



Ciência de Dados Quântica 22/23

Classical Machine Learning: Regression & Classification

LUÍS PAULO SANTOS

Material de Consulta

- ▶ Maria Schuld, Francesco Petruccione
“Machine Learning with Quantum Computers”
Cap. 2

- ▶ Luís G. Serrano
“grokking Machine Learning”
Caps. 2, 3, 4, 5, 6, 7 e apêndice B

An attempt on definition

“Machine learning can be seen as one particular flavour of AI, which **designs computer algorithms** that generalise **from patterns found in data** to make predictions in unknown situations.”

“The attraction of machine learning, therefore, lies in its capability to learn from **extremely large data sets** that only computers can process.”

[Schuld2021 – chapter 2]

An attempt on definition

“Machine learning **starts** with a very **general, agnostic mathematical model**, and **uses data to adapt it** to the case, a process called “**training**”. “

When **looking at the final model, we do not necessarily gain information on the physical mechanism or statistical relationships behind the data**, but consider it as a **black box** that has learned to **produce reliable outputs**.

[Schuld2021 – chapter 2]

Contrast this with **Explainable AI (XAI)**: https://en.wikipedia.org/wiki/Explainable_artificial_intelligence

Fundamental concepts

► **Problem:**

defined by an **unknown** function f^* that relates every point x in the **domain** \mathcal{X} to a value y in the **codomain** \mathcal{Y} :

$$f^*: \mathcal{X} \rightarrow \mathcal{Y}$$

► **Regression:**

y is continuous

► **Classification:**

y is discrete

Fundamental concepts

► Data set:

set of points \mathcal{D} in either \mathcal{X} or $\mathcal{X} \otimes \mathcal{Y}$ available to the machine learning algorithm

- ▶ $\mathcal{D} \in \mathcal{X} = \{x^1, \dots, x^M\}$: unlabelled data
- ▶ $\mathcal{D} \in \mathcal{X} \otimes \mathcal{Y} = \{(x^1, y^1), \dots, (x^M, y^M)\}$: labelled data

► Features:

\mathcal{X} is n -dimensional, meaning that $x \in \mathcal{X}$ is a vector with n elements

$$x = (x_1, x_2, \dots, x_n) \in \mathcal{X}$$

Fundamental concepts

- ▶ **Model:**

$f: \mathcal{X} \rightarrow \mathcal{Y}$, which “approximates” f^*

- ▶ **Model Family:**

$\mathcal{F} = \{f_\theta\}$, with $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$ and

$\theta = \{\theta_0, \theta_1, \dots, \theta_T\}$ is a T -dimensional vector that parameterizes f_θ

- ▶ **Loss:**

$L(f_\theta(x), f^*(x) = y)$ quantifies the disparity between the model prediction $f_\theta(x)$ and the true value $y = f^*(x)$

Fundamental concepts

- ▶ **Risk (aka generalization error):**

The risk of a model f_θ is the expected value of the loss over the whole domain \mathcal{X} :

$$R_{f_\theta} = \int_{\mathcal{X}} L(f_\theta(x), f^*(x)) dx$$

- ▶ **ML goal:**

find the model f_{θ^*} which minimizes the generalization error over \mathcal{X} :

$$f_{\theta^*} = \operatorname{argmin}_{f_\theta \in \mathcal{F}} \int_{\mathcal{X}} L(f_\theta(x), f^*(x)) dx$$

Fundamental concepts

- ▶ Empirical Risk :

the expected value of f_θ 's loss over \mathcal{D} :

$$\hat{R}_{f_\theta}(\mathcal{D}) = \frac{1}{M} \sum_{i=1}^M L(f_\theta(x^i), y^i)$$

- ▶ ML's practical goal:

find the model \hat{f}_θ which minimizes the empirical risk $\hat{R}_{f_\theta}(\mathcal{D})$:

$$\hat{f}_\theta = \operatorname{argmin}_{f_\theta \in \mathcal{F}} \sum_{i=1}^M L(f_\theta(x^i), y^i)$$

Classification of ML

- ▶ Supervised Learning : labelled data set
- ▶ Unsupervised Learning: unlabelled data set
- ▶ Reinforcement Learning: interactive based on feedback

Types of ML - Supervised Learning

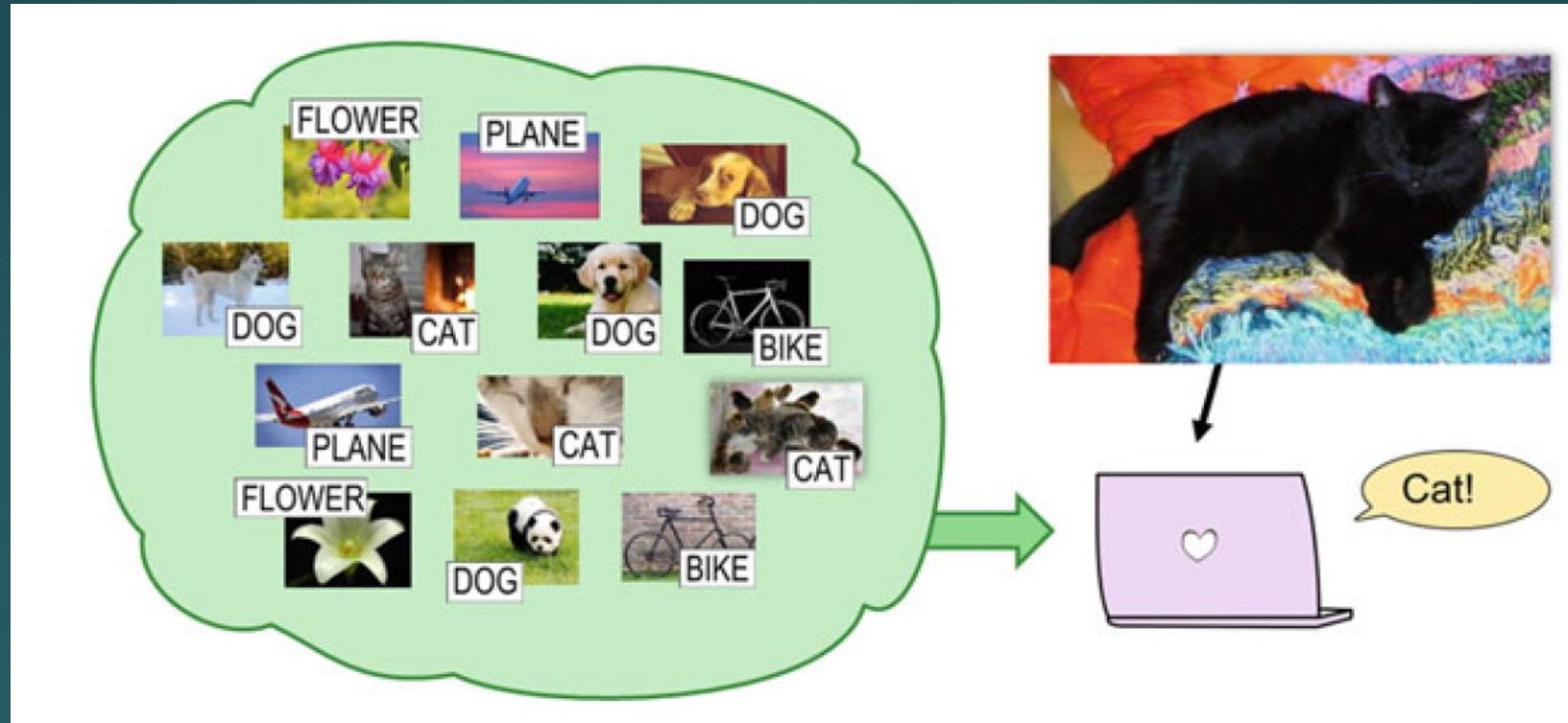
11

► Definition:

Given a **labelled data set** $\mathcal{D} \in \mathcal{X} \otimes \mathcal{Y} = \{(x^1, y^1), \dots, (x^M, y^M)\}$

as well as a **new unclassified** point $x \in \mathcal{X}$,
predict the corresponding output $y \in \mathcal{Y}$.

Supervised Learning Examples



[Schuld2021 – chapter 2]

Supervised Learning Examples

► Time Series

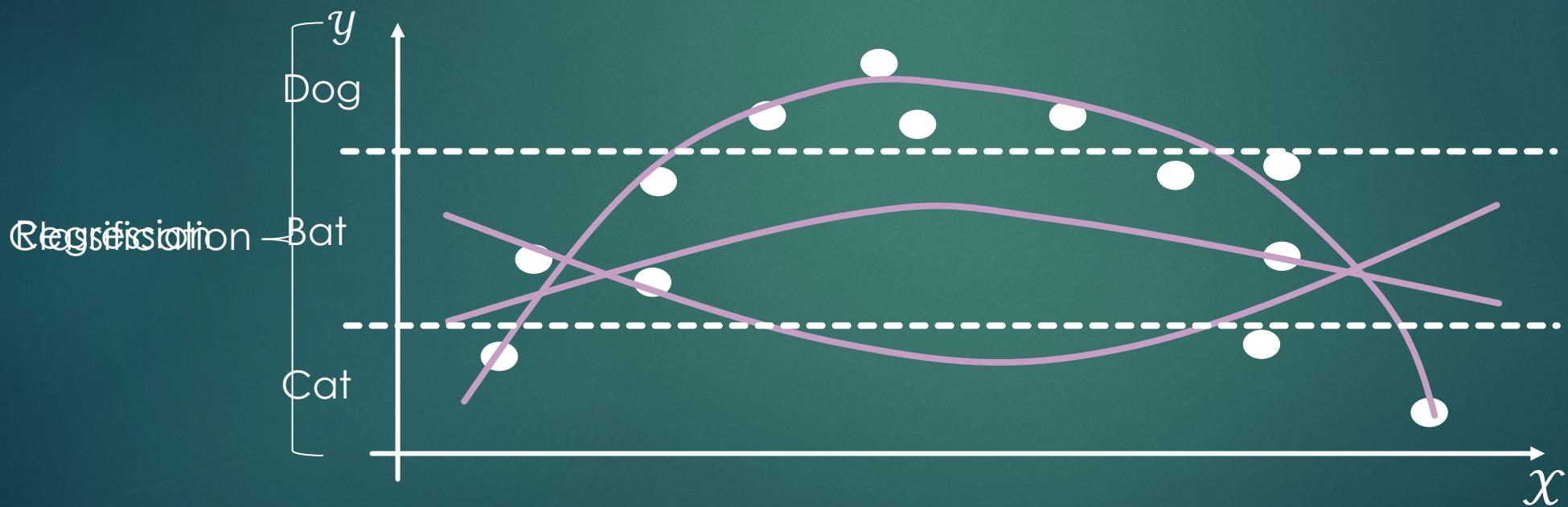
What will be the value of C
for time=20, given A and B ?



