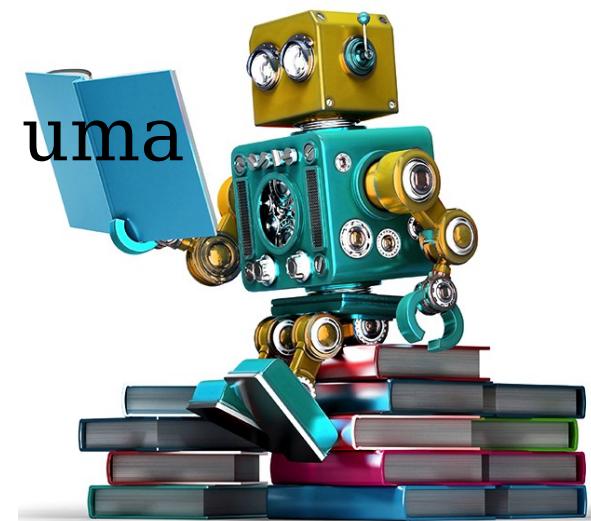
Germany

France



Exemplos de uso de técnicas de aprendizado de máquina

- Risco de empréstimo / seguro / plano de saúde
- Preços de imóveis / carros
- Produtos com maior chance de venda
- Sistemas de recomendação
- Evasão de cursos
- Probabilidade de desenvolver uma doença
- Sistema de busca
- Sistemas de segurança
- etc.



Top 10 Use Cases for Data Science & Machine Learning



HEALTHCARE:
Patient Diagnosis



FINANCE:
Fraud Detection



MANUFACTURING:
Anomaly Detection



RETAIL:
Inventory Optimization



GOVERNMENT:
Smarter Services



TRANSPORTATION:
Demand Forecasting



NETWORKS:
Intrusion Detection



E-COMMERCE:
Recommender Systems



MEDIA:
Interaction & Speed



EDUCATION:
Research Insight



<https://serenatadeamor.org/>

<https://medium.com/data-science-brigade>

Brazilian group develops an AI to help in public expenditures monitoring. Rosie, the robot's name, found more than 8.000 suspicious reimbursements from Brazilian congresspeople.

Despite all the corruption related news coming from Brazil, there is a movement for transparency in the country. Several bills signed in the last years put Brazil in the top of transparency rankings worldwide, specially when our former president Dilma Rousseff signed in 2011 the Access of Information Law, a Brazilian version for the american FOIA (Freedom of Information Act), which completed 50 years in 2016. It makes open data compulsory for all public bodies. Something similar has happened in some



A Operação Serenata de Amor criou a **Rosie** - uma inteligência artificial capaz de analisar cada pedido de reembolso dos deputados e identificar a probabilidade de ilegalidade.



R\$ 6.205,00

É o valor de uma nota que foi reembolsada, referente a

UMA REFEIÇÃO

30

**TANQUES DE
GASOLINA
COMPLETOS**

Temos um deputado que costuma gastar

R\$ 6.000,00

mensais em gasolina. Em média,
30 tanques por mês.



— 13

Dois deputados já pediram reembolso de 13 refeições feitas no mesmo dia

**2
1
9**

90%



42%

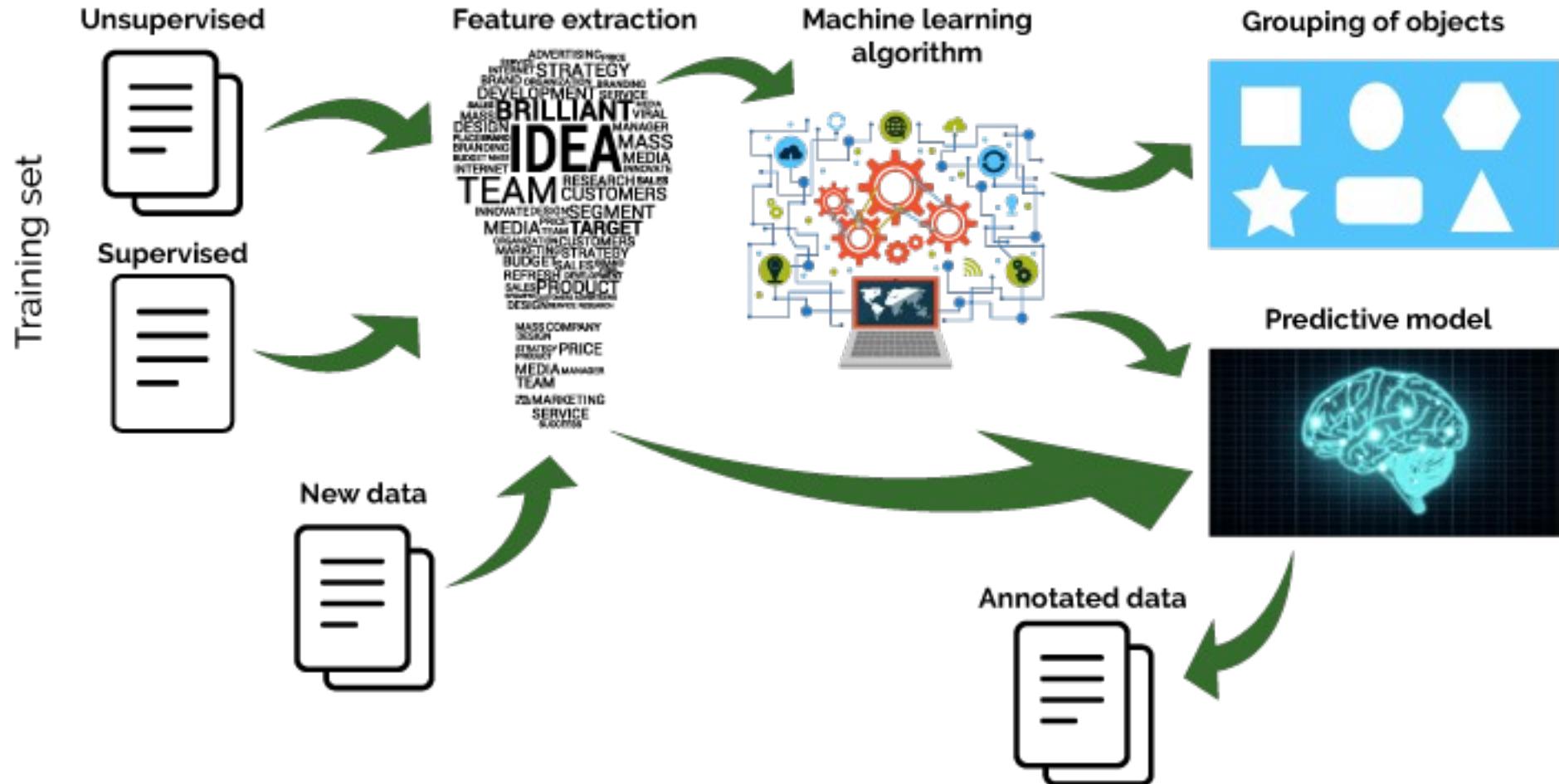
Um deputado já foi reembolsado por
**BEBIDA ALCÓOLICA
EM LAS VEGAS**



DEPUTADOS

Costumam a usar
**O VALOR MÁXIMO
PERMITIDO
MENSALMENTE**

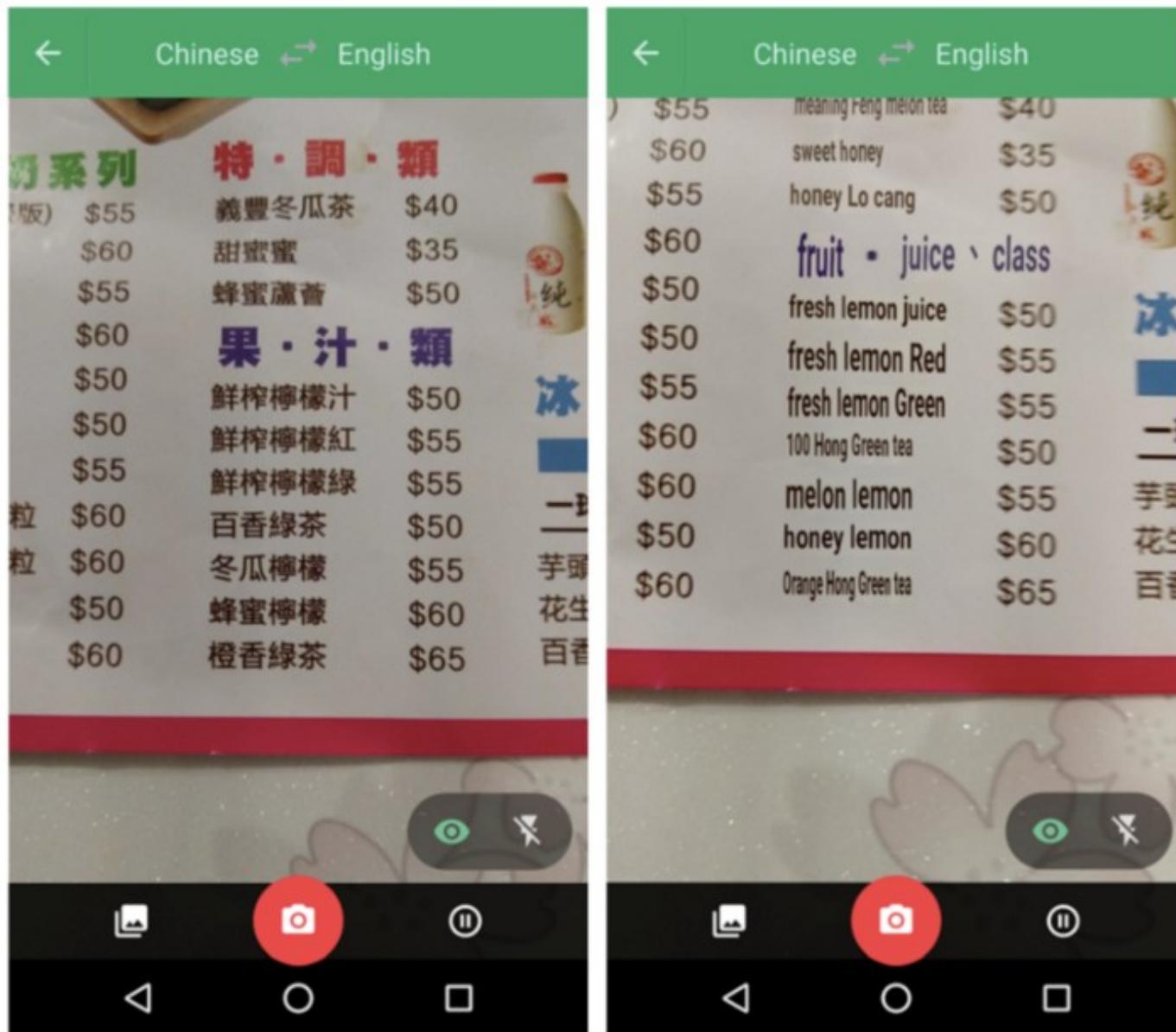
Machine Learning





See the full match at The International 2017, with Dendi (human) vs. OpenAI (bot), on [YouTube](#).

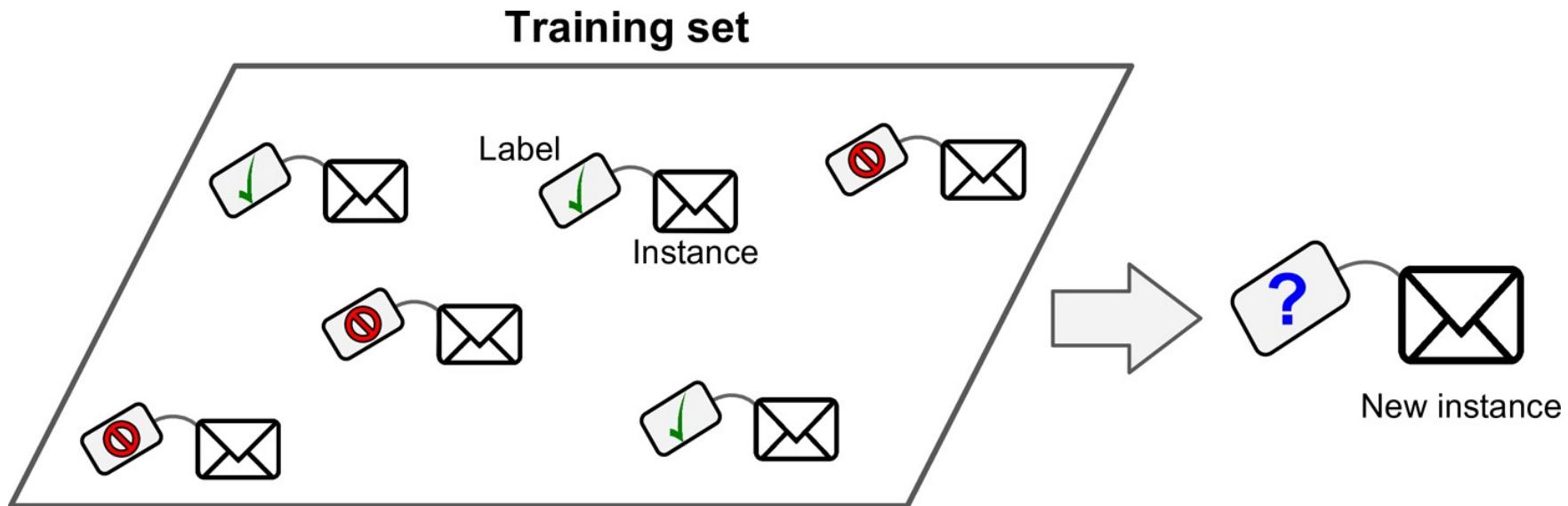
<https://www.youtube.com/watch?v=wiOopO9jTzW>



Google Translate overlaying English translations on a drink menu in real time using convolutional neural networks.

Usos de Machine Learning

Spam classification



Spam Detection

- data from 4601 emails sent to an individual (named George, at HP labs, before 2000). Each is labeled as **spam** or **email**.
- goal: build a customized spam filter.
- input features: relative frequencies of 57 of the most commonly

	george	you	hp	free	!	edu	remove
spam	0.00	2.26	0.02	0.52	0.51	0.01	0.28
email	1.27	1.27	0.90	0.07	0.11	0.29	0.01

Average percentage of words or characters in an email message equal to the indicated word or character. We have chosen the words and characters showing the largest difference between spam and email.

Identify the risk factors for prostate cancer



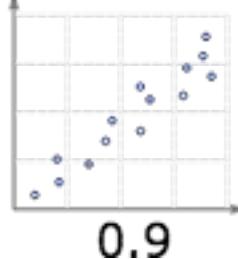
Scatter plot correlation

Correlação

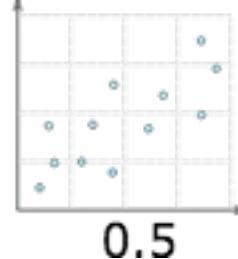
Perfect Positive Correlation



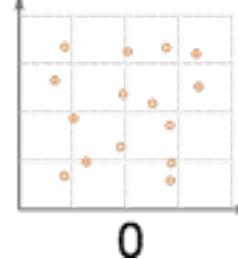
High Positive Correlation



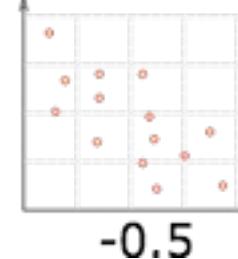
Low Positive Correlation



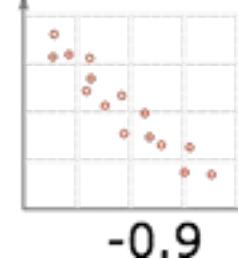
No Correlation



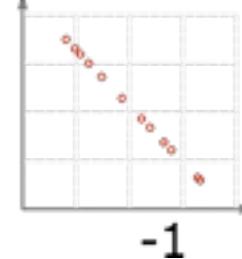
Low Negative Correlation



High Negative Correlation

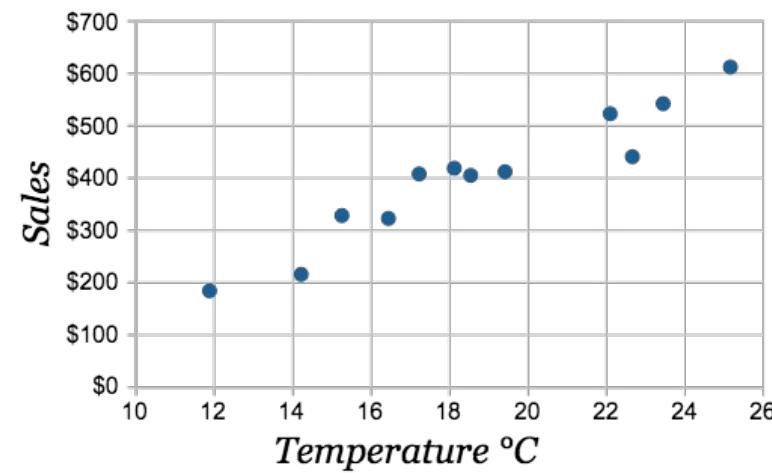


Perfect Negative Correlation



Ice Cream Sales vs Temperature

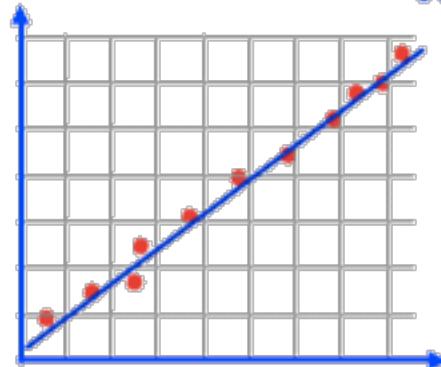
Temperature °C	Ice Cream Sales
14.2°	\$215
16.4°	\$325
11.9°	\$185
15.2°	\$332
18.5°	\$406
22.1°	\$522
19.4°	\$412
25.1°	\$614
23.4°	\$544
18.1°	\$421
22.6°	\$445
17.2°	\$408



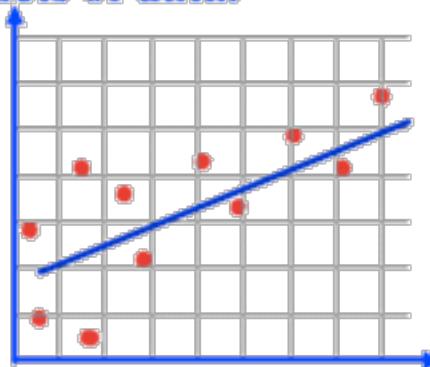
the correlation is 0.9575

SCATTERPLOTS & CORRELATION

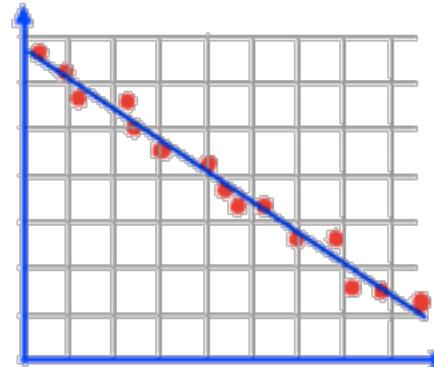
Correlation - indicates a relationship (connection) between two sets of data.