Supervised Learning

14

- ▶ Learns a model from the training data set which is used to predict unseen data points



Supervised Learning Algorithms

- ▶ Support-vector machines.
- ▶ Linear regression.
- ▶ Random forest.
- ▶ Naive Bayes.
- ▶ Linear discriminant analysis.
- ▶ Decision trees.
- ▶ K-nearest neighbor algorithm.
- ▶ Artificial Neural Networks

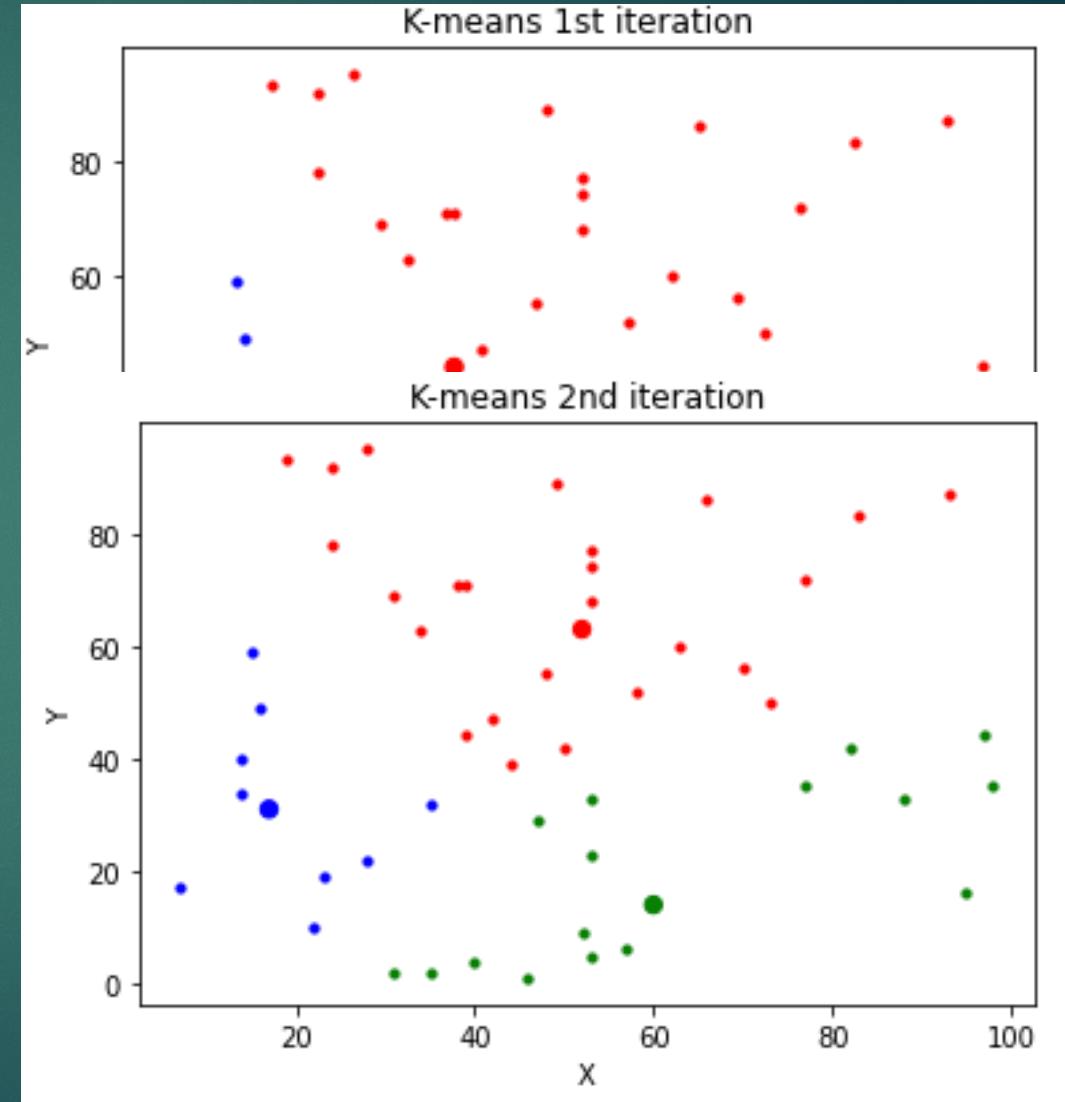
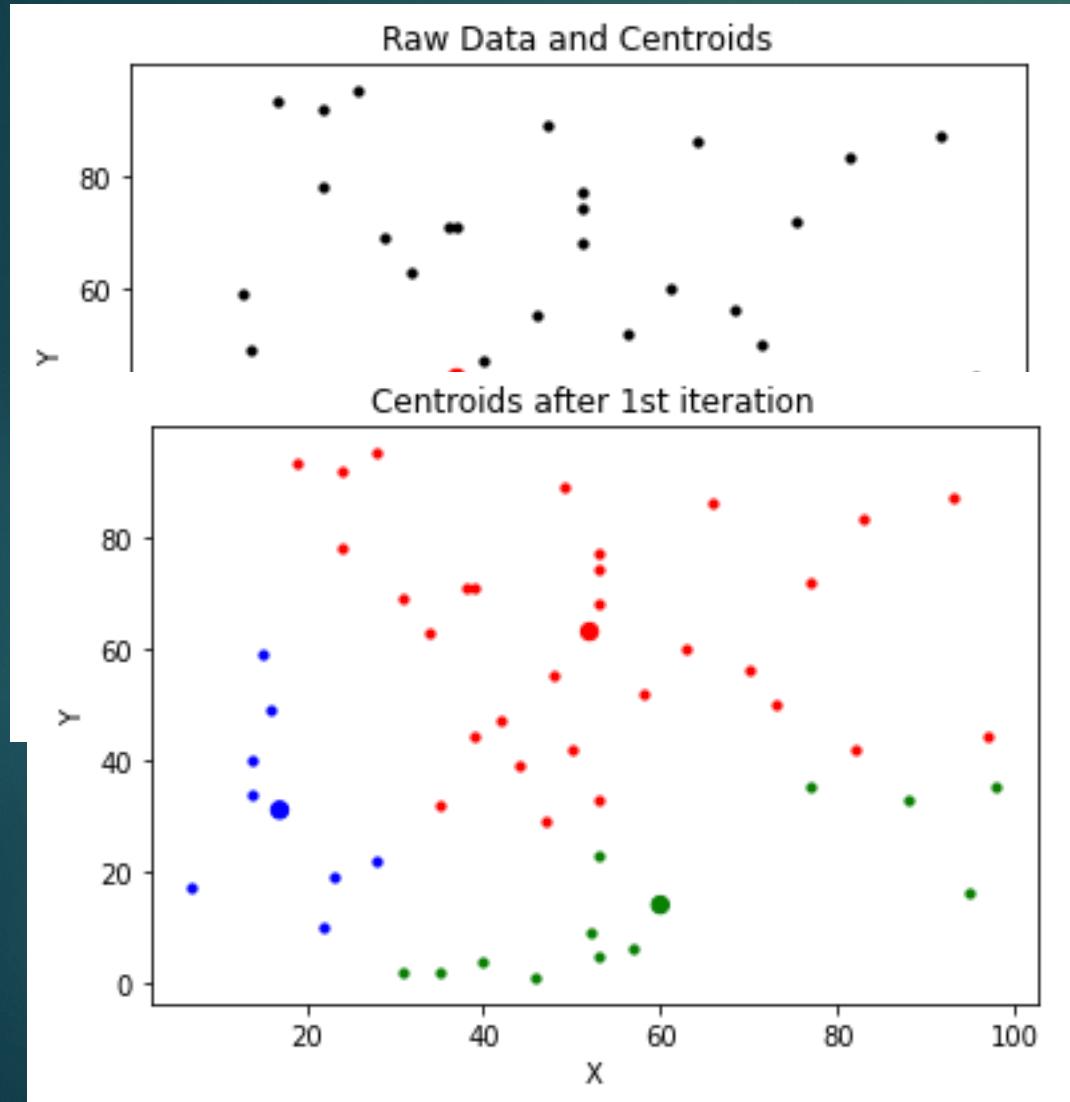
► **Definition:**

Given an **unlabelled data set** $\mathcal{D} \in \mathcal{X} = \{x^1, \dots, x^M\}$ drawn from a probability distribution $p(x)$:

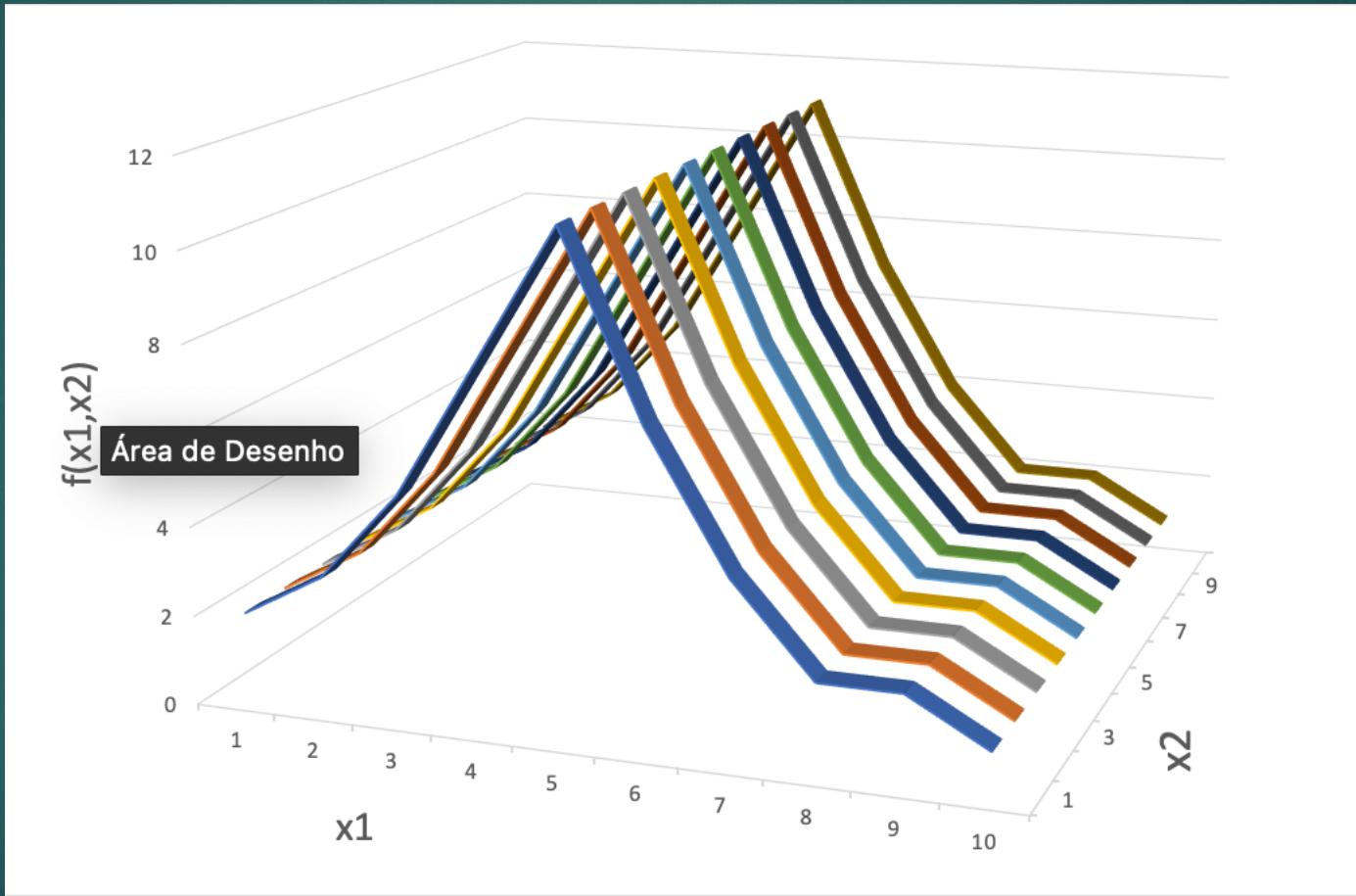
- ▶ group the points according to some similarity measure – clustering
- ▶ eliminate/combine features in \mathcal{X} – dimensionality reduction
- ▶ **draw a new sample** from $p(x)$ – generative approach

Unsupervised Learning - Clustering

17



Unsupervised Learning – Dimensionality Reduction



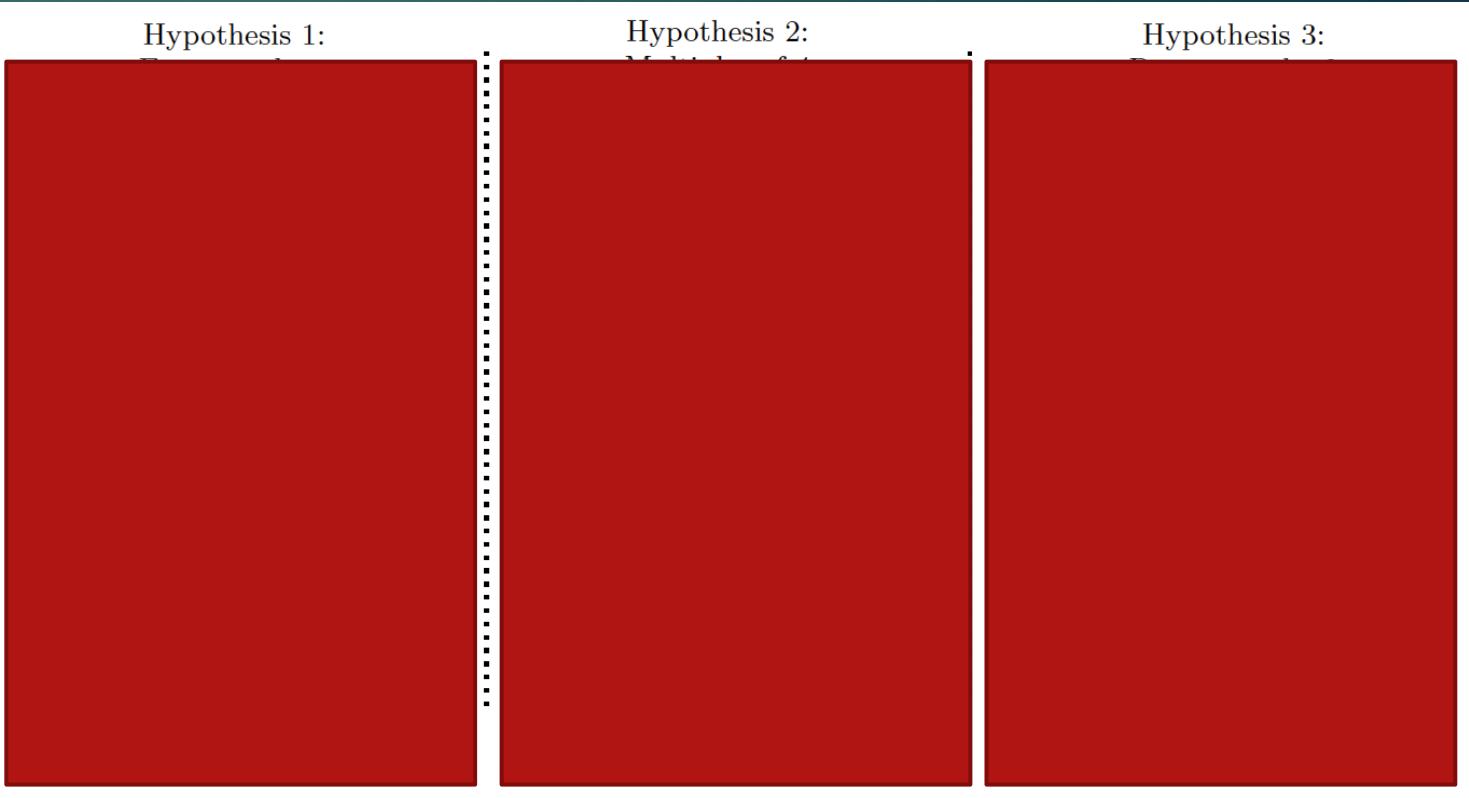
Unsupervised Learning - generative

19

► Hypothesis guessing:

In the range $i \in \{1 \dots 100\}, i \in \mathbb{N}$,
the values $\{4, 16, 36, 100\}$ are
given according to a
(unknown) rule.

Generate another sample with
the same rule.



Unsupervised Learning - generative

20

DALL-E



A photo of a white fur monster standing in a purple room

Chat GPT

podes escrever a receita de bacalhau à Brás no estilo de Luís de Camões?

"Oh, bacalhau, peixe digno de um prato
Que no mar do Norte é encontrado
E agora, em nossa terra, será trato
Com o tempero e o azeite profundo.

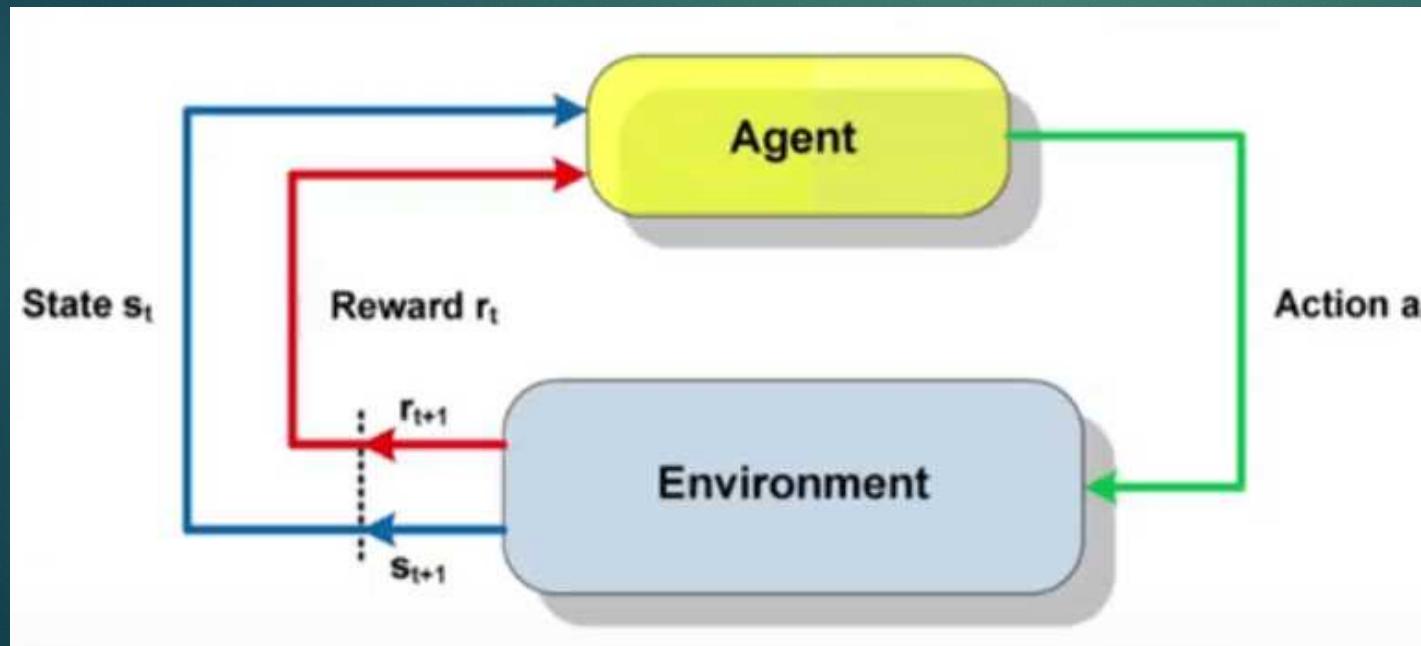
Com batatas, cebolas, alho e pimentões
Se misturarão em uma panela fervente
E em seu interior, ovos serão batidos
Até se tornarem cozidos.