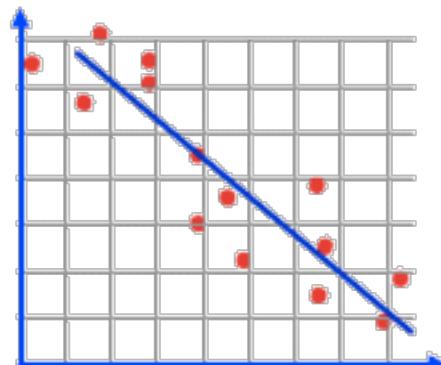
Strong positive correlation



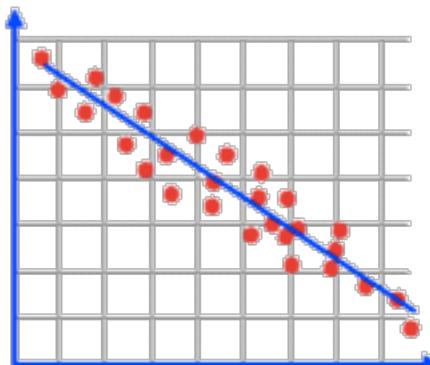
Weak positive correlation



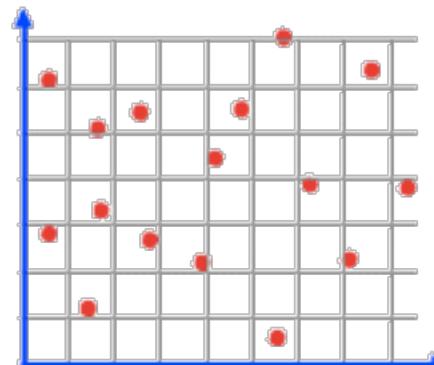
Strong negative correlation



Weak negative correlation

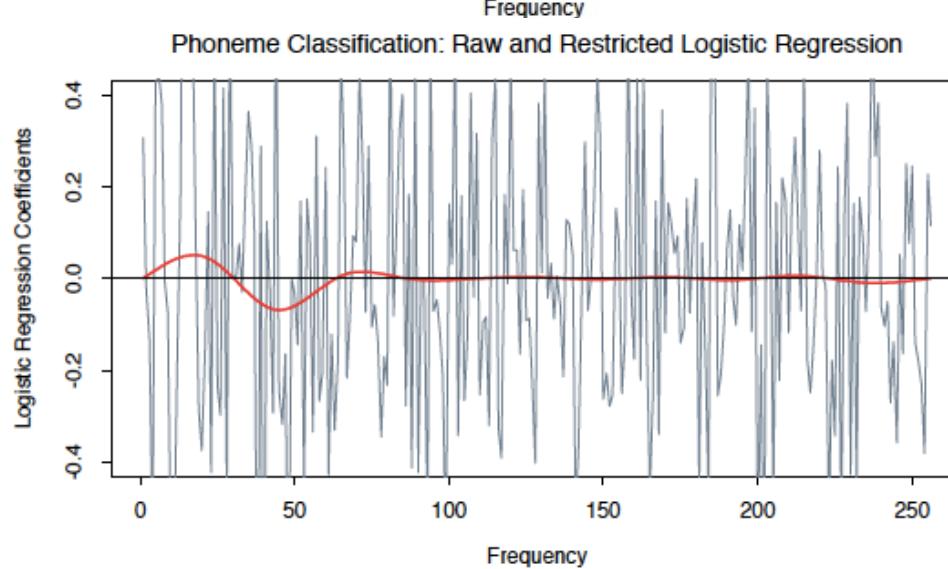
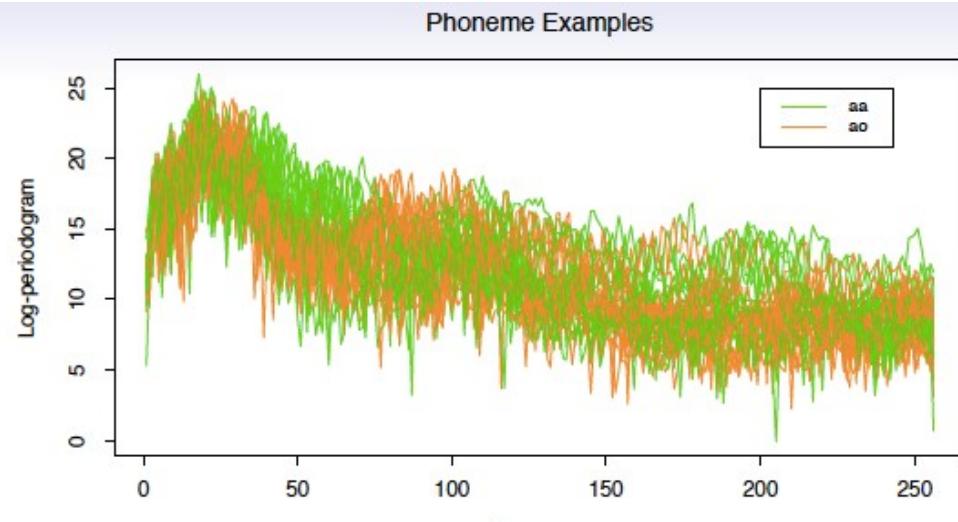


Moderate negative correlation

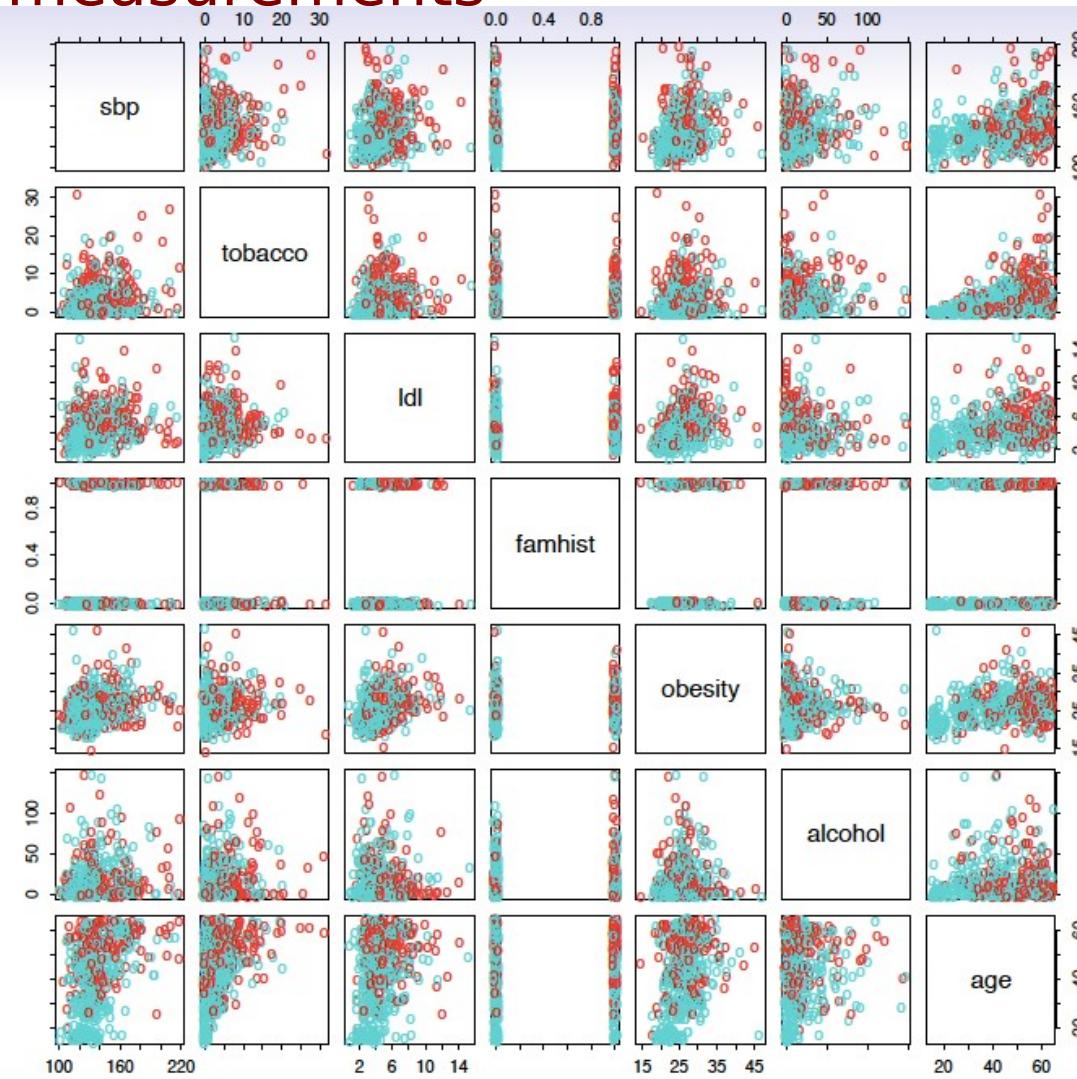


No correlation

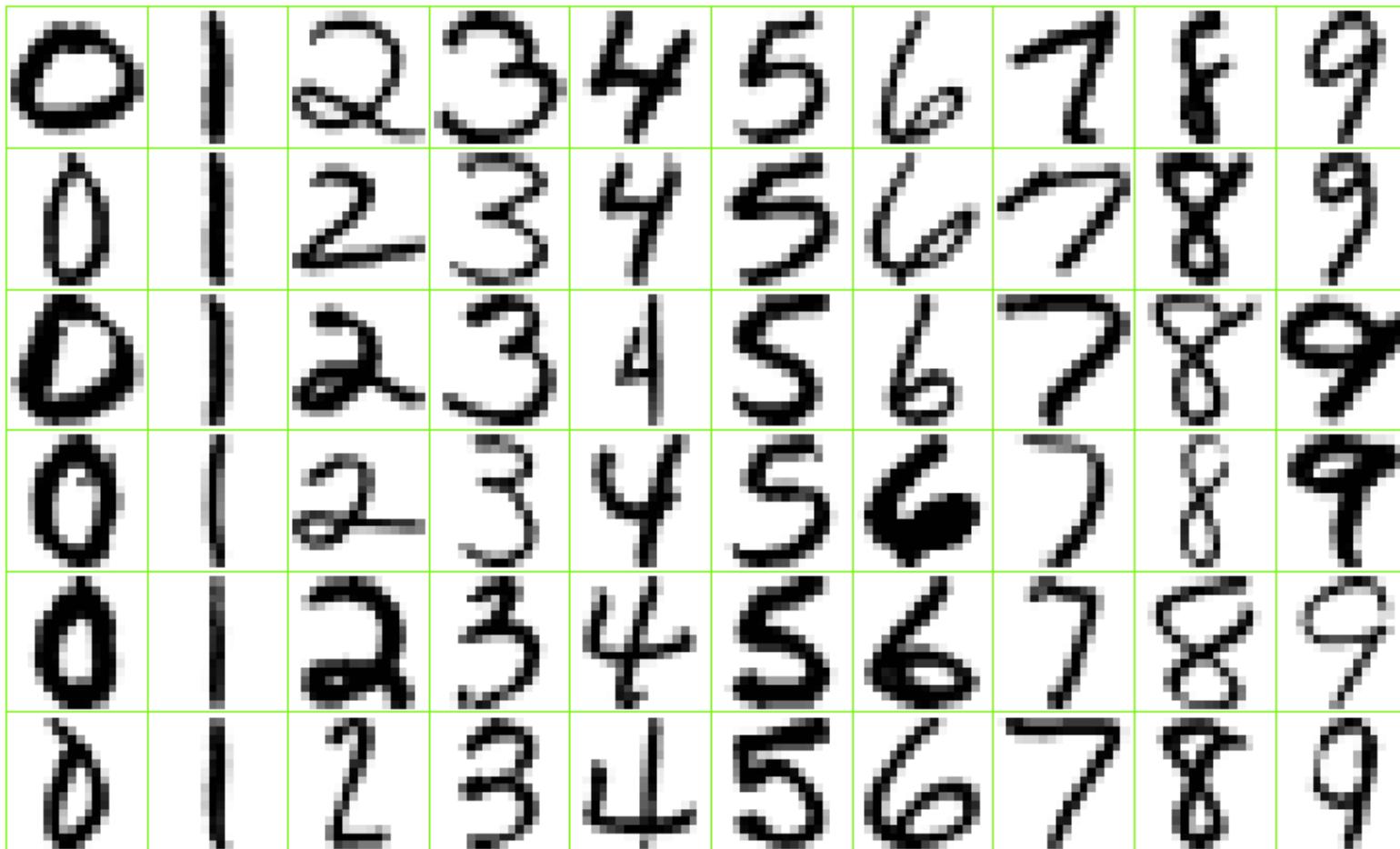
Classify a recorded phoneme based on a log-periodogram



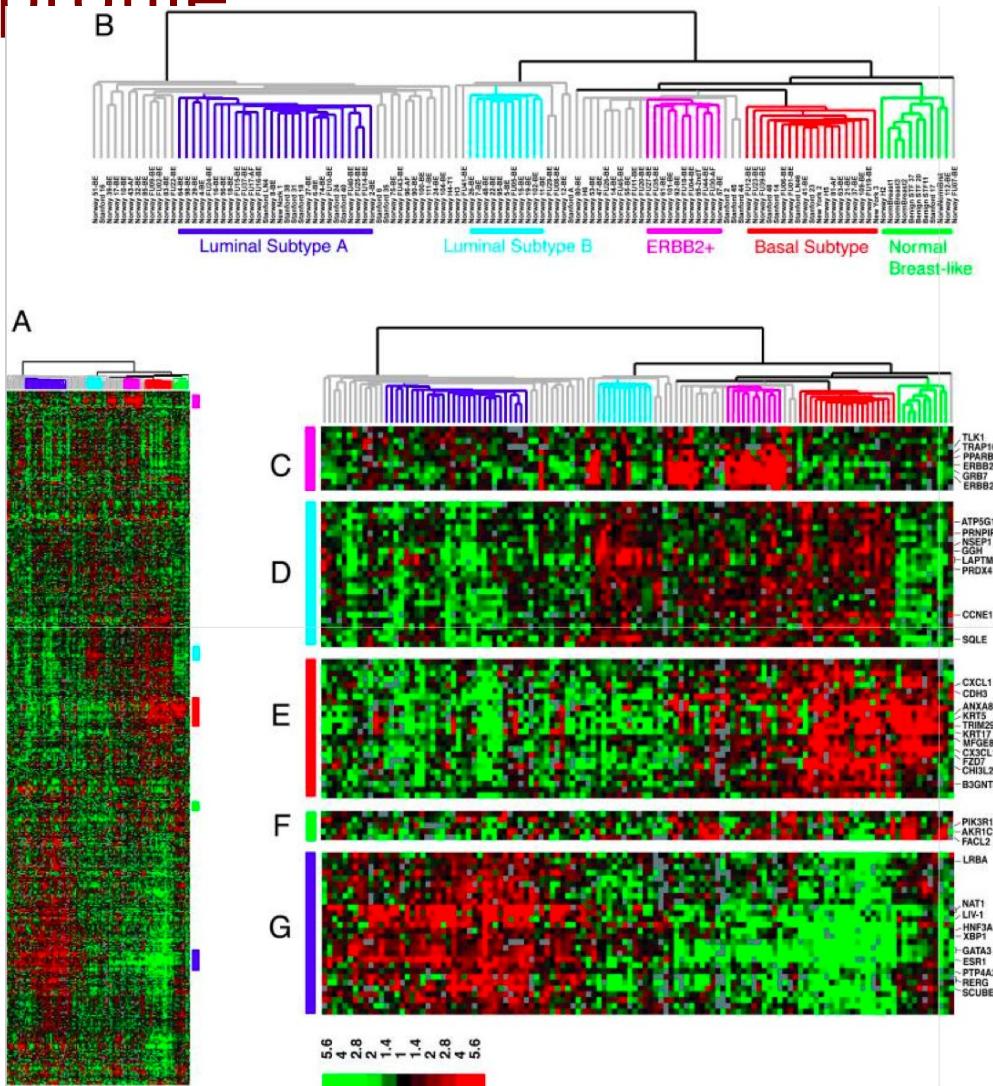
Predict whether someone will have a heart attack on the basis of demographic, diet and clinical measurements



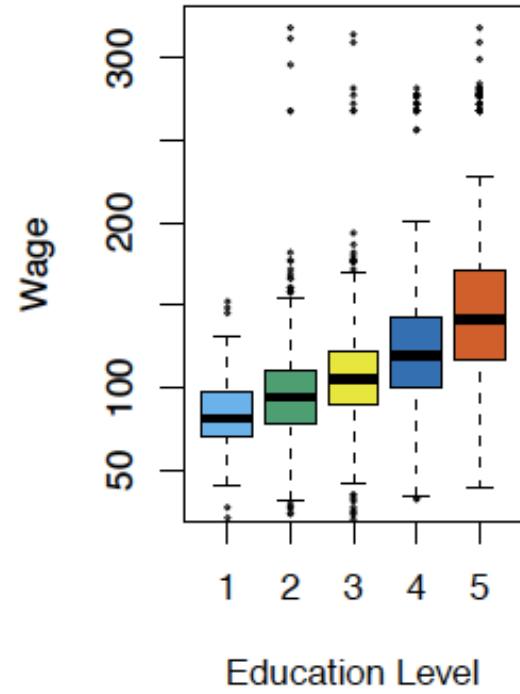
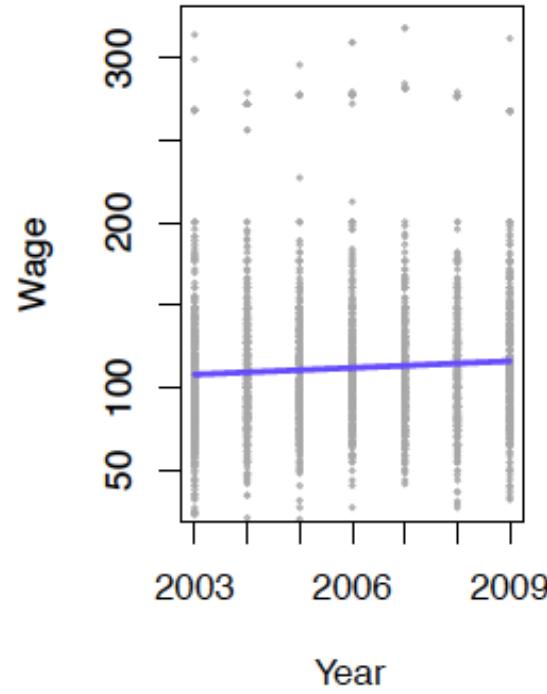
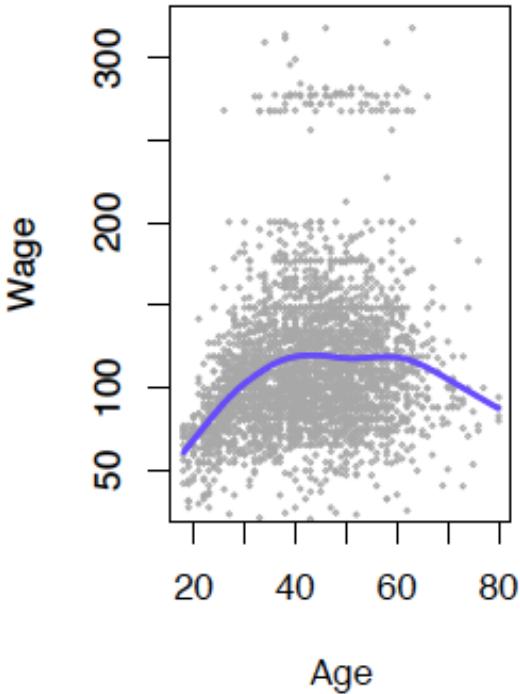
Identify the numbers in a handwritten zip code □ multiclass classification