E assim, a mistura será apresentada
Com sabor forte e aroma que encanta
E quem a experimentar, será conquistado
Pela iguaria de virtude tanta.

Sal e pimenta, temperos indispensáveis
Comprarão ainda mais brilho à refeição
E farão que o bacalhau, verdadeiro rei
Seja degustado com alegria e devoção."

Types of ML - Reinforcement Learning

- ▶ The agent learns good strategies by trial and error.
Every time it acts on the environment it receives a reward.
A good strategy maximizes the total received reward.
- ▶ The agent might learn a policy, indicating, for each state, which action will maximize the total (expected) reward.



Regression

Supervised Learning: Regression

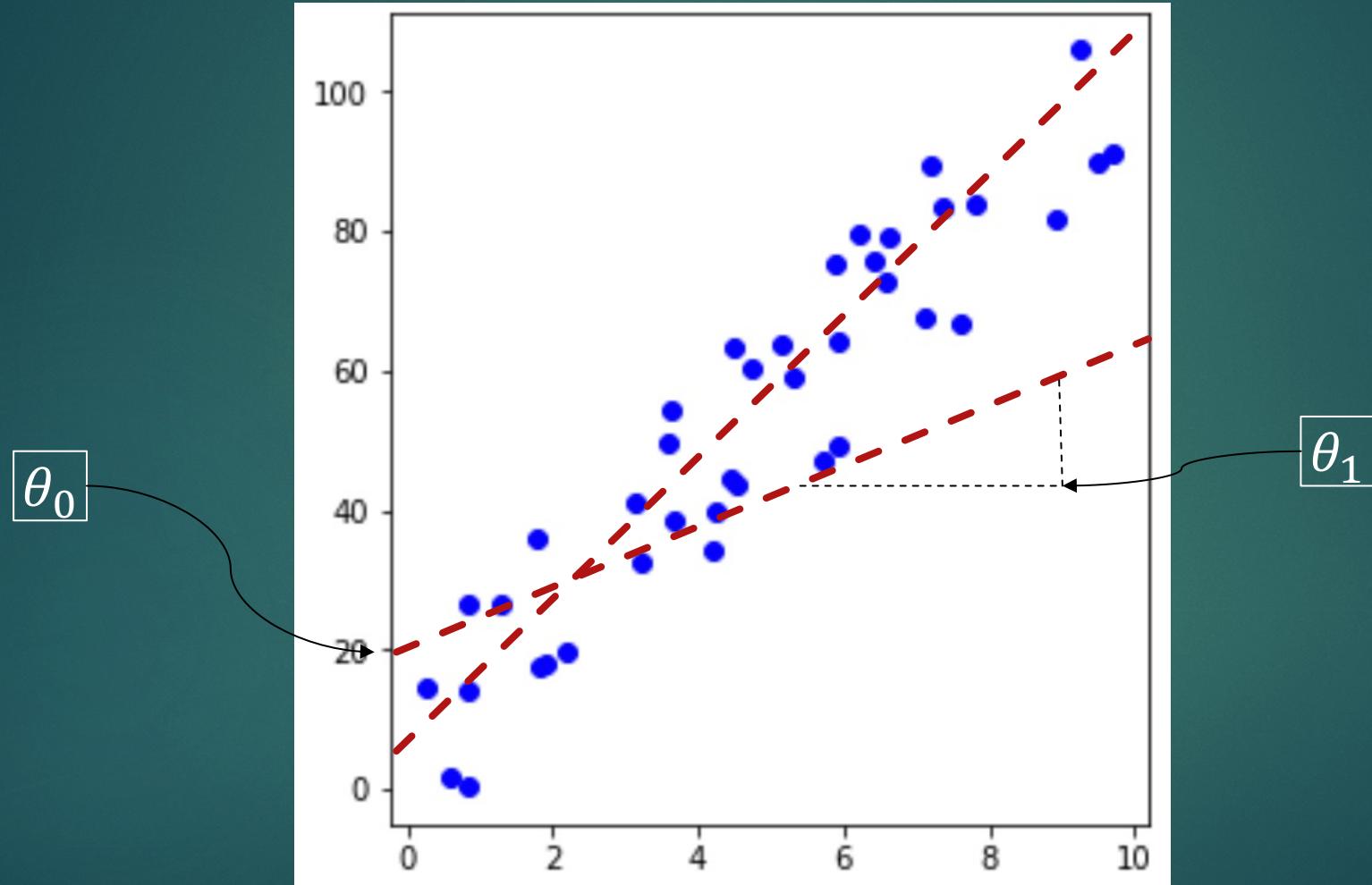
23

- ▶ Problem:
given a labelled dataset $\mathcal{D} \in \mathcal{X} \otimes \mathcal{Y}$, taken from an unknown function $f^*: \mathcal{X} \rightarrow \mathcal{Y}$, **learn** a function $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$ which approximates f^* and allows predicting $y = \hat{f}(x)$ for previously unseen (x, y) , i.e. $(x, y) \notin \mathcal{D}$
- ▶ Learnable functions f_θ are parameterized with θ being a vector of parameters. **For linear regression** $y = f_\theta(x) = \theta_0 + \sum_{i=1}^n \theta_i * x_i$, with the number of parameters $T = n + 1$, n being the number of features
- ▶ For uni-dimensional data ($n=1$) and linear regression lines are used:

$$y = f_\theta(x) = \theta_0 + \theta_1 * x_1$$

Supervised Learning: Regression

24

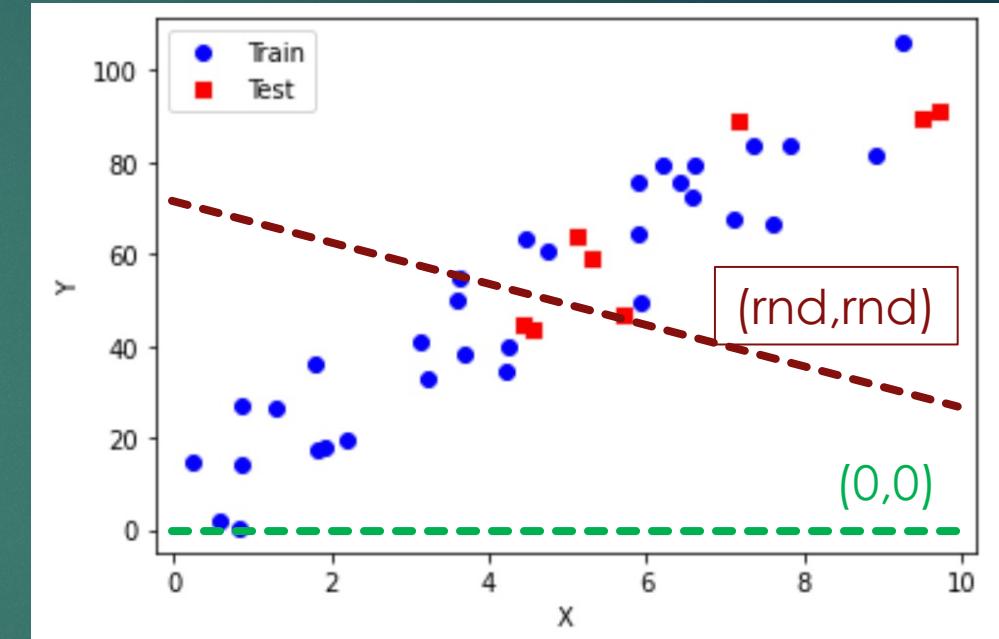
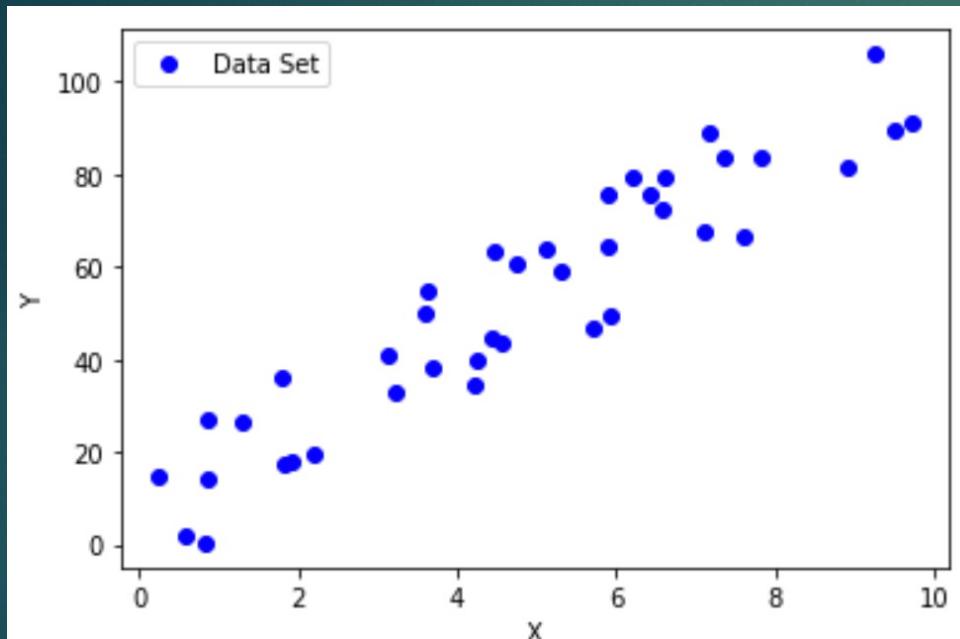


Linear Regression: algorithm

1. `Dtrain, Dtest = separate (D, fraction)`
2. `θ = Initialize ()`
3. Repeat many times
 1. take a subset D_i from D_{train}
 2. $\theta_{new} = update_theta (D_i, \theta)$
 3. $\theta = \theta_{new}$
4. verify (D_{test}, θ)

Linear Regression: algorithm

1. Dtrain, Dtest = separate (D, 0.8)



2. $\theta = \text{Initialize}()$

Optimization: gradient descent

```
theta_new = update_theta (D, theta)
```

- ▶ The goal is to learn a model that minimizes the empirical risk $\hat{R}_{f_\theta}(\mathcal{D})$:

$$\hat{f}_\theta = \operatorname{argmin}_{f_\theta \in \mathcal{F}} \sum_{i=1}^M L(f_\theta(x^i), y^i)$$

- ▶ The loss quantifies the distance between each prediction and the true value. Let's use the square error as the example:

$$L(f_\theta(x^i), y^i) = (f_\theta(x^i) - y^i)^2 = (\theta_0 + \theta_1 * x^i - y^i)^2$$

- ▶ The goal is to find θ that minimizes the objective function:

$$\sum_{i=1}^M (\theta_0 + \theta_1 * x^i - y^i)^2$$

Optimization: gradient descent

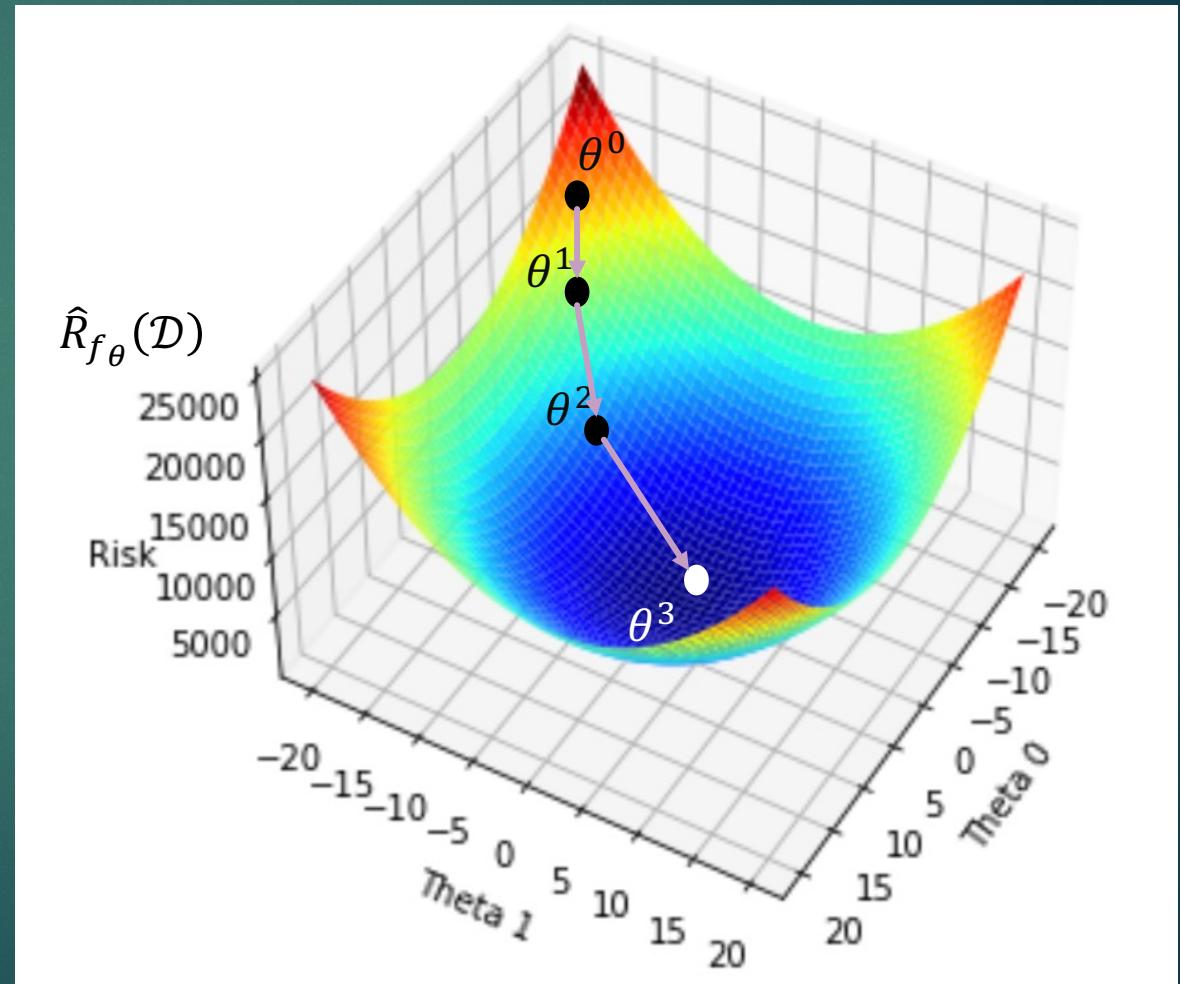
`theta_new = update_theta (D_i, theta)`

- The goal is to find θ that minimizes the objective function:

$$\hat{R}_{f_\theta}(\mathcal{D}) = \sum_{i=1}^M (\theta_0 + \theta_1 * x^i - y^i)^2$$

- At each iteration change θ slightly along the direction of maximum change

$$\theta^{t+1} = \theta^t - \eta * \nabla \sum_{i=1}^M (\theta_0^t + \theta_1^t * x^i - y^i)^2$$



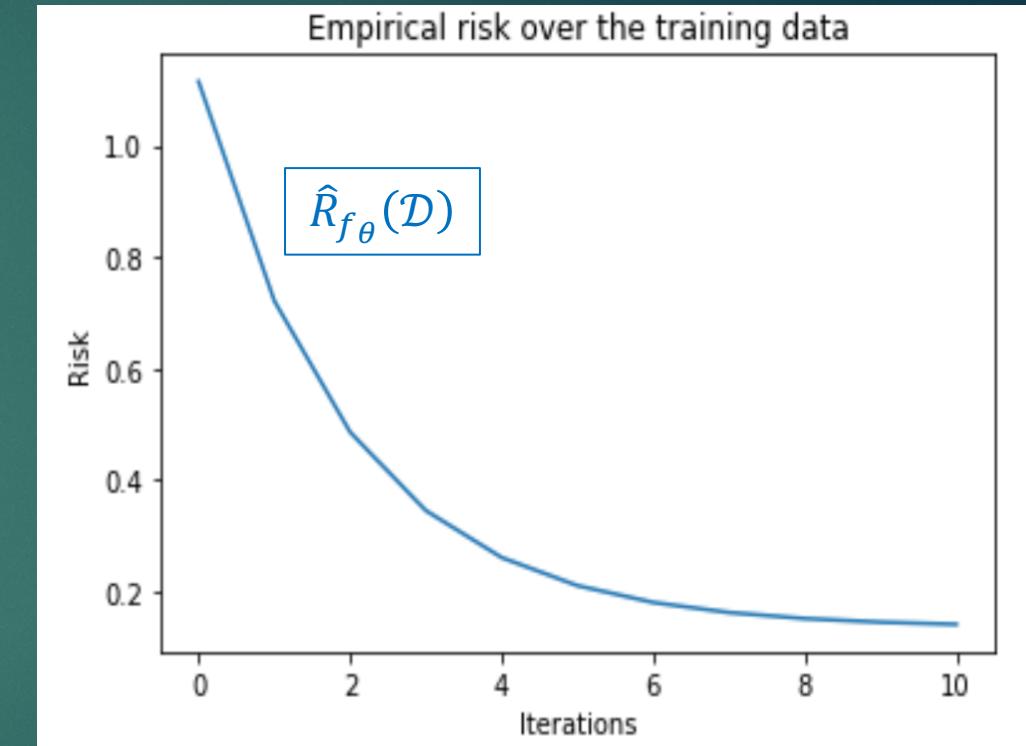
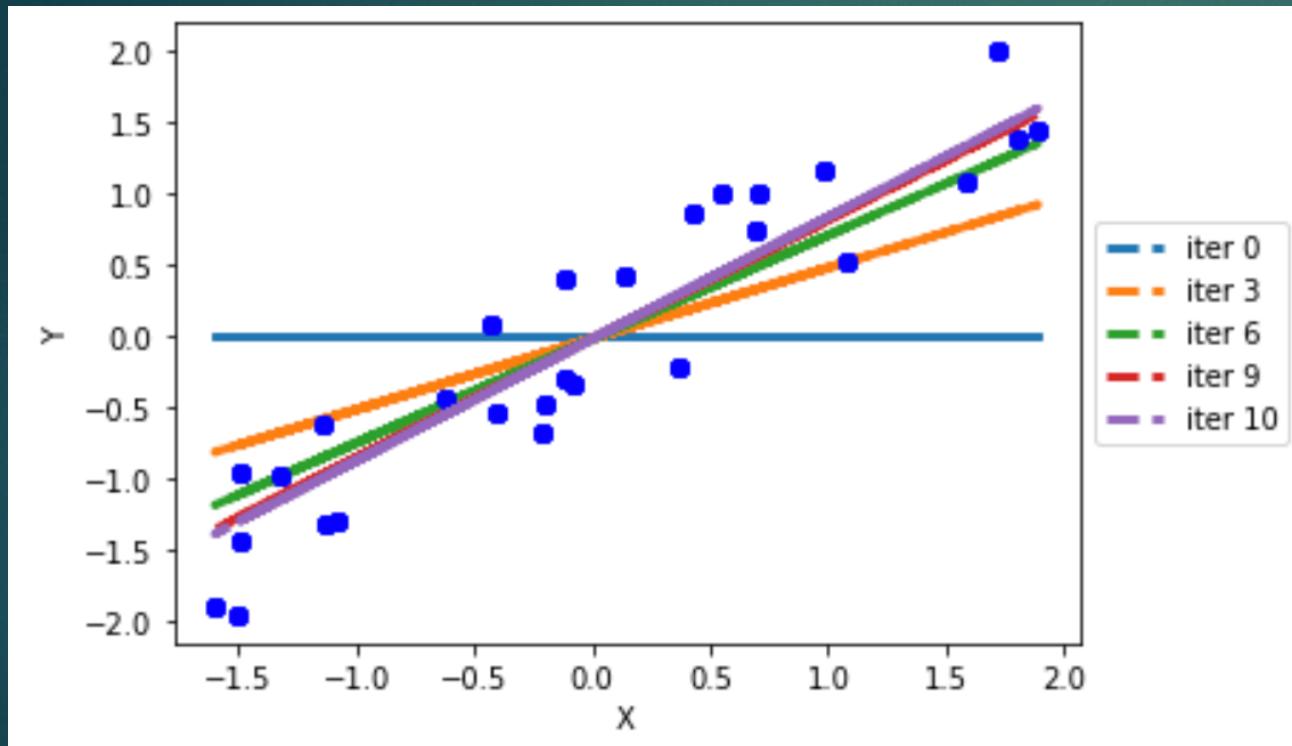
Optimization: gradient descent