Classify a tissue sample into one of several cancer classes, based on a gene expression profile

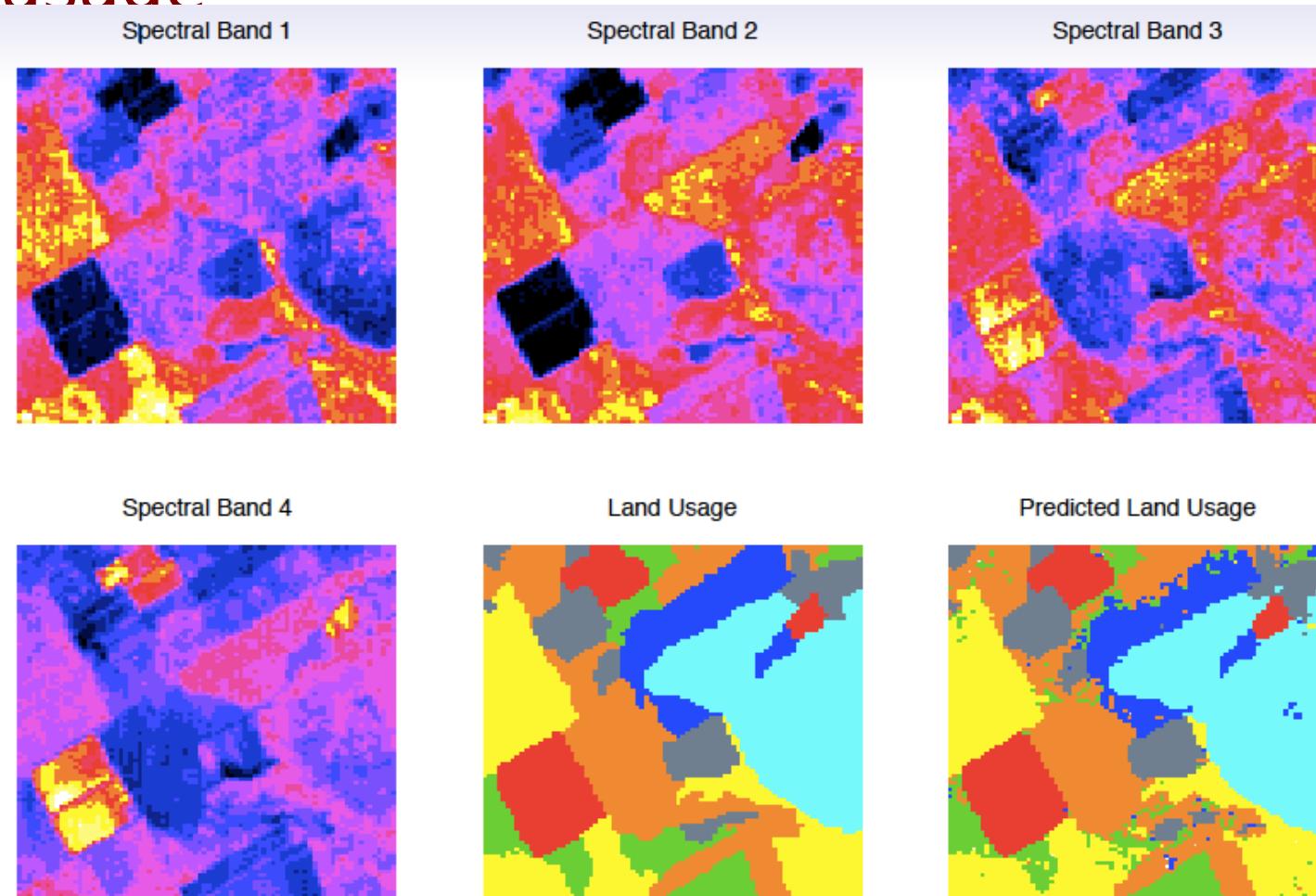


Establish the relationship between salary and demographic variables in population survey data

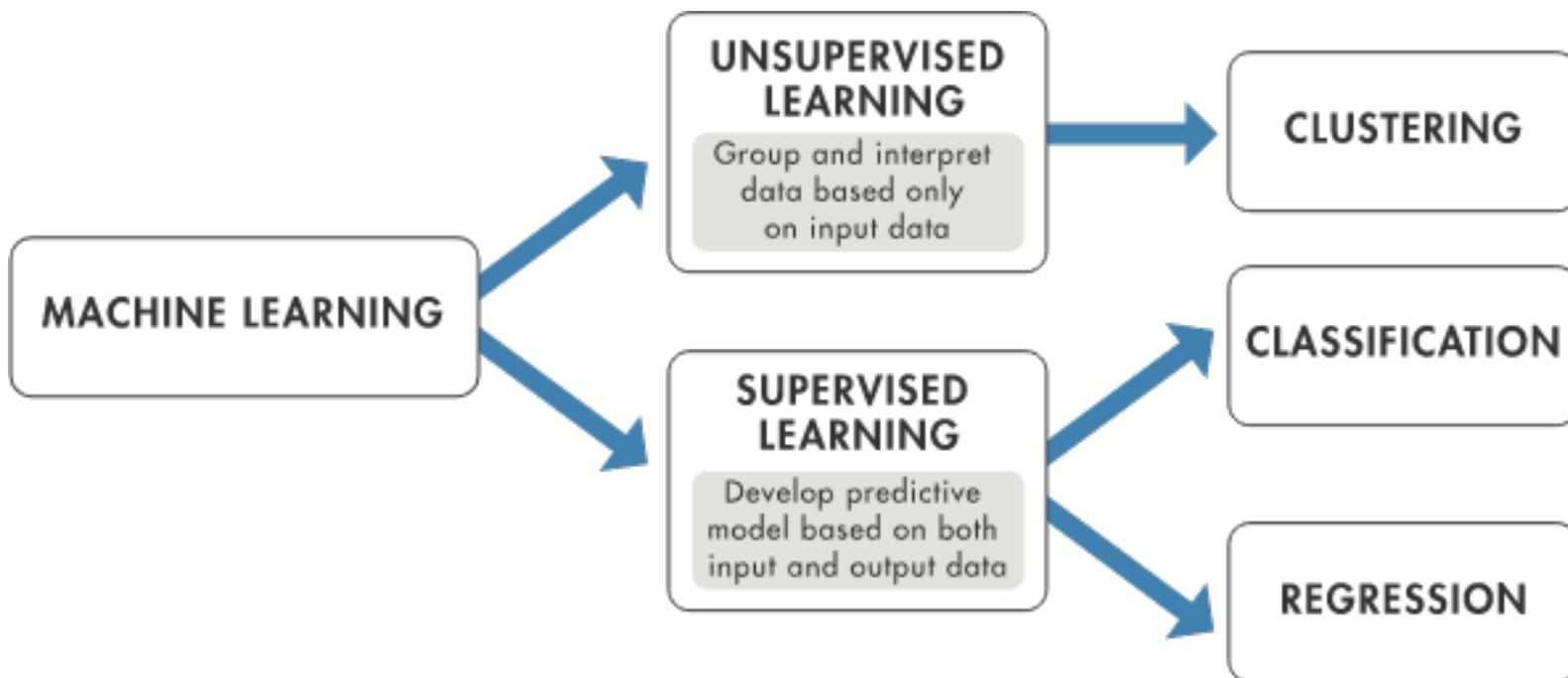


Income survey data for males from the central Atlantic region of the USA in 2009.

Classify the pixels in a LANDSAT image, by usage



Usage $\in \{\text{red soil, cotton, vegetation stubble, mixture, gray soil, damp gray soil}\}$



Machine learning \subseteq artificial intelligence

ARTIFICIAL INTELLIGENCE

Design an intelligent agent that perceives its environment and makes decisions to maximize chances of achieving its goal.

Subfields: vision, robotics, machine learning, natural language processing, planning, ...

MACHINE LEARNING

Gives "computers the ability to learn without being explicitly programmed" (Arthur Samuel, 1959)

SUPERVISED LEARNING

Classification, regression

UNSUPERVISED LEARNING

Clustering, dimensionality reduction, recommendation

REINFORCEMENT LEARNING

Reward maximization

Types of machine learning

Supervised Learning

- Labeled data
- Direct feedback
- Predict outcome/future

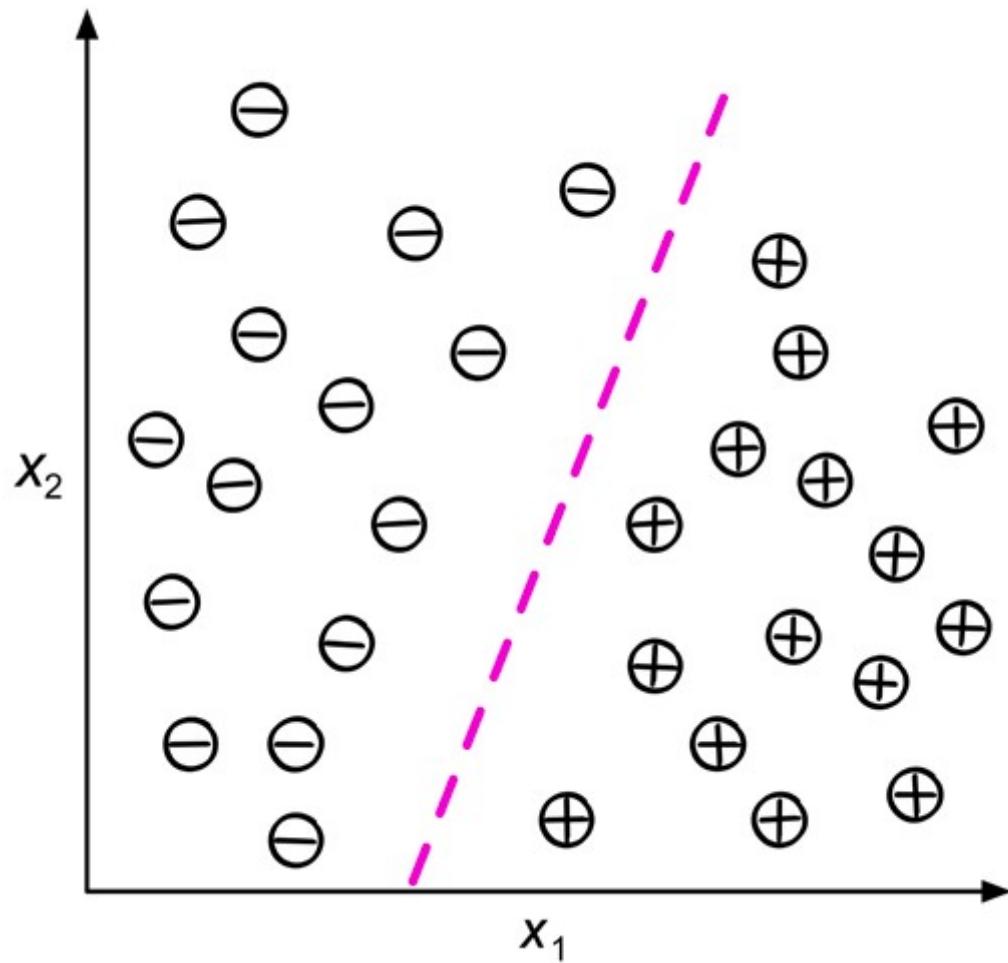
Unsupervised Learning

- No labels/targets
- No feedback
- Find hidden structure in data

Reinforcement Learning

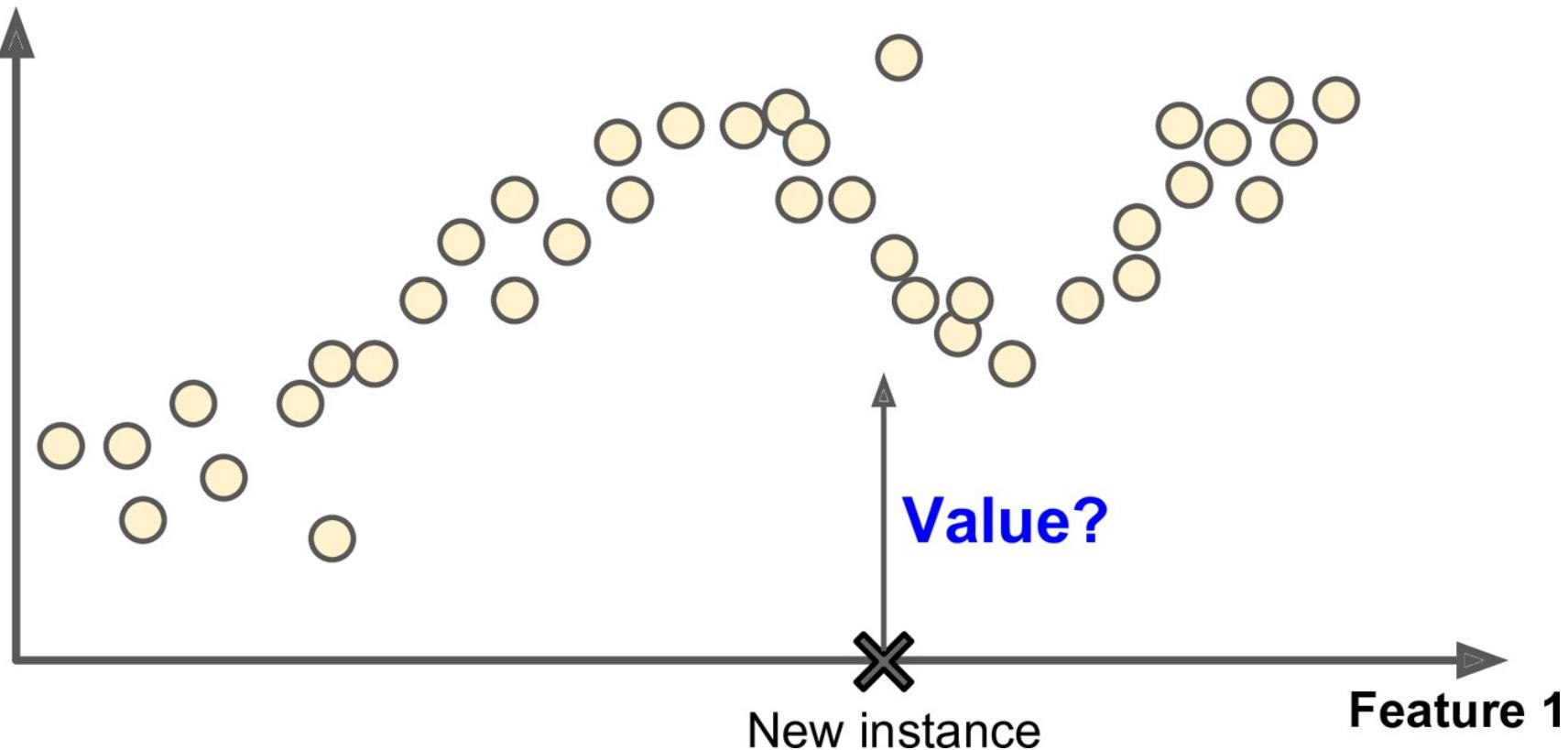
- Decision process
- Reward system
- Learn series of actions

Binary classification

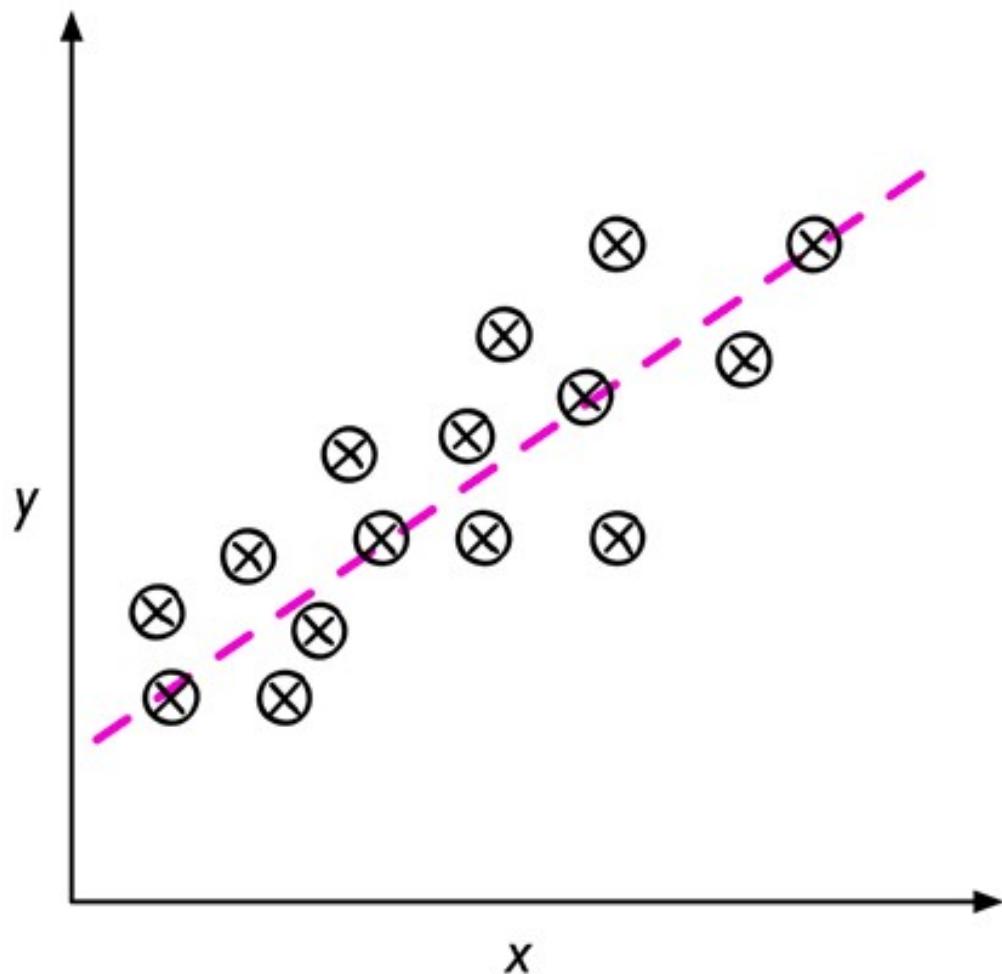


Regression

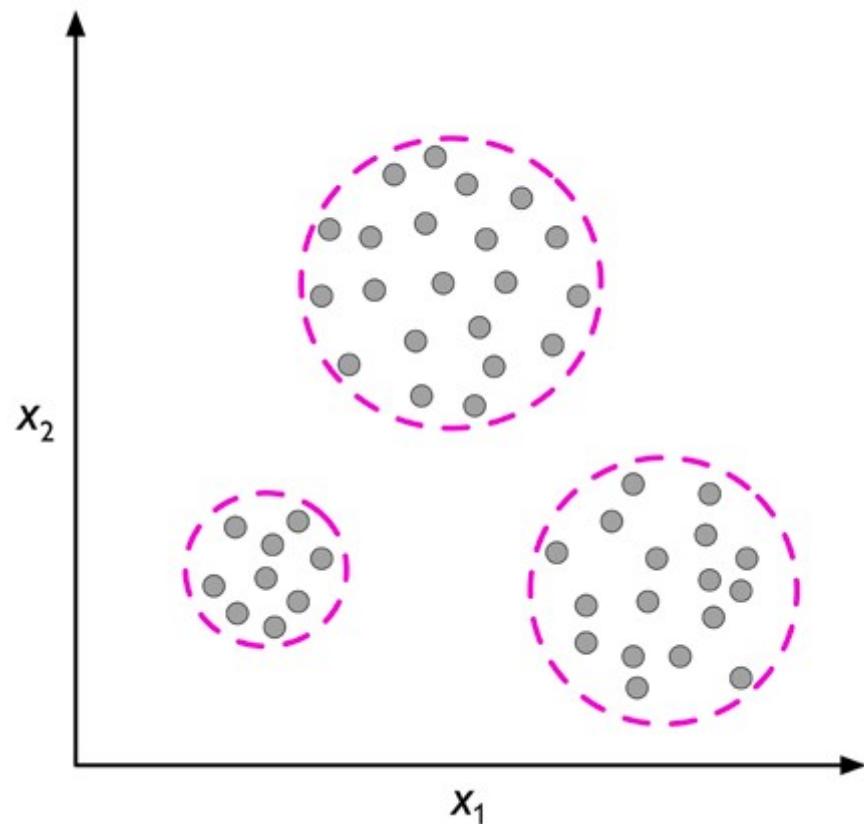
Value



Regression



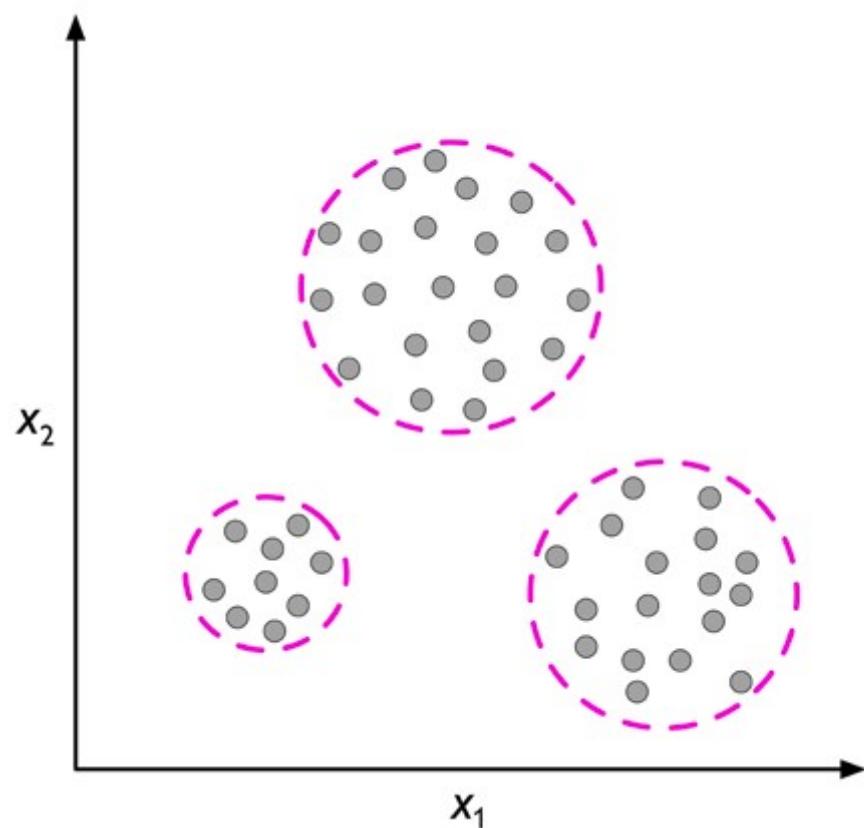
Clustering



Exploratory data analysis technique that allows to organize a pile of information into meaningful subgroups (clusters) without having any prior knowledge of their group memberships.

Each cluster defines a group of objects that share a certain degree of similarity but are more dissimilar to objects in other clusters.

Clustering



Clustering can be applied to organizing unlabeled data into three distinct groups based on the similarity of their features x_1 and x_2 .

For example, it allows marketers to discover customer groups based on their interests, in order to develop distinct marketing programs.

Dimensionality reduction

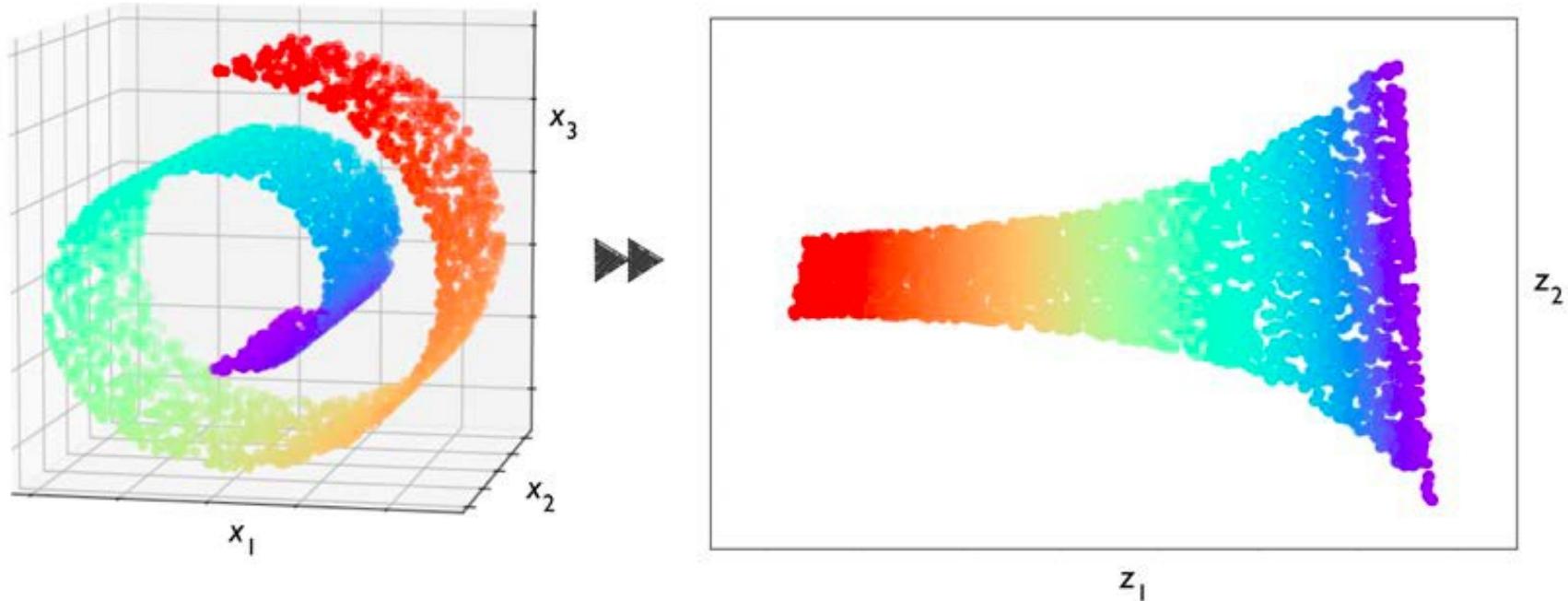
Commonly used approach in feature preprocessing to **remove noise from data**, which can also degrade the predictive performance of certain algorithms, and **compress the data** onto a smaller dimensional subspace while **retaining most of the relevant information**.

It also be useful for visualizing data:

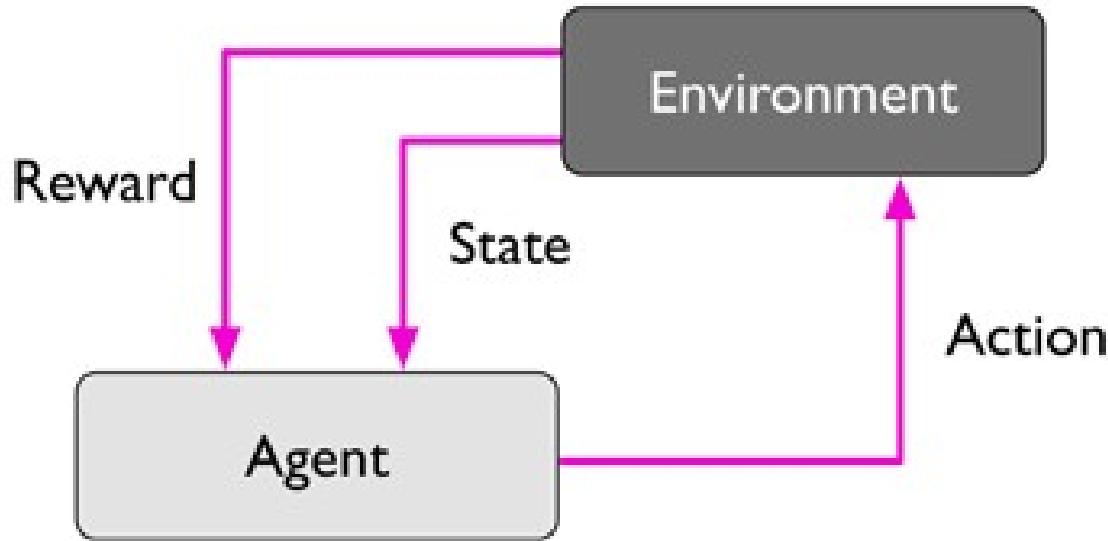
- A high-dimensional feature set can be projected onto one-, two-, or three-dimensional feature spaces in order to visualize it via 3D or 2D scatterplots or histograms.

Dimensionality reduction

Nonlinear dimensionality reduction applied to compress a 3D Swiss Roll onto a new 2D feature subspace:

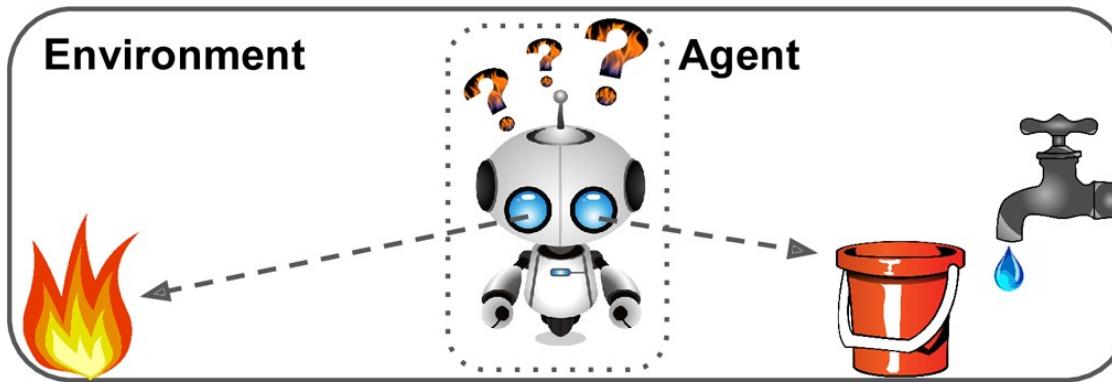


Reinforcement learning

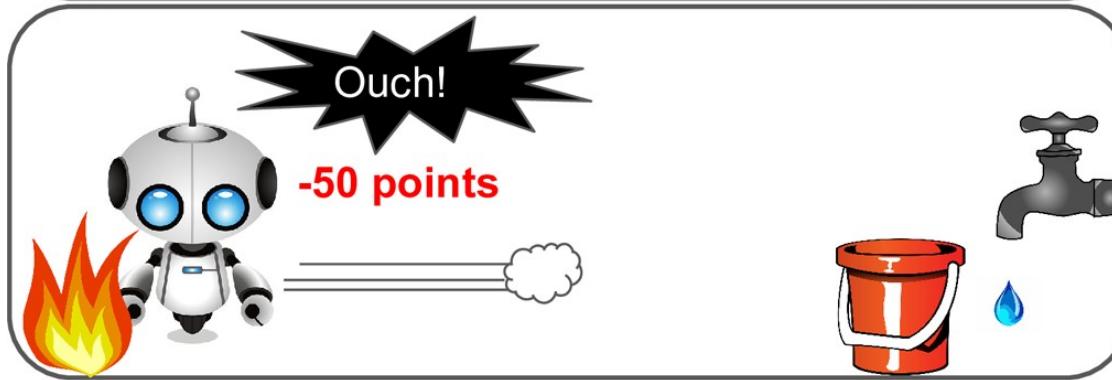


Reward can be defined as win or lose at the end of the game.

Reinforcement learning



- 1 Observe
- 2 Select action using policy



- 3 Action!
- 4 Get reward or penalty



- 5 Update policy (learning step)
- 6 Iterate until an optimal policy is found

Reinforcement learning



Reinforcement learning

DeepMind's AlphaGo program learned its winning policy by analyzing millions of games, and then playing many games against itself.

Note that learning was turned off during the games against the champion; AlphaGo was just applying the policy it had learned.

Deep Reinforcement Learning

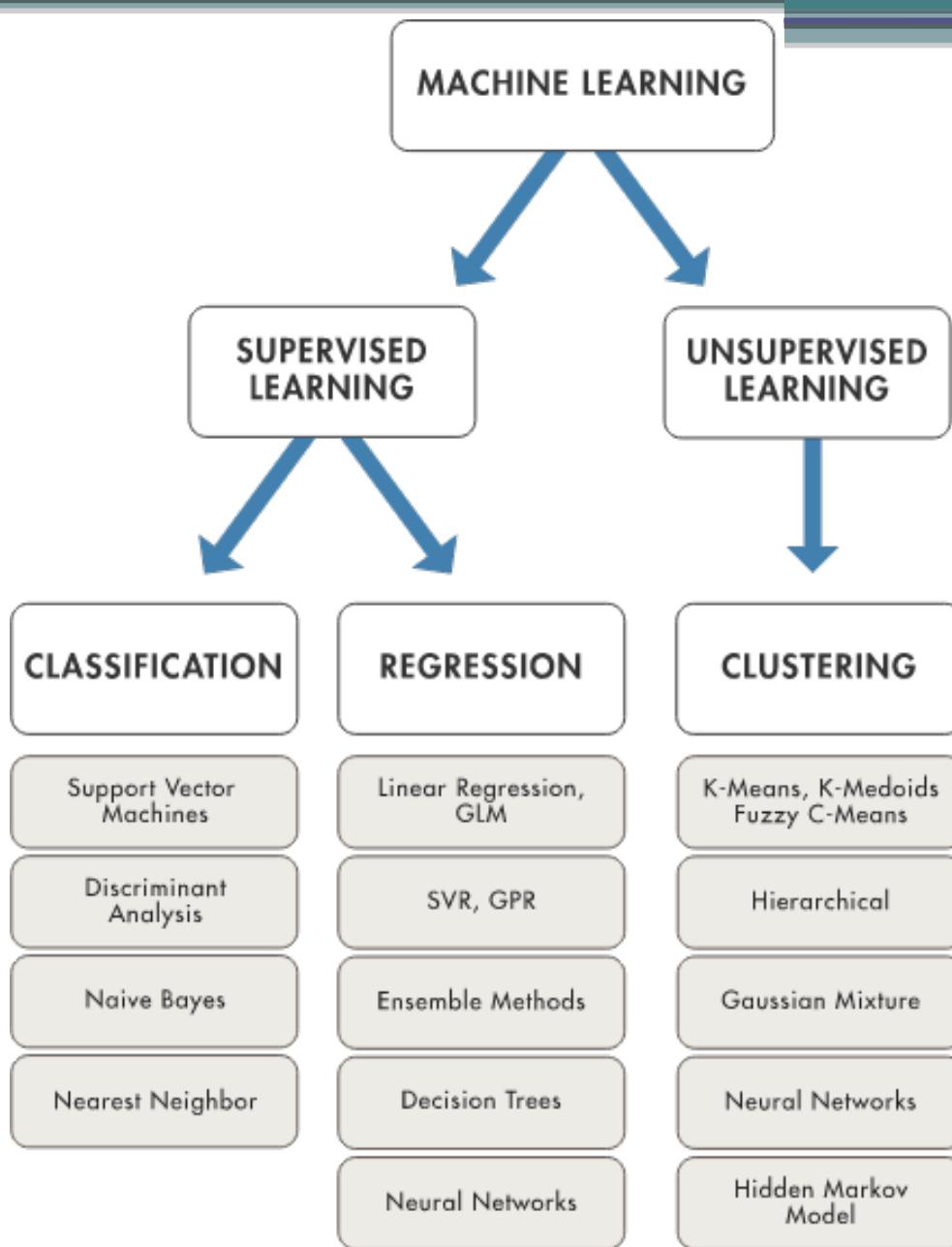
2015 – AlphaGo learned from moves based on the historical tournament data.

2017 – AlphaGo Zero - AlphaGo Zero's strategies were self-taught.

2018 - AlphaZero, a modified version of AlphaGo Zero, gained superhuman abilities at chess and shogi. Like AlphaGo Zero, AlphaZero learned solely through self-play.

2018 – AlphaFold - protein-folding, one of the toughest problems in science.

2019 – AlphaStar - knowledge equivalent to 200 years of playing time.