29

$$\nabla \sum_{i=1}^M (\theta_0^t + \theta_1^t * x^i - y^i)^2 = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \end{pmatrix} = \begin{pmatrix} 2 \sum_{i=1}^M (\theta_0^t + \theta_1^t * x^i - y^i) \\ 2 \sum_{i=1}^M (\theta_0^t + \theta_1^t * x^i - y^i) * x^i \end{pmatrix}$$

$$= \begin{pmatrix} 2 \sum_{i=1}^M (f_{\theta^t}(x^i) - y^i) \\ 2 \sum_{i=1}^M (f_{\theta^t}(x^i) - y^i) * x^i \end{pmatrix}$$

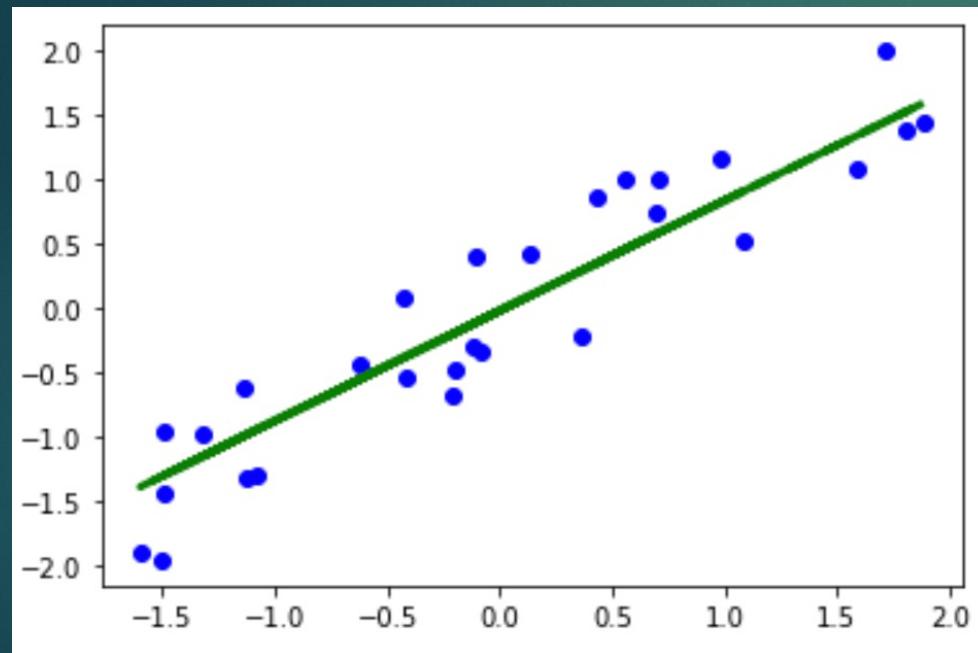
Linear Regression



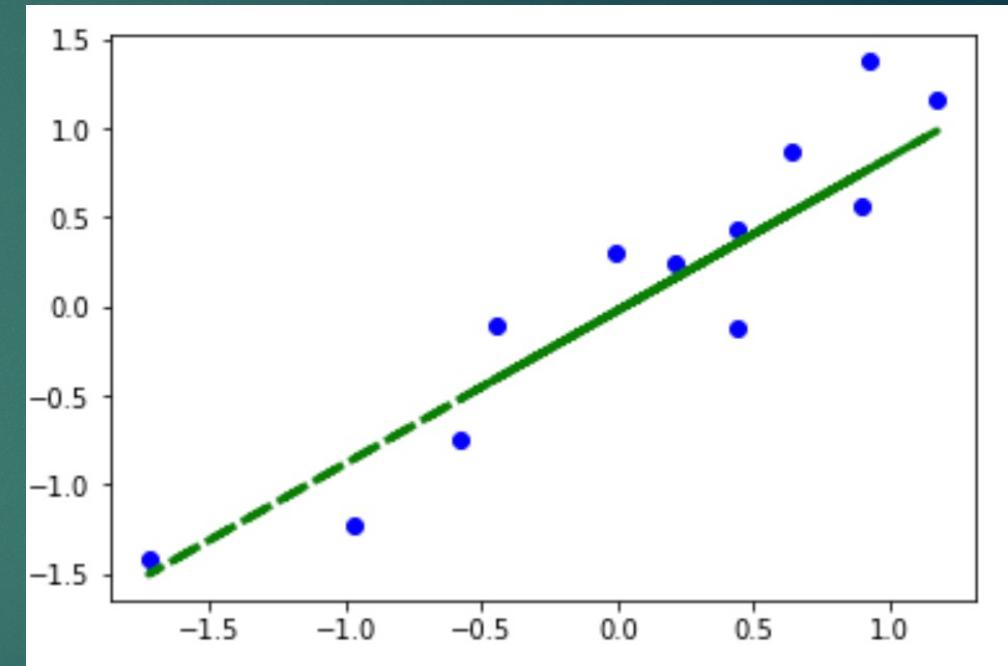
Linear Regression: testing

4. verify (D_{test} , θ)

train data set



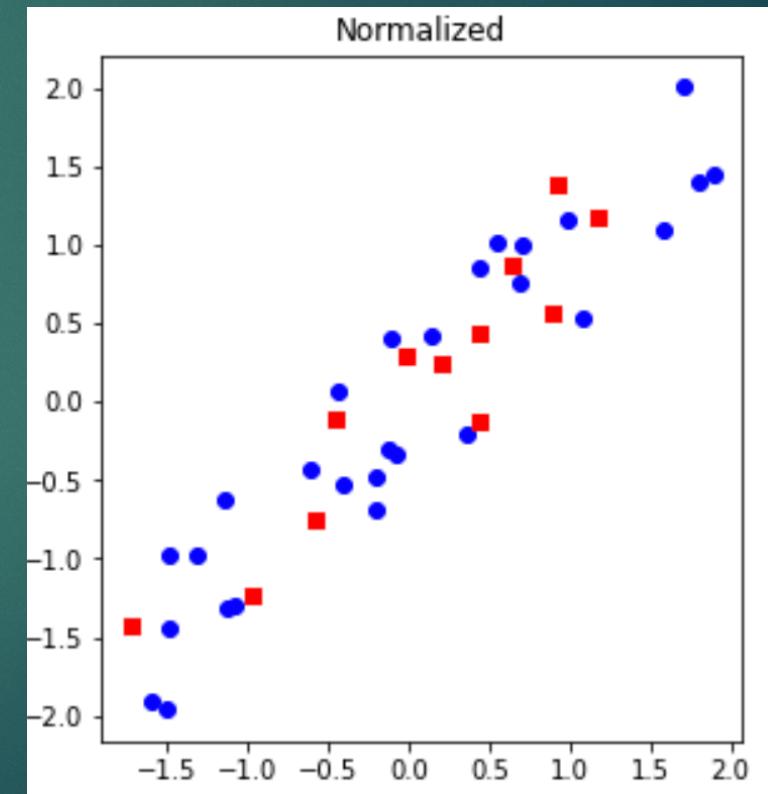
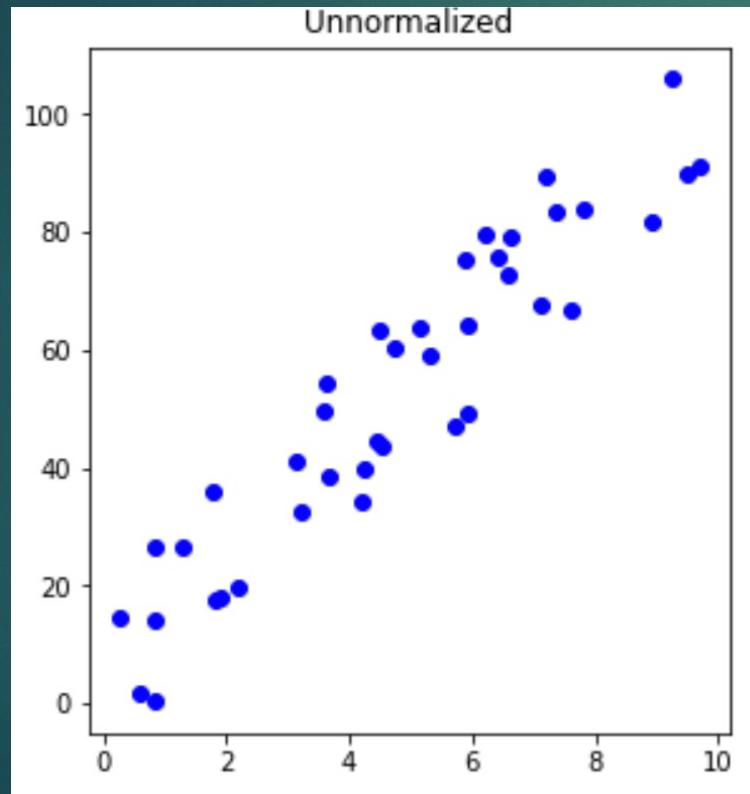
test data set