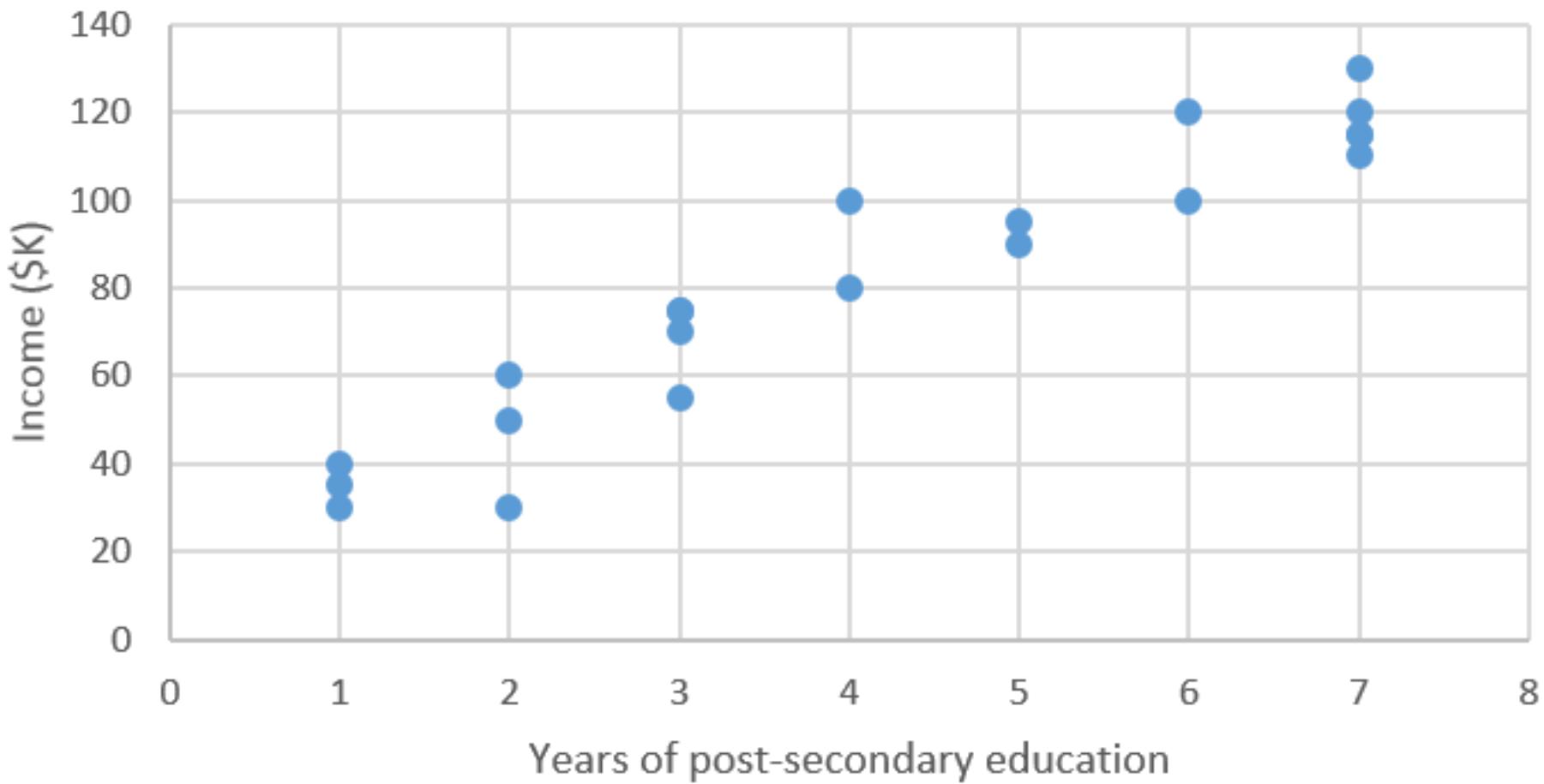


Supervised Learning: Regression

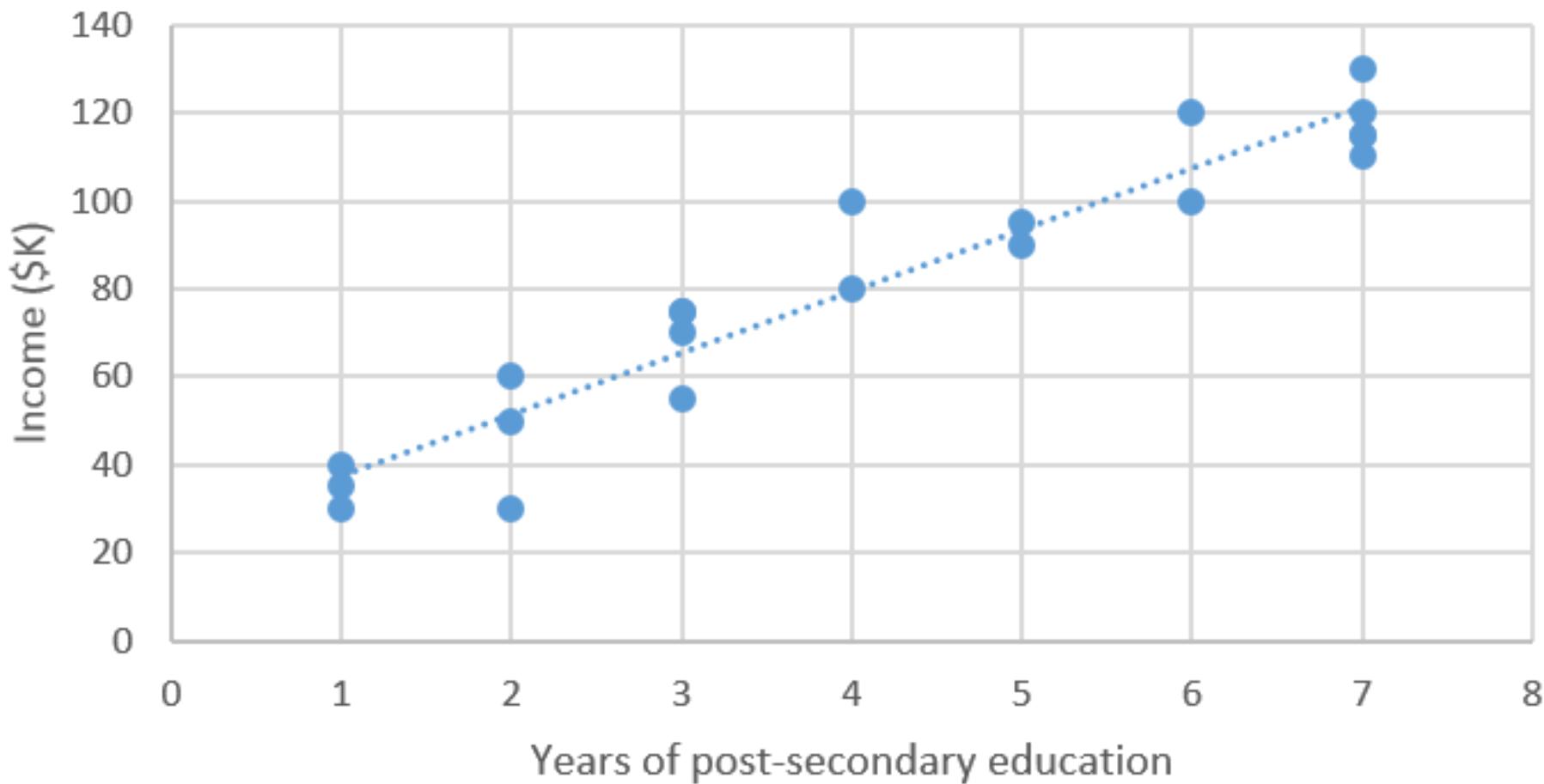
Observation #	Years of Higher Education (X)	Income (Y)
1	4	\$80,000
2	5	\$91,500
3	0	\$42,000
4	2	\$55,000
...
N	6	\$100,000

1	4	???
2	6	???

Income



Income



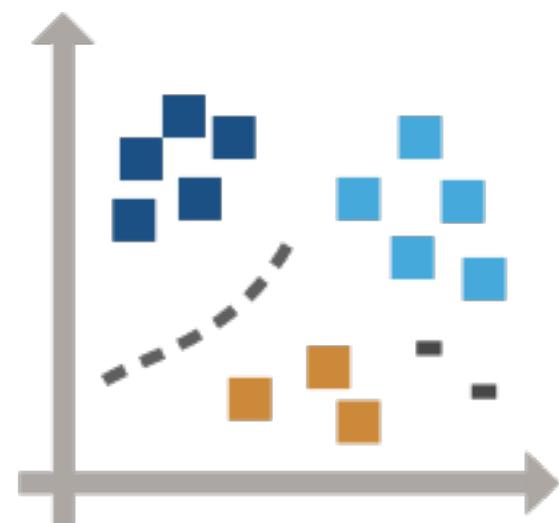
Supervised Learning: Classification

Observation #	Input image (X)	Label (Y)
1		"dog"
2		"cat"
3		"dog"
...
N		"dog"
test set	1	
	2	

Aprendizagem não supervisionada



Clustering
Patterns in
the Data



Unsupervised learning algorithms

Clustering

- k-Means
- Hierarchical Cluster Analysis (HCA)
- Expectation Maximization

Visualization and dimensionality reduction

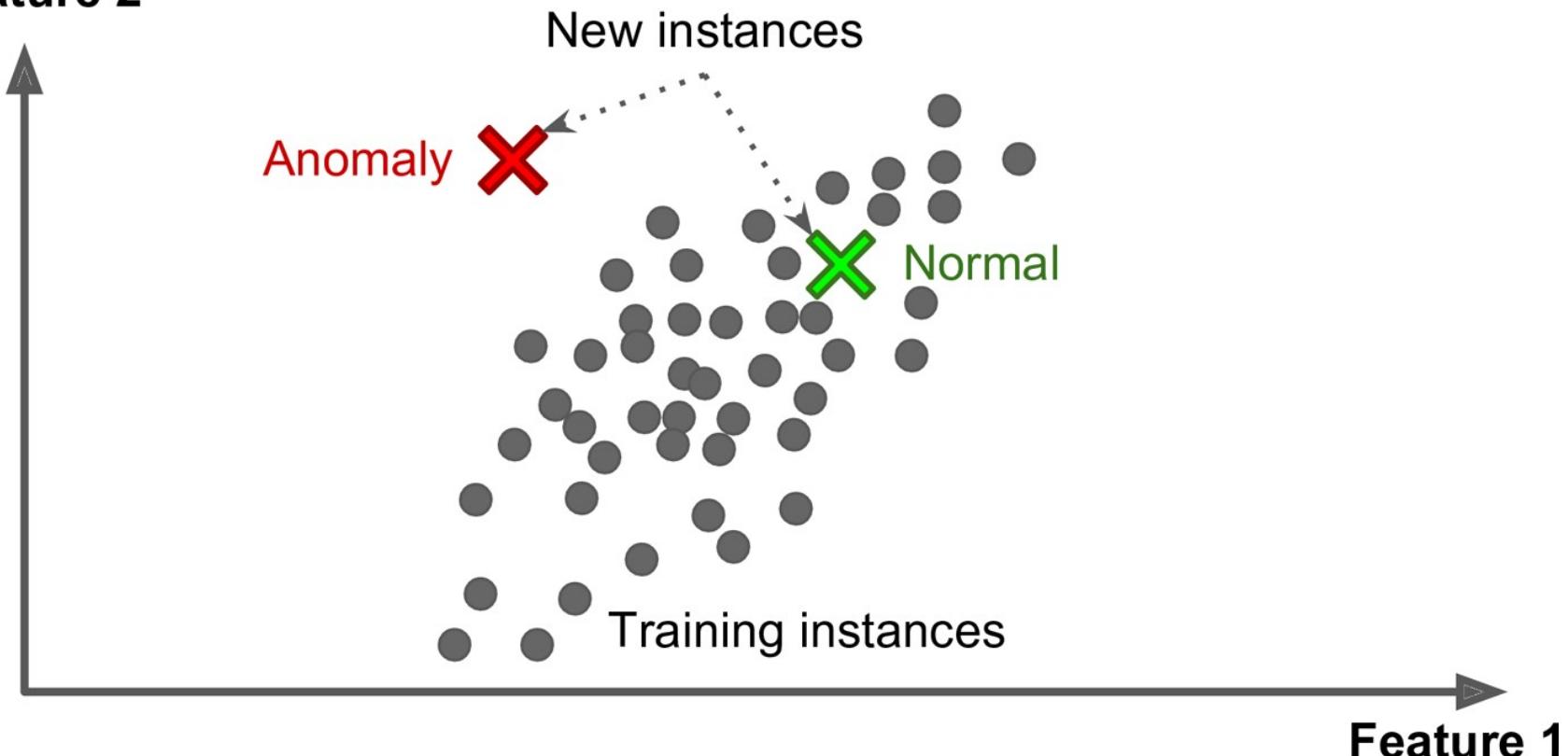
- Principal Component Analysis (PCA)
- Kernel PCA
- Locally-Linear Embedding (LLE)
- t-distributed Stochastic Neighbor Embedding (t-SNE)

Association rule learning

- Apriori
- Eclat

Unsupervised task - anomaly detection

Feature 2



Unsupervised task - anomaly detection

Examples:

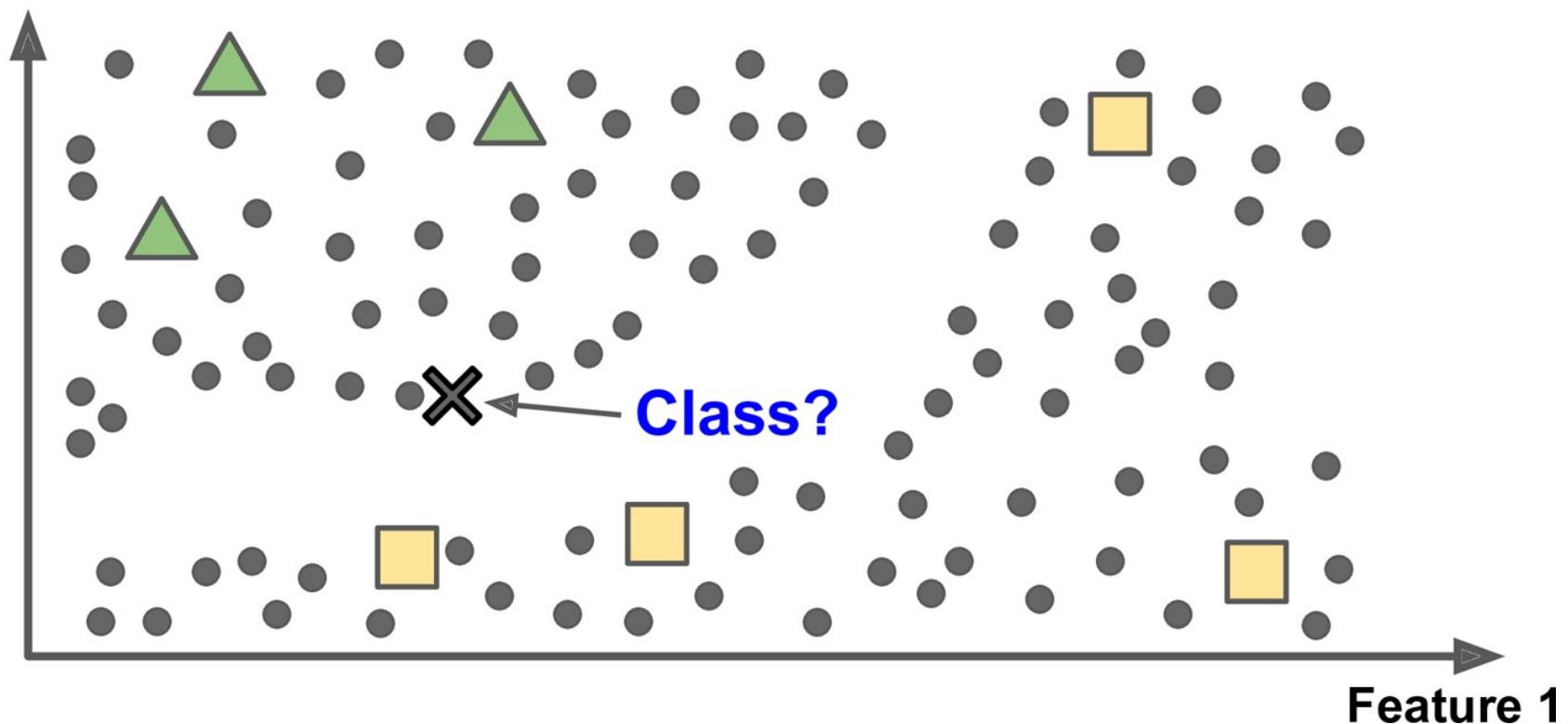
- Detecting unusual credit card transactions to prevent fraud;
- Catching manufacturing defects;
- Automatically removing outliers from a dataset before feeding it to another learning algorithm.

The system is trained with normal instances, and when it sees a new instance it can tell whether it looks like a normal one or whether it is likely an anomaly.

Semisupervised learning

Some algorithms can deal with partially labeled training data, usually a lot of unlabeled data and a little bit of labeled data.

Feature 2



Semisupervised learning

Photo-hosting services, such as Google Photos:

- **Unsupervised learning (clustering)**
 - Once you upload all your family photos to the service, it automatically recognizes that the same person A shows up in photos 1, 5, and 11, while another person B shows up in photos 2, 5, and 7.
- **Supervised learning**
 - Tell the system who these people are. Just one label per person, and it is able to name everyone in every photo, which is useful for searching photos.