MSE = 0.0989

Linear Regression: data shifting and normalization

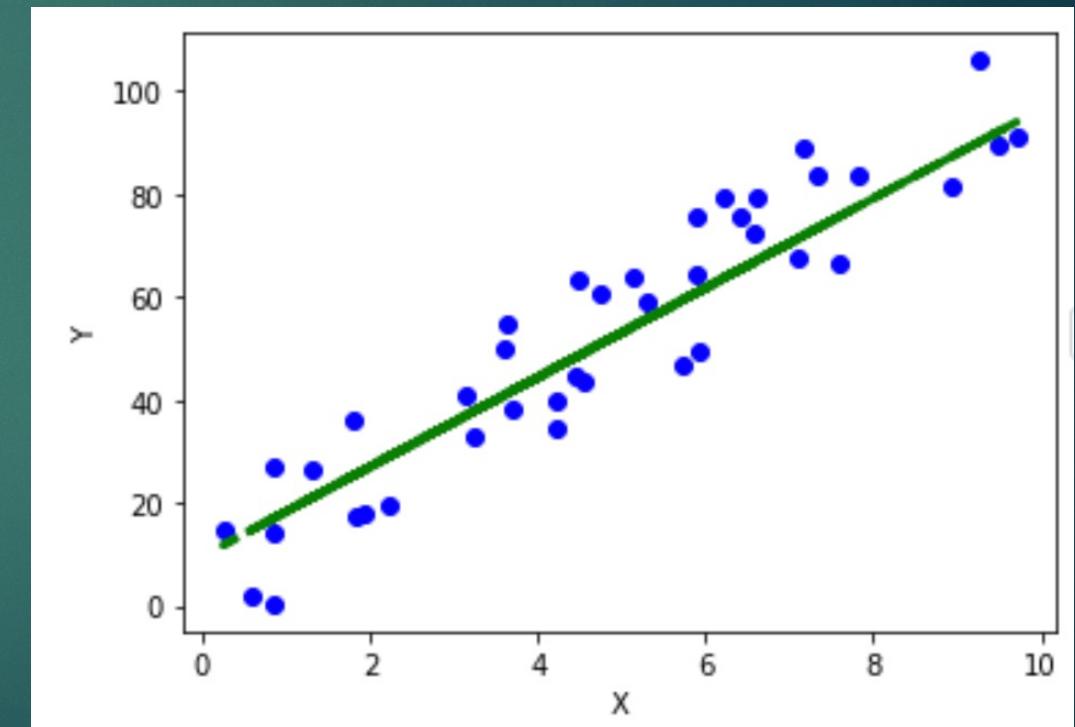
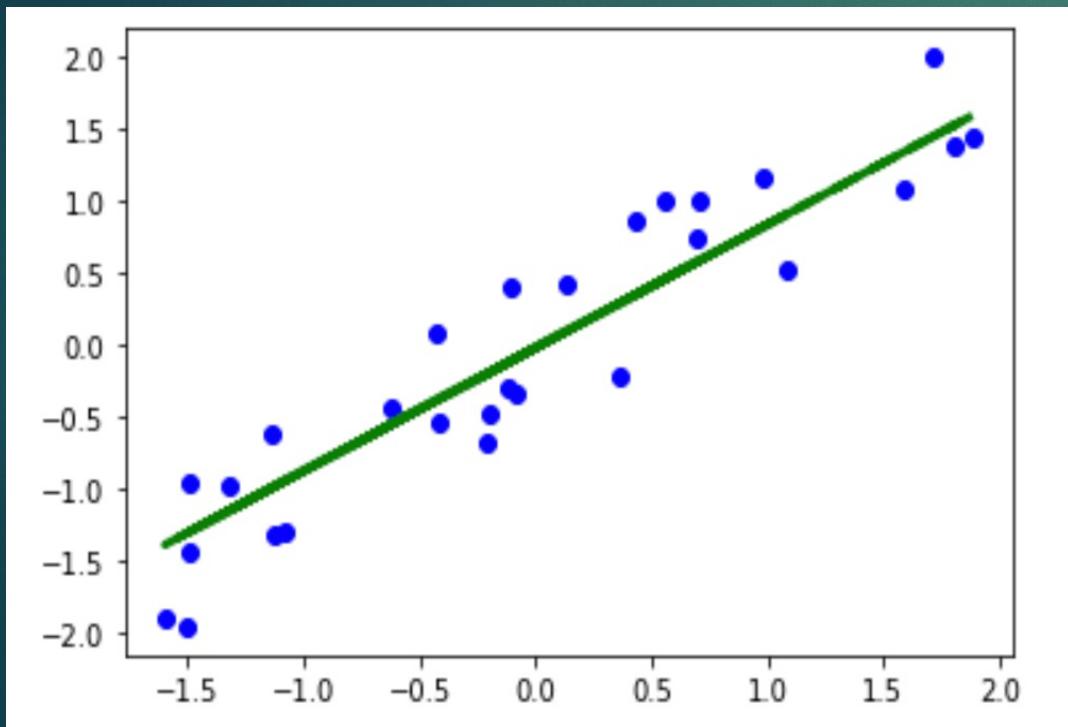
$$x = \frac{(x - \bar{x})}{\sigma_x} \quad y = \frac{(y - \bar{y})}{\sigma_y} \quad \forall (x, y) \in \mathcal{D}$$



Linear Regression: data shifting and normalization

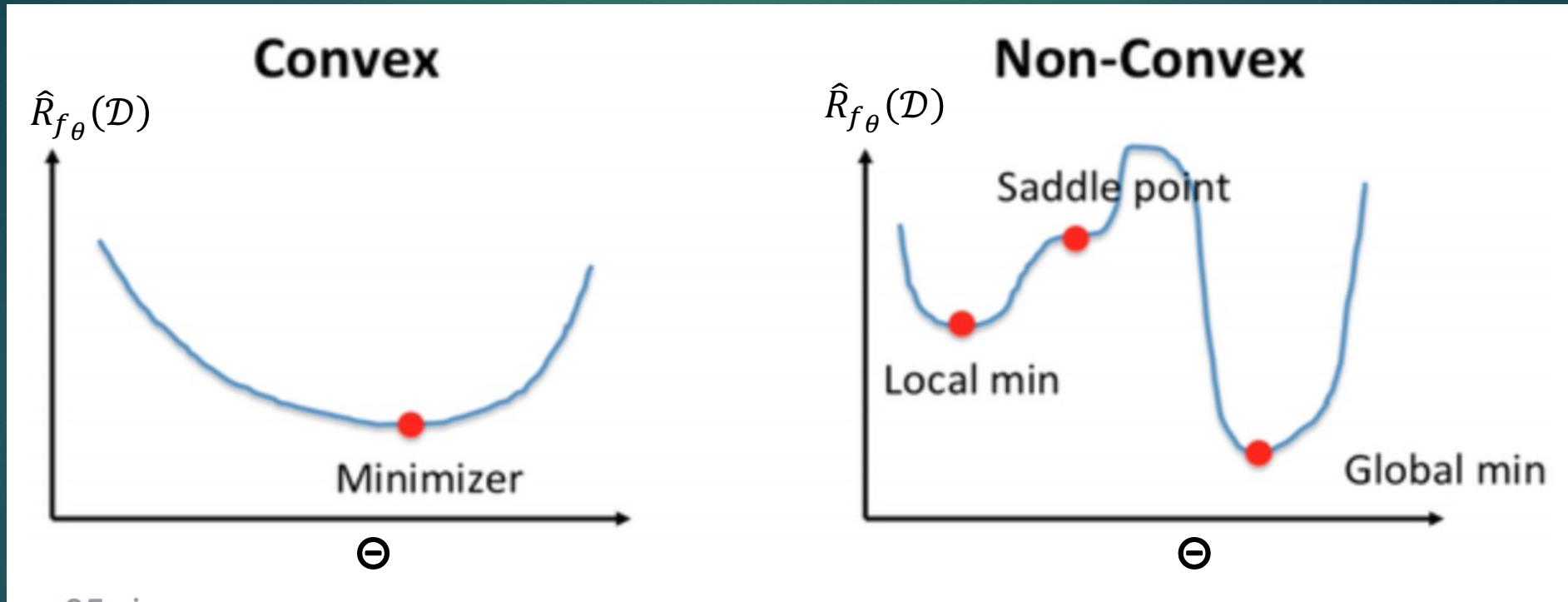
$$\theta_0 = \frac{\sigma_y(\theta_{0,n} - \theta_{1,n} * \bar{x})}{\sigma_x} + \bar{y}$$

$$\theta_1 = \frac{\sigma_y * \theta_{1,n}}{\sigma_x}$$



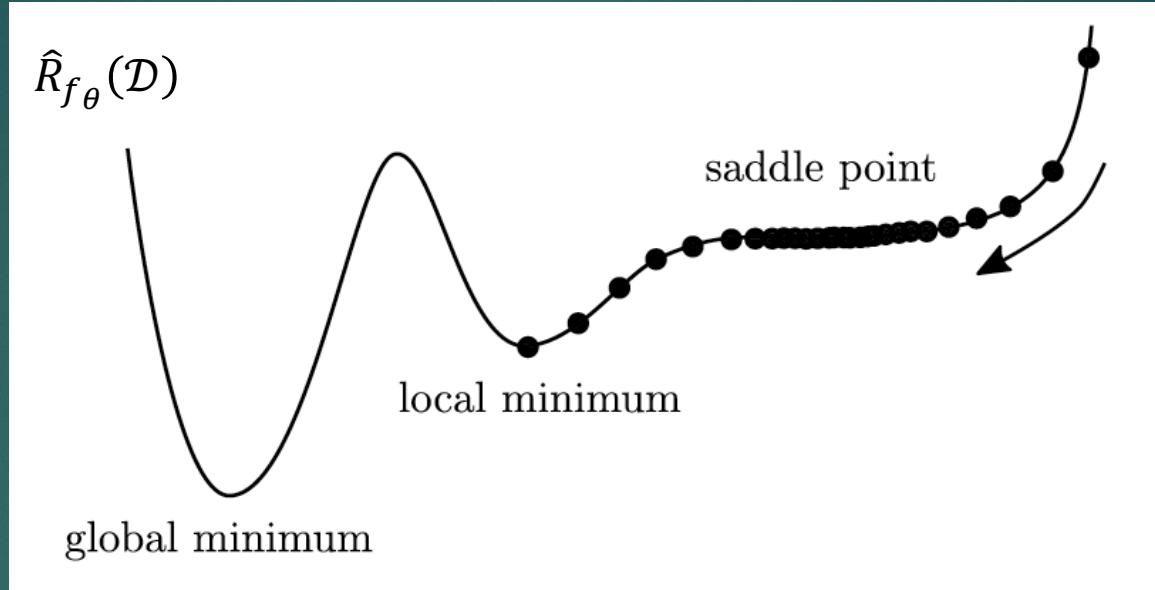
Optimization: convexity and local minima

Convexity: non-convex objective functions are much harder to optimize



Many machine learning problems are high-dimensional (large number of parameters) non convex

Optimization: local minima



- ▶ clever initialization of θ
- ▶ random restart
- ▶ clever choice of the learning rate η
- ▶ stochastic gradient descent: use a stochastic subset of the training data at each step
- ▶ natural gradient: use geometric information on the curvature of the landscape

Multivariate Linear Regression

► Loss

$$L(f_{\theta}(x), y) = (f_{\theta}(x) - y)^2 = \left(\theta_0 + \sum_{i=1}^n (\theta_i * x_i) - y \right)^2$$