O prêmio Netflix

- competition started in October 2006. Training data is ratings for 18, 000 movies by 400, 000 Netflix customers, each rating between 1 and 5.
- training data is very sparse— about 98% missing.
- objective is to predict the rating for a set of 1 million customer-movie pairs that are missing in the training data.
- Netflix's original algorithm achieved a root MSE of 0.953. The first team to achieve a 10% improvement wins one million

O prêmio Netflix



Netflix Prize

COMPLETED

[Home](#) [Rules](#) [Leaderboard](#) [Update](#)

Leaderboard

Showing Test Score. [Click here to show quiz score](#)

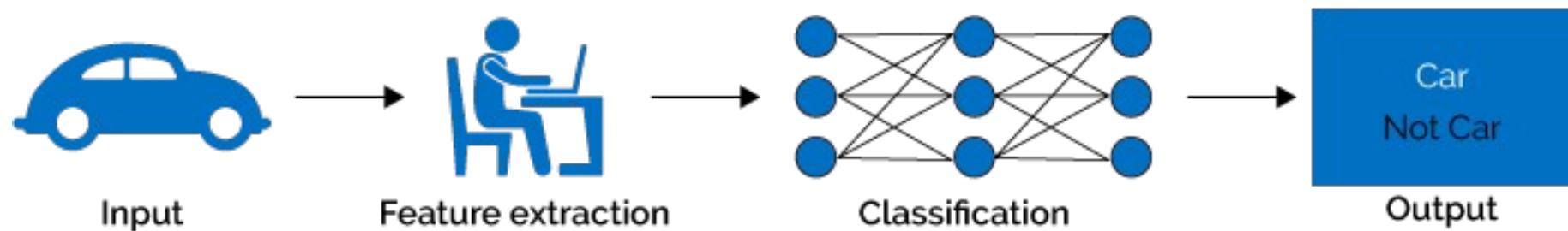
Display top leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries !	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

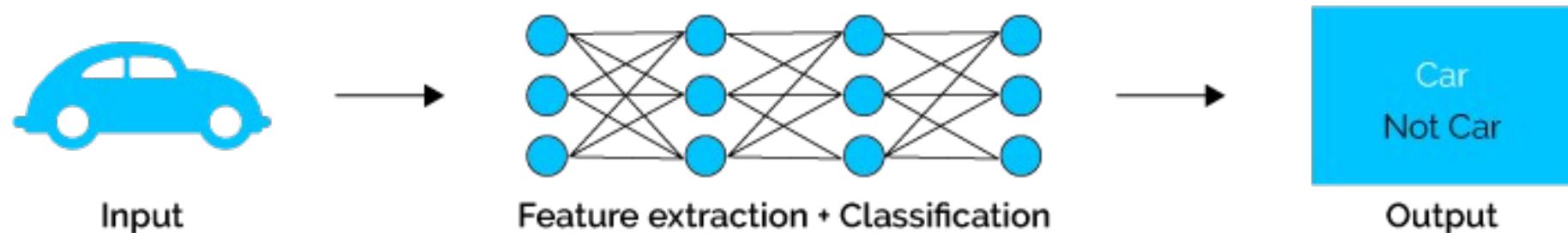
Statistical Learning versus Machine Learning

- Machine learning arose as a subfield of Artificial Intelligence.
- Statistical learning arose as a subfield of Statistics.
- There is much overlap — both fields focus on supervised and unsupervised problems:
 - Machine learning has a greater emphasis on **large scale** applications and **prediction accuracy**.
 - Statistical learning emphasizes **models** and their interpretability, and **precision** and **uncertainty**.

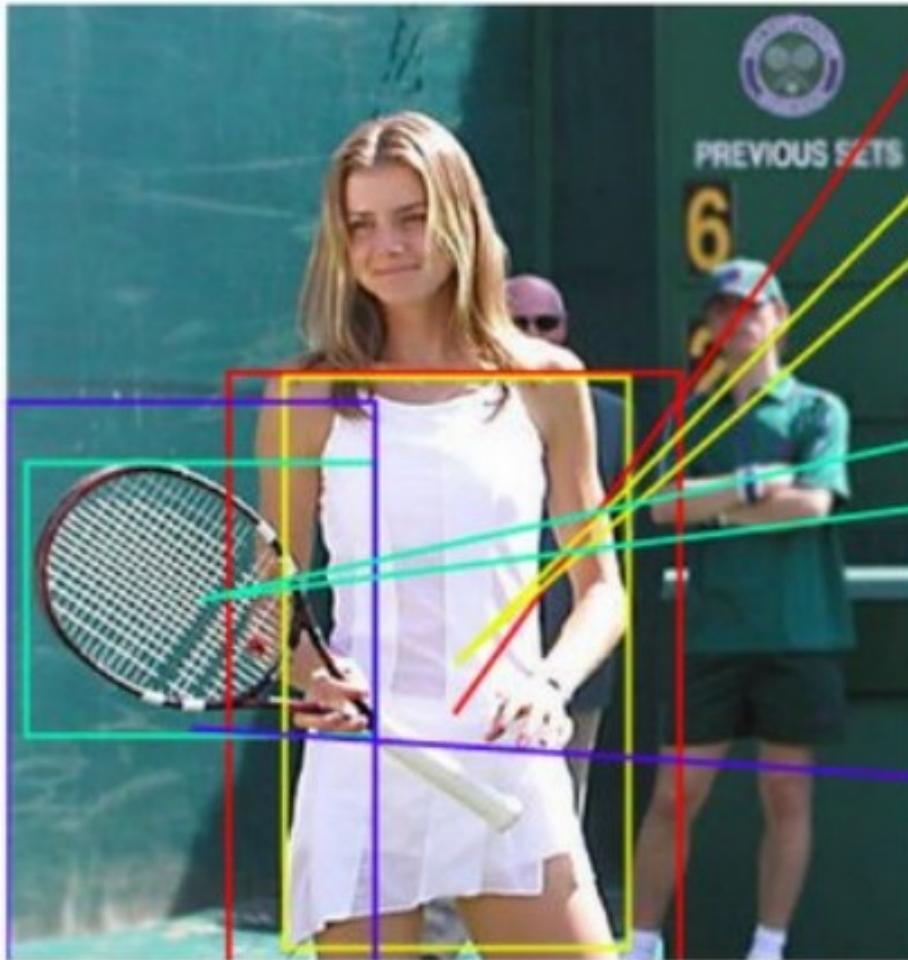
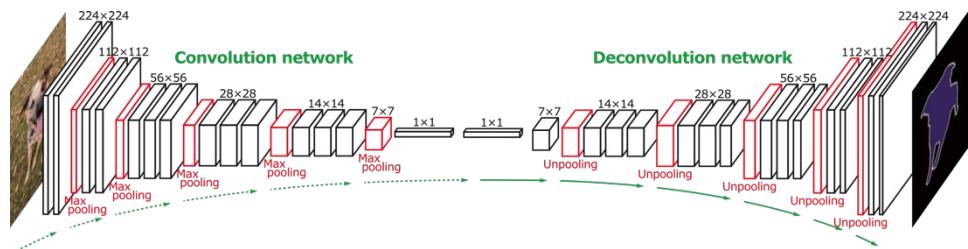
Machine Learning



Deep Learning



Deep Learning



1.12 woman

-0.28 in

1.23 white

1.45 dress

0.06 standing

-0.13 with

3.58 tennis

1.81 racket

0.06 two

0.05 people

-0.14 in

0.30 green

-0.09 behind

-0.14 her

Image,
automatically
annotated by
Deep Learning.

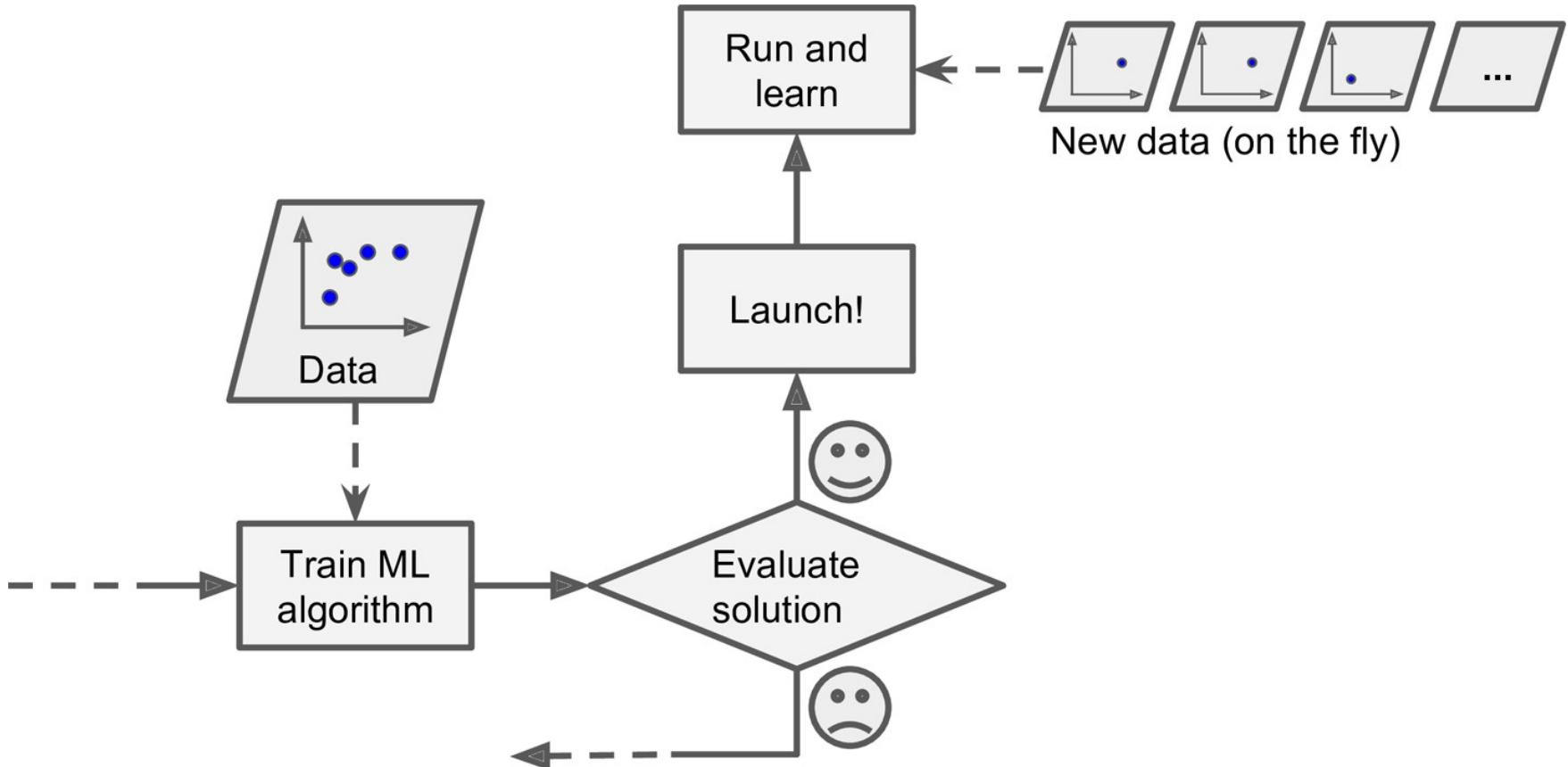
Batch learning

- In batch learning, the system is incapable of learning incrementally: it must be trained using all the available data.
- This will generally take a lot of time and computing resources, so it is typically done offline (offline learning).
- The whole process of training, evaluating, and launching a Machine Learning system can be automated.
 - Batch learning system can adapt to change.

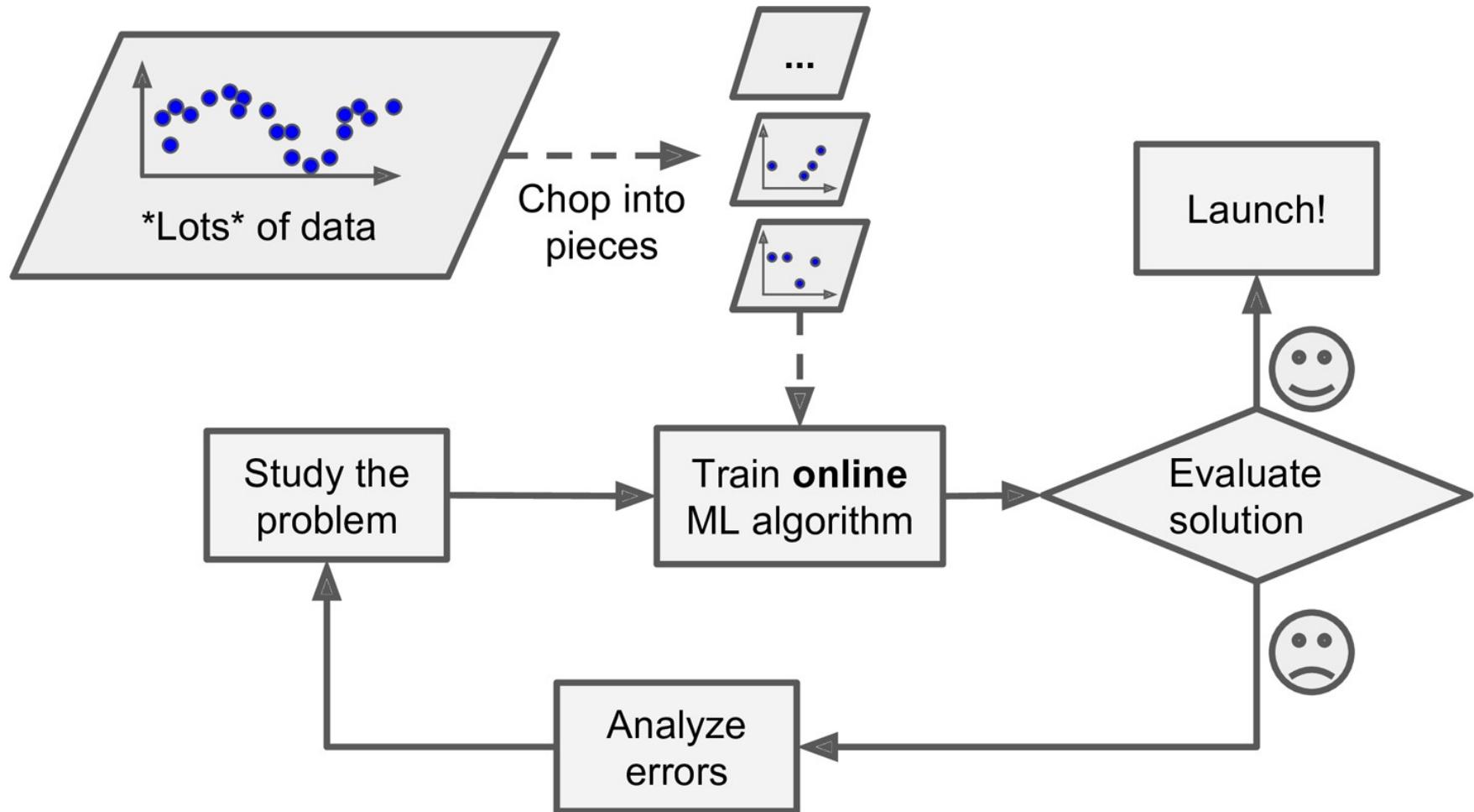
Online learning

- In online learning, you train the system incrementally by feeding it data instances sequentially, either individually or by small groups called mini-batches.
- Each learning step is fast and cheap, so the system can learn about new data on the fly, as it arrives

Online learning



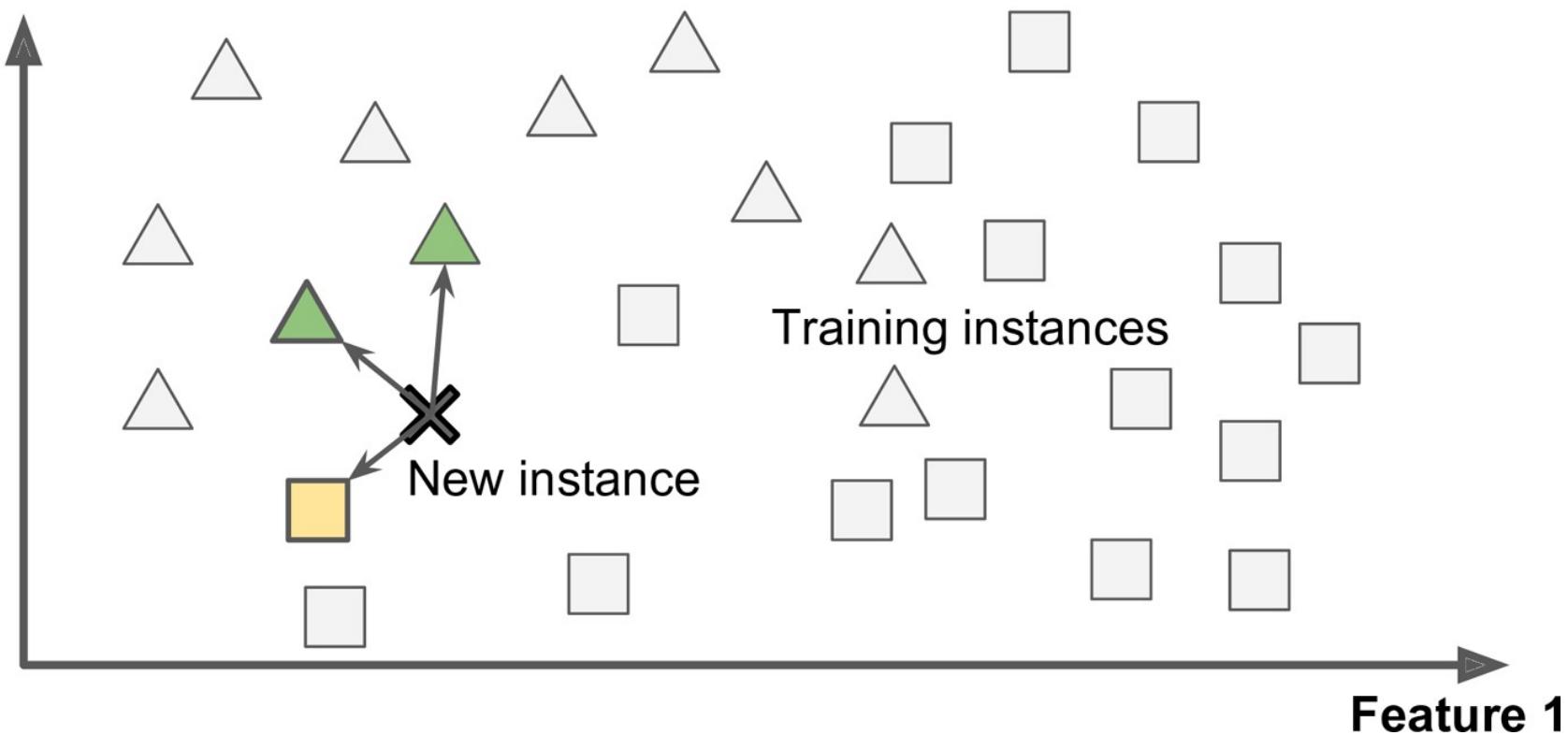
Online learning to handle huge datasets



Instance-based learning

The system learns the examples by heart, then generalizes to new cases using a similarity measure.

Feature 2



Instance-based learning

A (very basic) similarity measure between two emails could be to count the number of words they have in common.

The system would flag an email as spam if it has many words in common with a known spam email.

Model-based learning

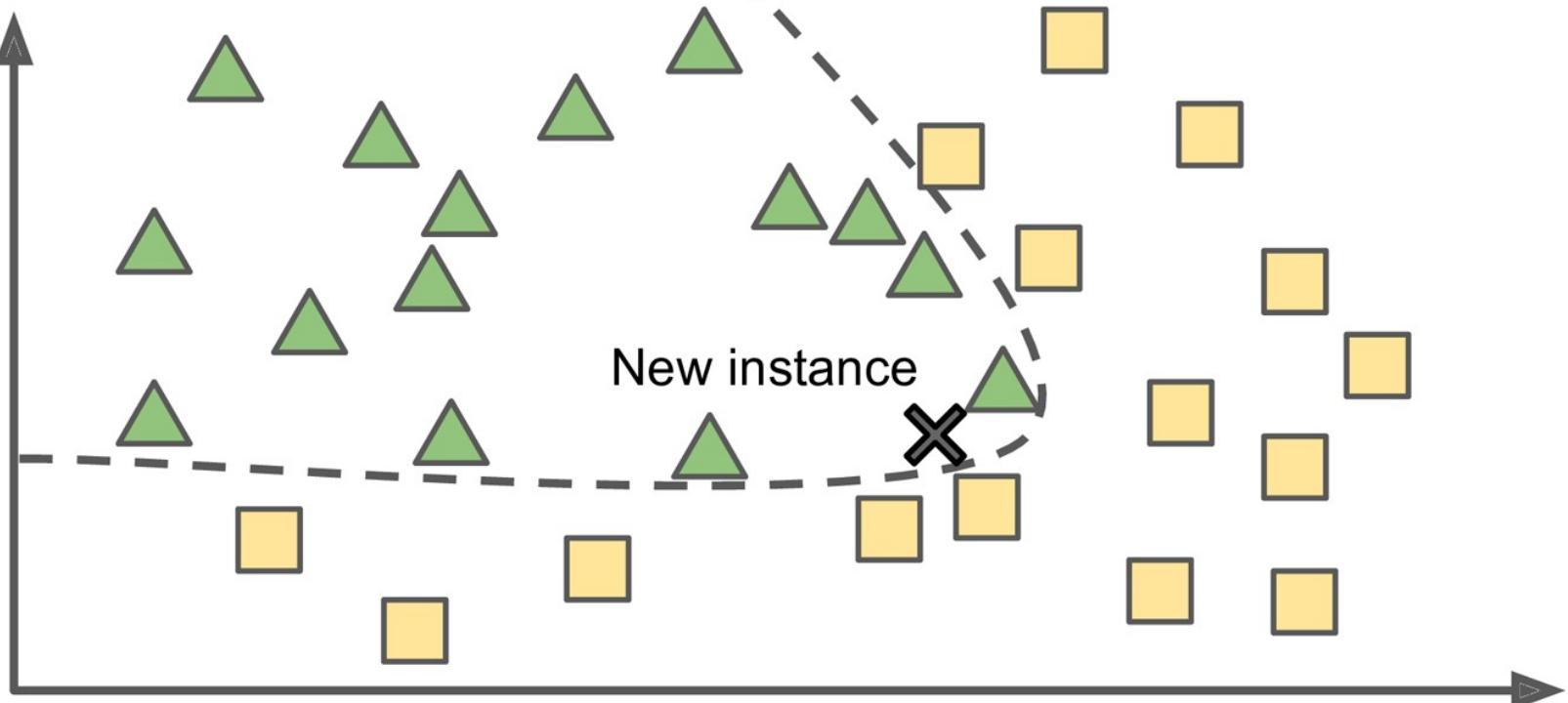
Build a model of these examples, then use that model to make predictions.

Feature 2

Model

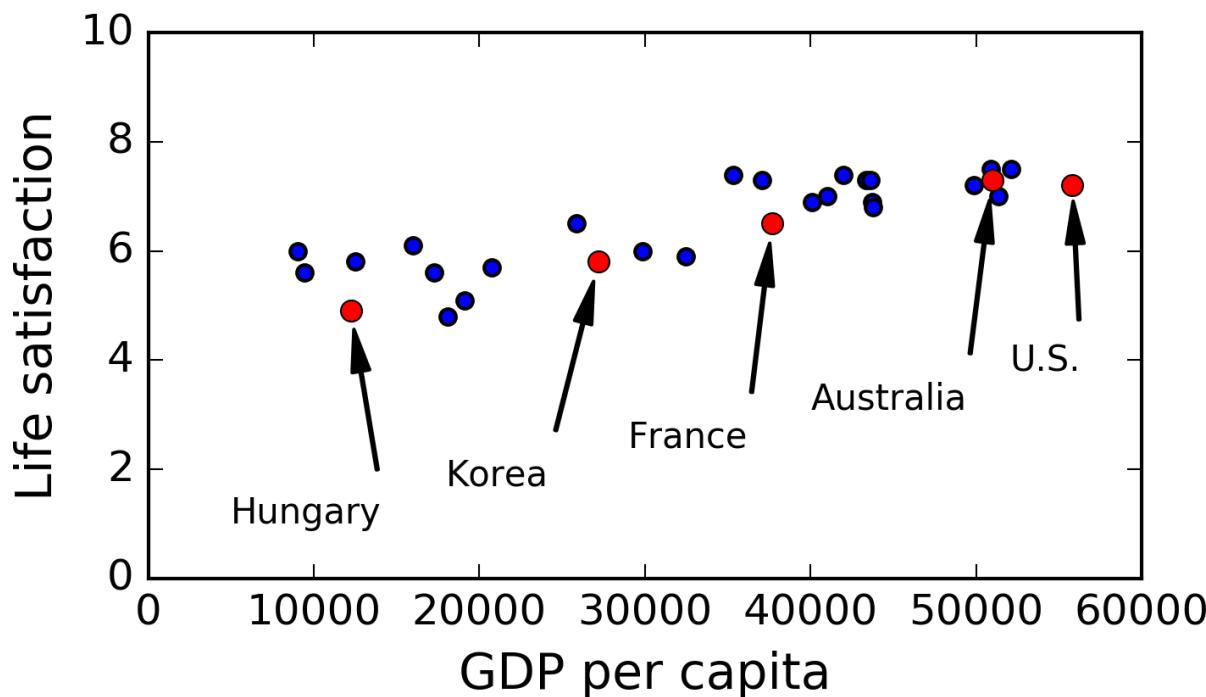
New instance

Feature 1

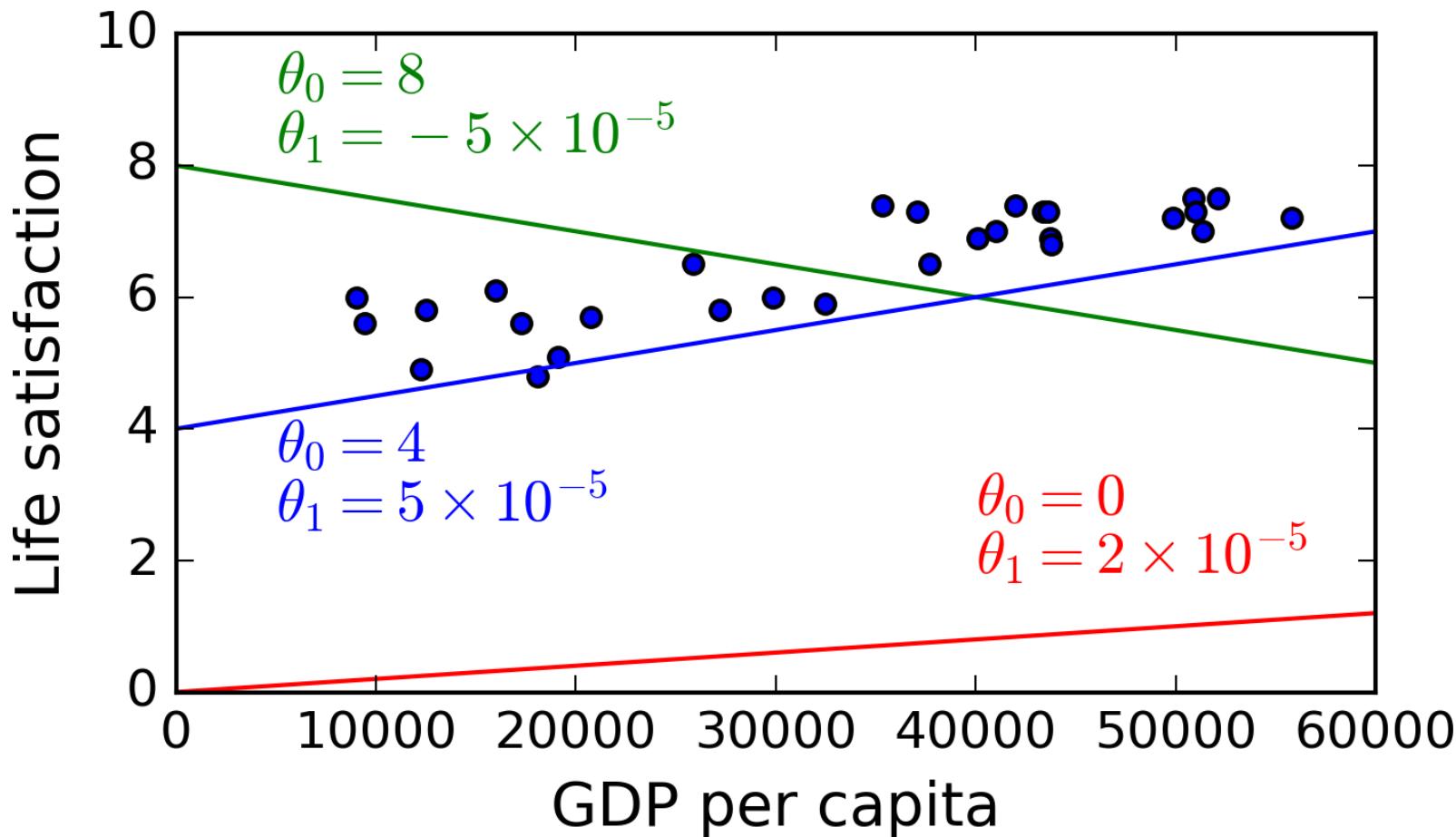


Model-based learning

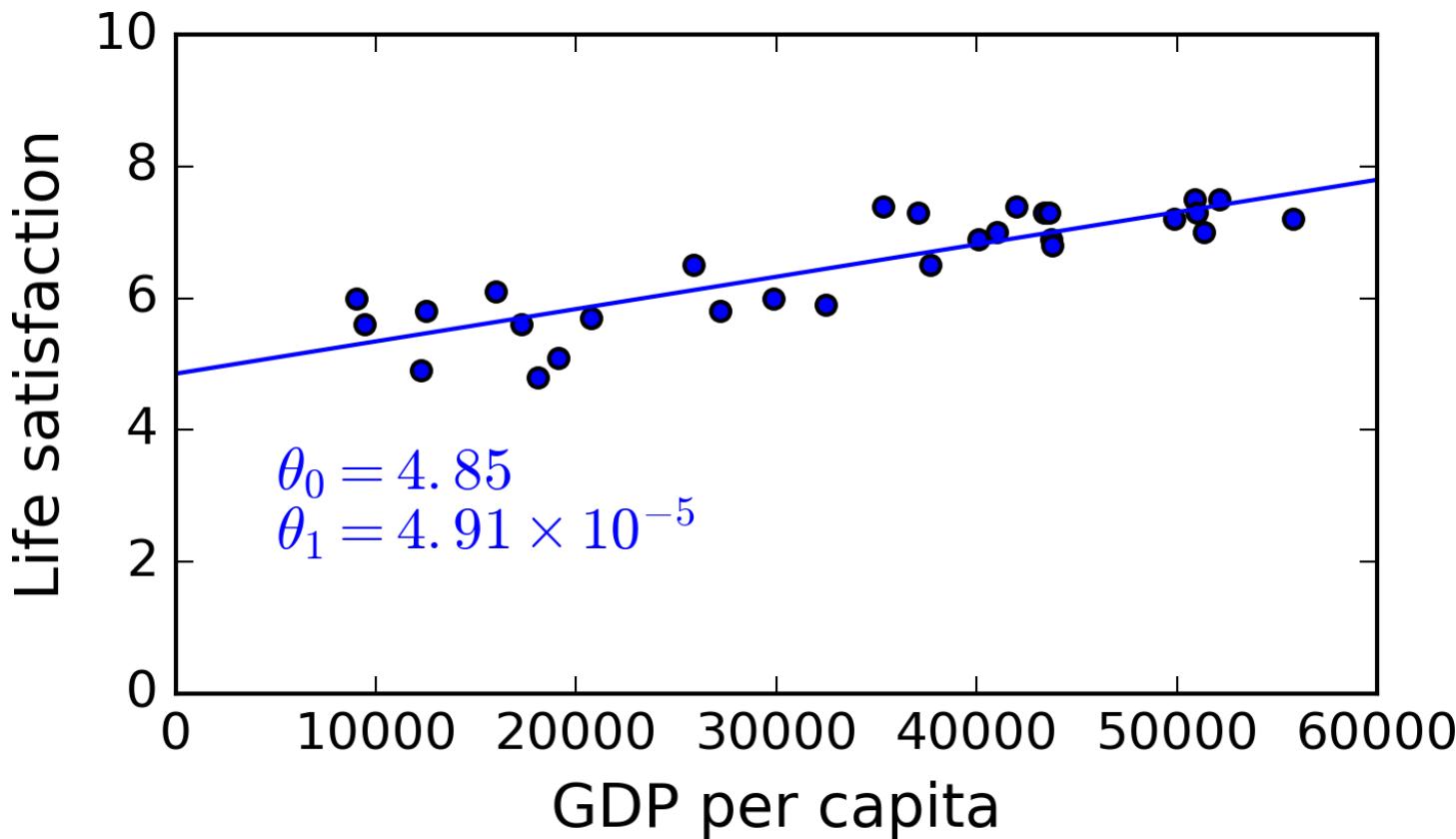
Country	GDP per capita (USD)	Life satisfaction
Hungary	12,240	4.9
Korea	27,195	5.8
France	37,675	6.5
Australia	50,962	7.3
United States	55,805	7.2



Model-based learning



Model-based learning



Cyprus's GDP per capita, find \$22,587, and then apply your model and find that life satisfaction is likely to be somewhere around 5.96.

Instance based learning

Slovenia has the closest GDP per capita to that of Cyprus (\$20,732), and since Slovenians' life satisfaction is 5.7, you would have predicted a life satisfaction of 5.7 for Cyprus.

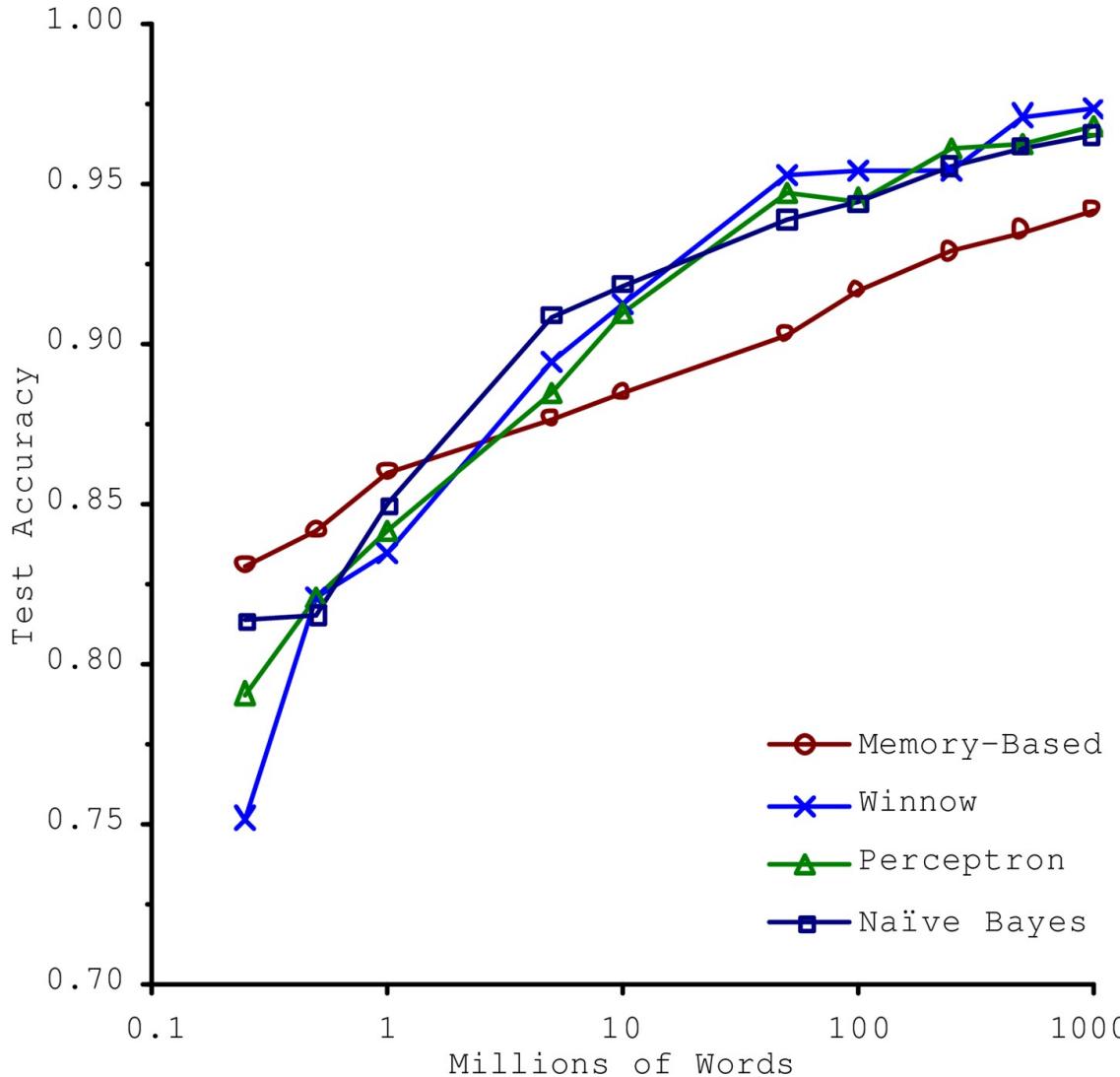
If you zoom out a bit and look at the two next closest countries, you will find Portugal and Spain with life satisfactions of 5.1 and 6.5, respectively.

Averaging these three values, you get 5.77, which is pretty close to your model-based prediction.

This simple algorithm is called k-Nearest Neighbors regression (in this example, $k = 3$).

Desafios

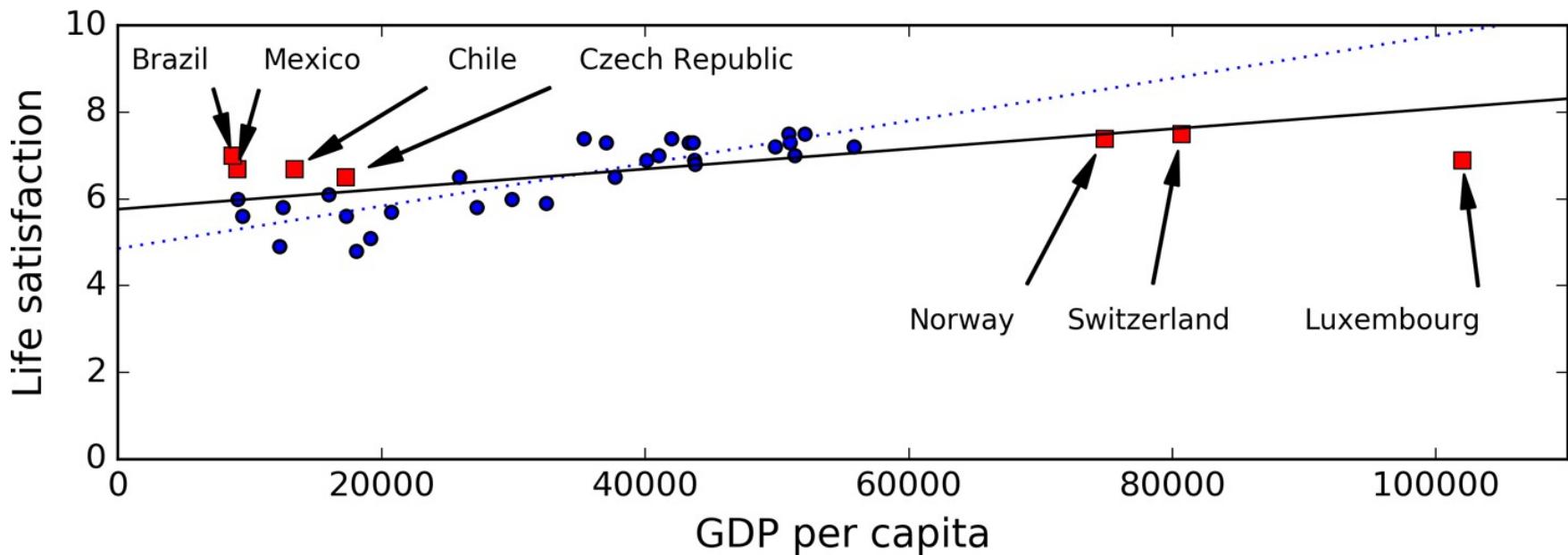
- Insufficient Quantity of Training Data



Desafios

- Nonrepresentative Training Data

A more representative training sample

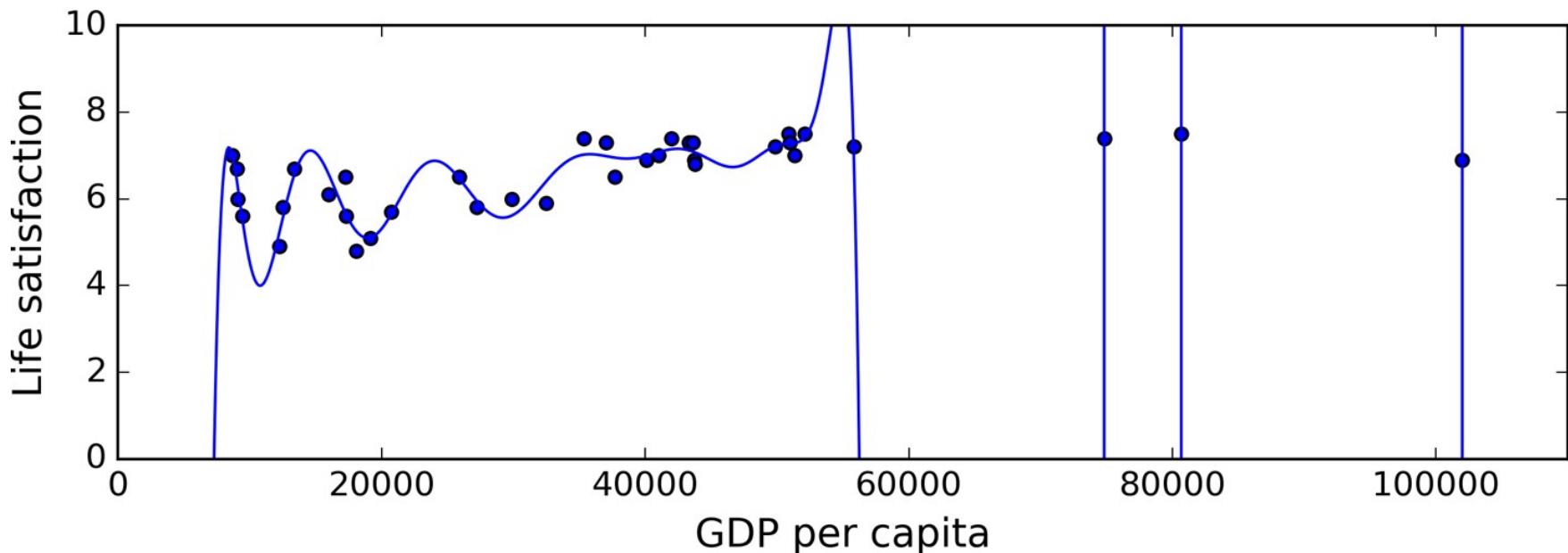


Desafios

- Poor-Quality Data
- Irrelevant Features
 - Saying: garbage in, garbage out
 - Feature engineering:
 - Feature selection: selecting the most useful features to train on among existing features.
 - Feature extraction: combining existing features to produce a more useful one (dimensionality reduction algorithms can help).
 - Creating new features by gathering new data.

Desafios

- Overfitting
 - the model performs well on the training data, but it does not generalize well.
 - Overfitting happens when the model is too complex relative to the amount and noisiness of the training data.

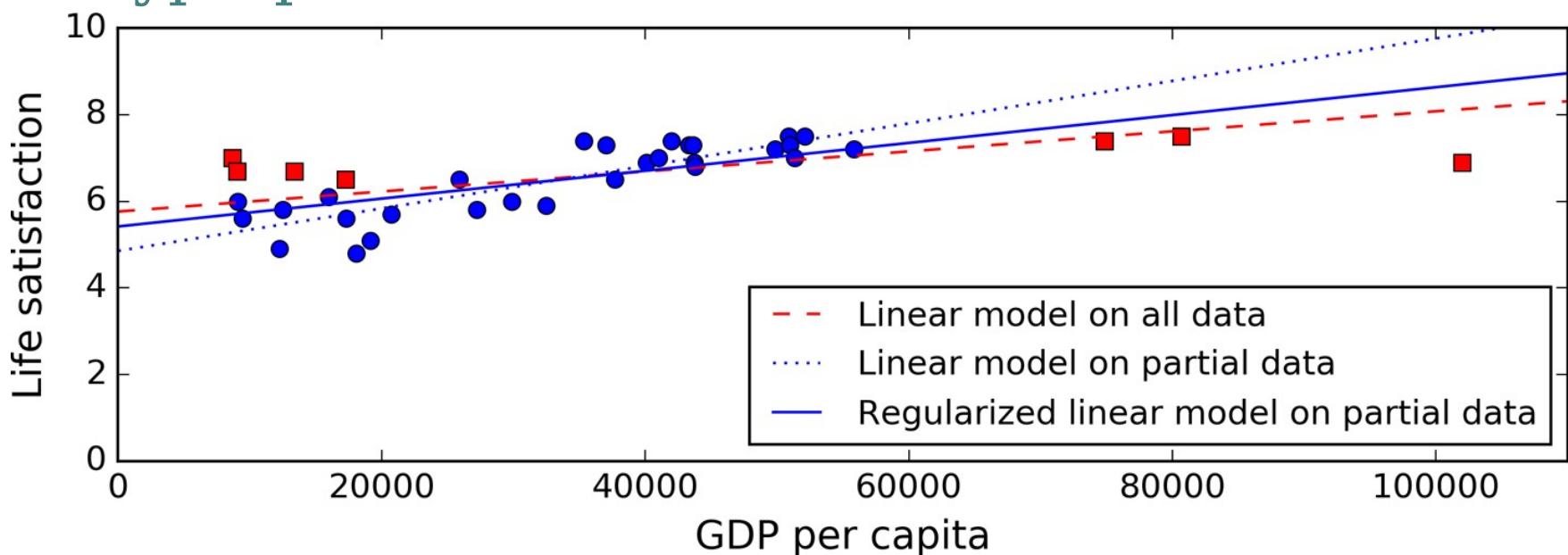


Desafios

- Overfitting – Soluções
 - To simplify the model by selecting one with fewer parameters, by reducing the number of attributes in the training data or by constraining the model (e.g. regularization constraint);
 - To gather more training data;
 - To reduce the noise in the training data (e.g., fix data errors and remove outliers).

Desafios

- Overfitting – Soluções
 - Using a regularization constraint hyperparameter



Desafios

- Underfitting
 - it occurs when your model is too simple to learn the underlying structure of the data.
 - reality is just more complex than the model, so its predictions are bound to be inaccurate, even on the training examples.
 - Soluções:
 - Selecting a more powerful model, with more parameters;
 - Feeding better features to the learning algorithm (feature engineering);
 - Reducing the constraints on the model (e.g., reducing the regularization hyperparameter).

Terminology and notations

Samples

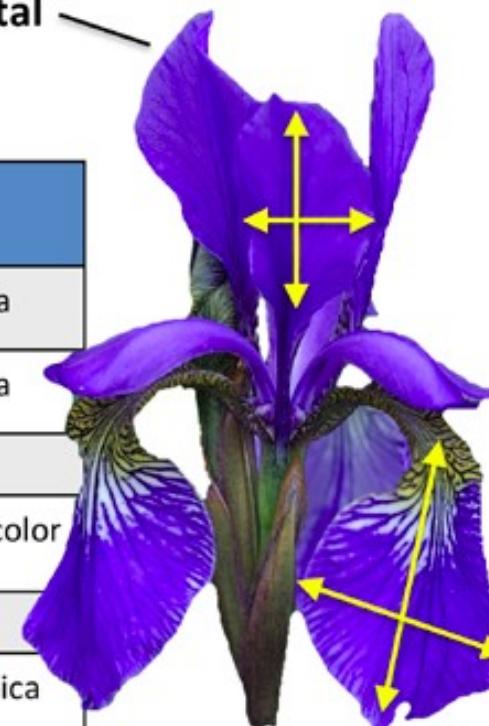
(instances, observations)

	Sepal length	Sepal width	Petal length	Petal width	Class label
1	5.1	3.5	1.4	0.2	Setosa
2	4.9	3.0	1.4	0.2	Setosa
...					
50	6.4	3.5	4.5	1.2	Versicolor
...					
150	5.9	3.0	5.0	1.8	Virginica

Features

(attributes, measurements, dimensions)

Petal



Sepal

Class labels
(targets)

Iris dataset

150 samples and four features:

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & x_4^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & x_4^{(2)} \\ \vdots & \vdots & \vdots & \vdots \\ x_1^{(150)} & x_2^{(150)} & x_3^{(150)} & x_4^{(150)} \end{bmatrix}$$

Row vector:

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_1^{(i)} & x_2^{(i)} & x_3^{(i)} & x_4^{(i)} \end{bmatrix}$$

Column vector:

$$\mathbf{x}_j = \begin{bmatrix} x_j^{(1)} \\ x_j^{(2)} \\ \vdots \\ x_j^{(150)} \end{bmatrix}$$

Iris dataset

Target variables (class labels) as a 150-dimensional column vector:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \dots \\ y^{(150)} \end{bmatrix} \left(y \in \{\text{Setosa, Versicolor, Virginica}\} \right)$$

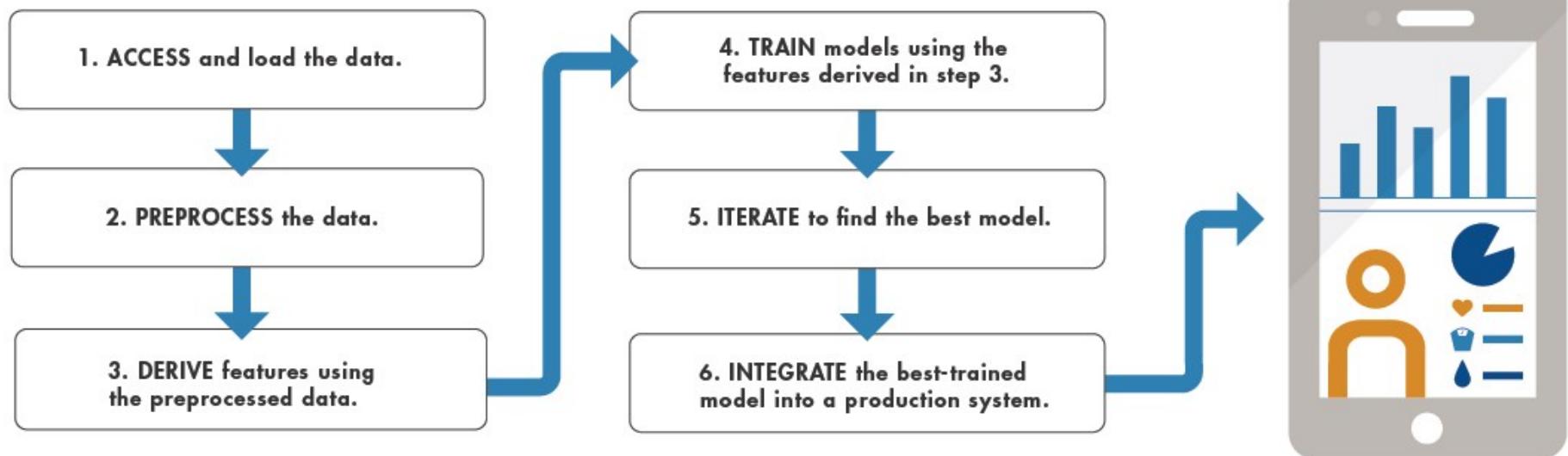
Workflow 1



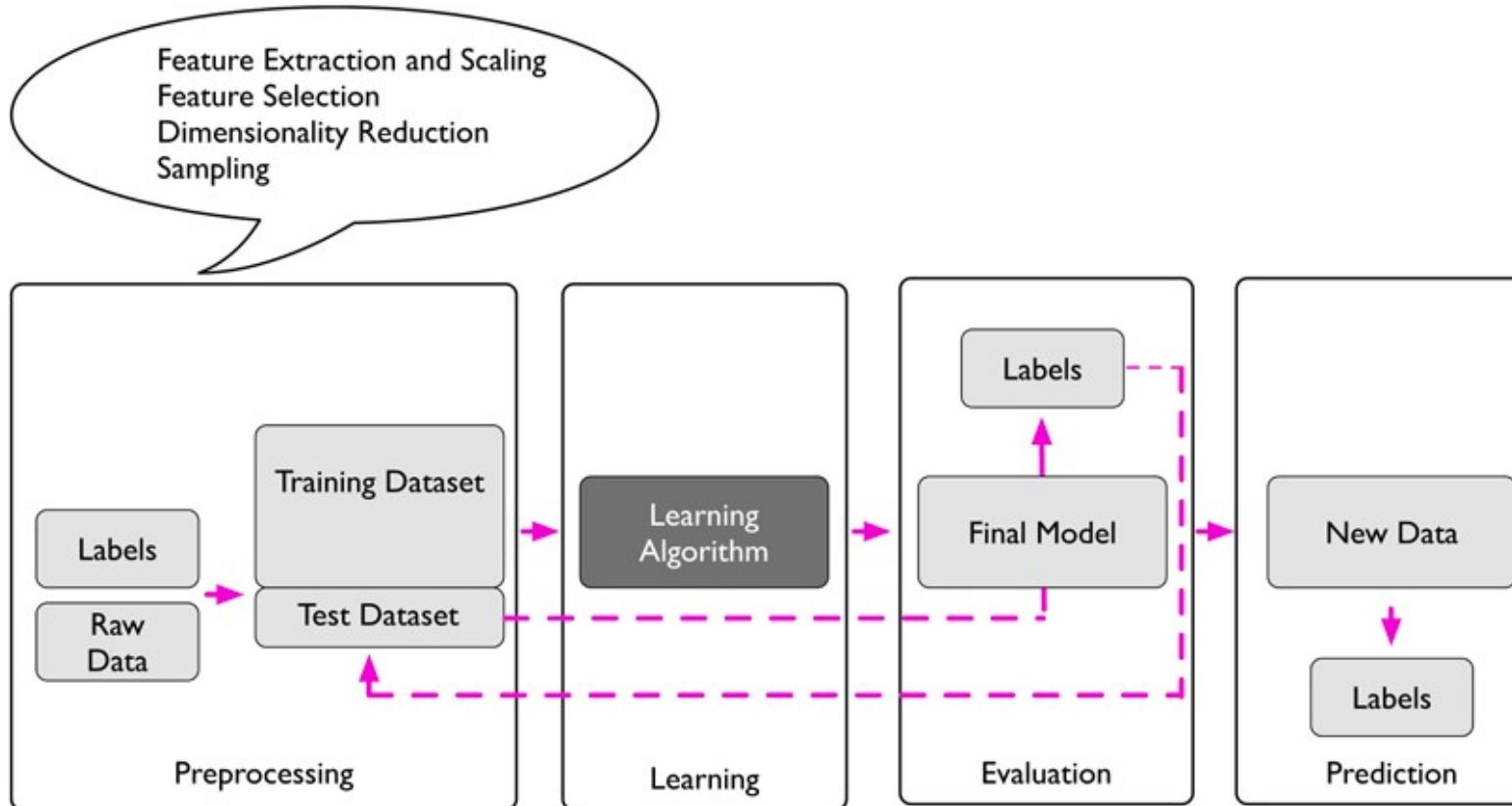
1. Definição do problema
2. Análise Exploratória dos Dados
3. Preparar os dados
4. Avaliar Algoritmos
5. Melhorar os resultados
6. Apresentar os resultados



Workflow 2



Roadmap for building machine learning systems



getting data into shape

Model Selection
Cross-Validation
Performance Metrics
Hyperparameter Optimization

Ecossistema Python



Projeto Ciência de Dados na Prática (cdp)

Material didático gratuito

- GitHub
 - <http://bit.ly/cdpgithub>
- YouTube
 - <http://bit.ly/cdpvideos>
- Facebook
 - <http://bit.ly/cdpface>



This repository page displays the following information:

- 49 commits
- 1 branch
- 0 releases
- 2 contributors

Recent activity:

- Merge branch 'master' of https://github.com/ciencia-de-dados-pratica/... by jvitorc17 - 16 days ago
- Adicionando alterações by minicurso-data-science - 27 days ago
- added tutorials by praticas - 4 months ago
- Initial commit by .gitignore - 8 months ago
- Adicionando introdução ao pandas by 01-introducao-ciencia-dados.ipynb - 6 months ago
- Adicionando introdução ao pandas by 02-instalacao_do_ambiente.ipynb - 6 months ago
- new notebooks by 03-jupyter_notebook.ipynb - 4 months ago
- new notebooks by 04-introducao-python-parte-1.ip... - 4 months ago

The YouTube channel page for 'Ciência de Dados Prática' shows the following details:

- 39 subscribers
- Navigation menu: HOME, VIDEOS, PLAYLISTS, CHANNELS, DISCUSSION, ABOUT
- Popular uploads:
 - #24 - Análise de dados Parte I (Ciência de Dados Prática) - 3 views • 3 days ago
 - #23 - Introdução a Análise de Dados (Ciência de Dados Prática) - 2 views • 5 days ago
 - #22 - Iniciando Análise de Dados a partir de dados de Dados (Ciência de Dados Prática) - 1 view • 1 week ago
- Introdução ao ambiente Jupyter Notebook (71 views • 5 months ago)
- Introdução a Análise de Dados a partir de dados de Dados (71 views • 5 months ago)

Obrigado!
Dúvidas, comentários, sugestões?

Regis Pires Magalhães
regismagalhaes@ufc.br



UFC