► Empirical Risk

$$\hat{R}_{f_{\theta}}(\mathcal{D}) = \frac{1}{M} \sum_{j=1}^M L(f_{\theta}(x^{(j)}), y^{(j)}) = \frac{1}{M} \sum_{j=1}^M \left(\theta_0 + \sum_{i=1}^n (\theta_i * x_i^{(j)}) - y^{(j)} \right)^2$$

► Gradient

$$\nabla \hat{R}_{f_{\theta}}(\mathcal{D}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_{T-1}} \end{pmatrix} = \begin{pmatrix} 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) \\ 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) * x_1^{(j)} \\ \vdots \\ 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) * x_{T-1}^{(j)} \end{pmatrix}$$

Polynomial Regression

► Loss

$$L(f_{\theta}(x), y) = (f_{\theta}(x) - y)^2 = \left(\theta_0 + \sum_{i=1}^{T-1} (\theta_i * x^i) - y \right)^2$$

x to the power of i

► Empirical Risk

$$\hat{R}_{f_{\theta}}(\mathcal{D}) = \frac{1}{M} \sum_{j=1}^M L(f_{\theta}(x^{(j)}), y^{(j)}) = \frac{1}{M} \sum_{j=1}^M \left(\theta_0 + \sum_{i=1}^{T-1} (\theta_i * x^{(j)i}) - y^{(j)} \right)^2$$

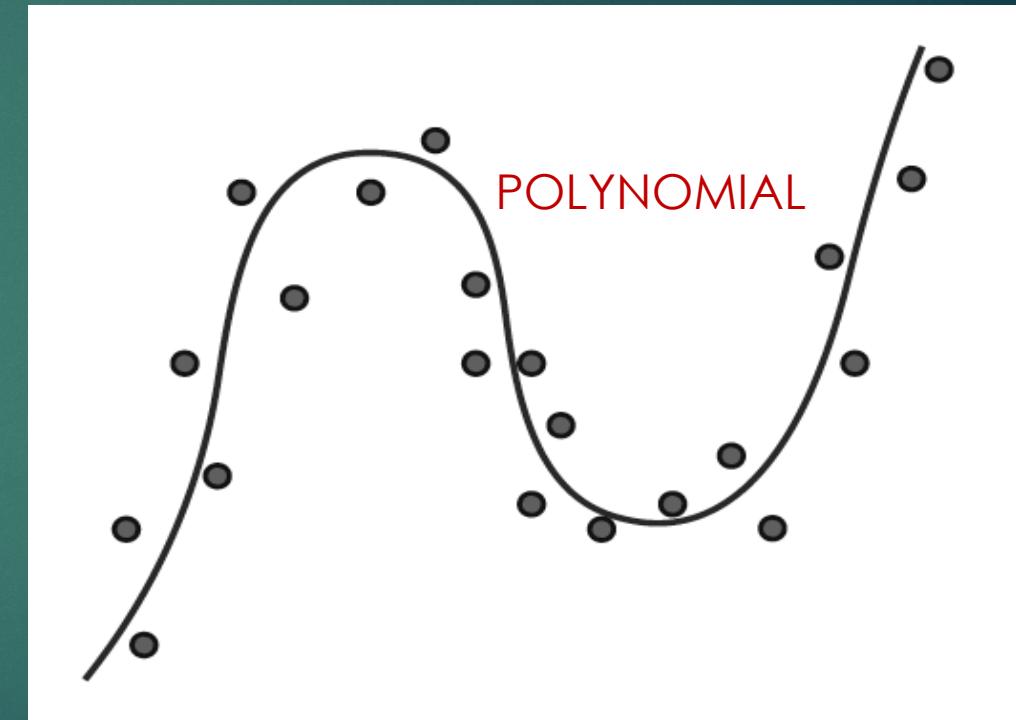
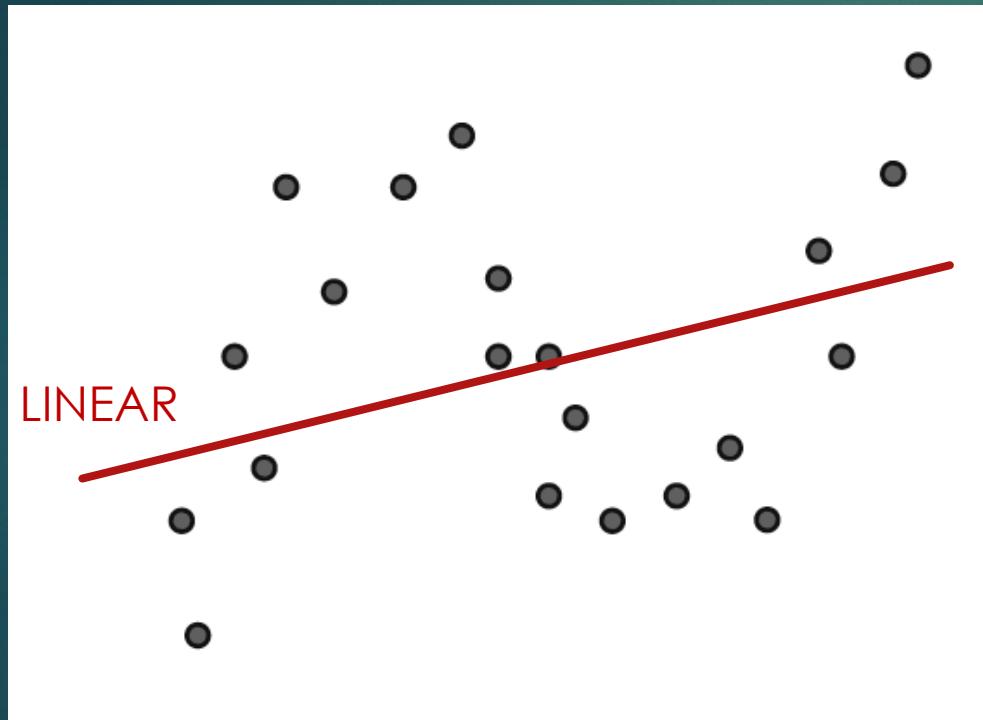
x^(j) to the power of i

► Gradient

$$\nabla \hat{R}_{f_{\theta}}(\mathcal{D}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_{T-1}} \end{pmatrix} = \begin{pmatrix} 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) \\ 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) * x^{(j)1} \\ \vdots \\ 2 \sum_{j=1}^M (f_{\theta}(x^{(j)}) - y^{(j)}) * x^{(j)T-1} \end{pmatrix}$$

Polynomial Regression

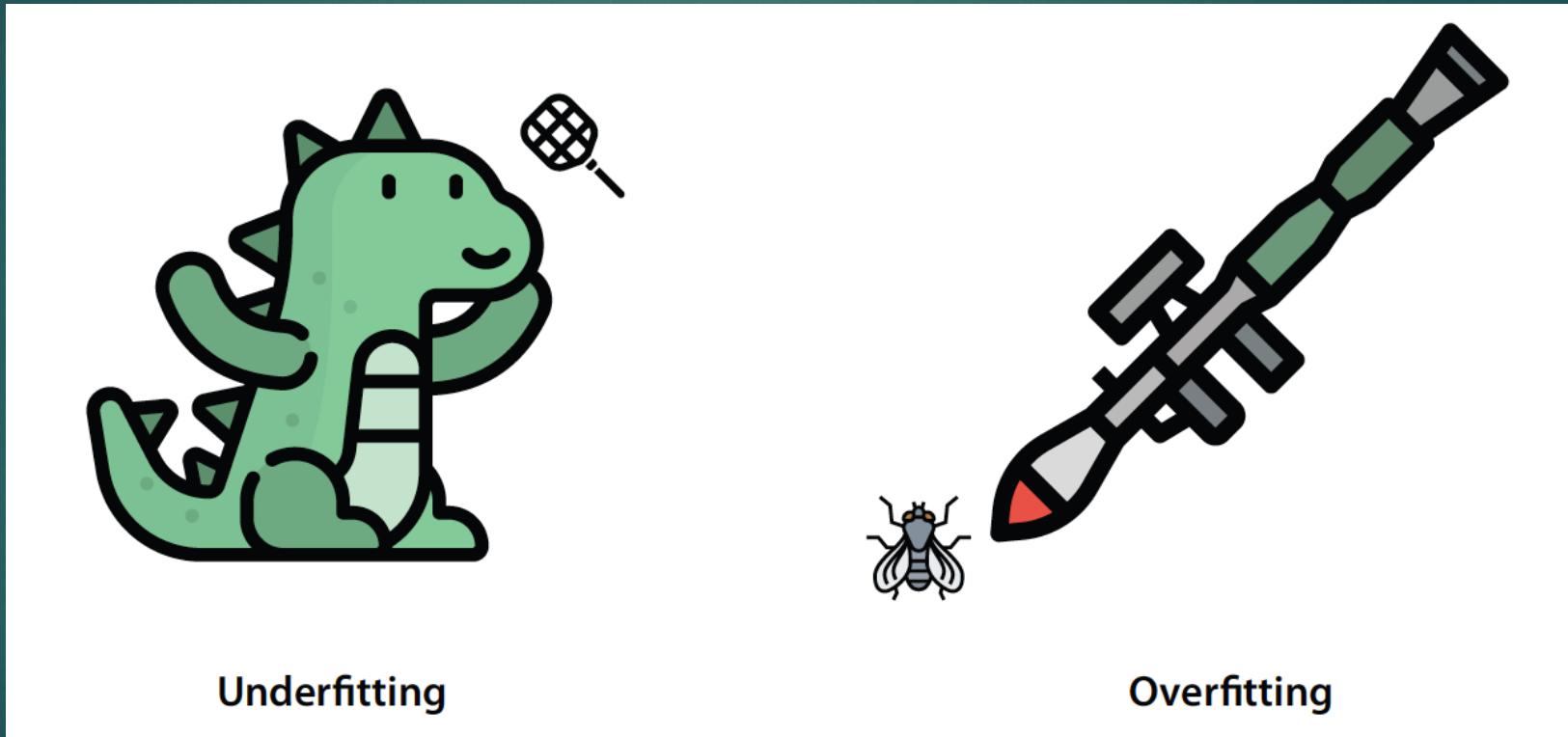
38



Model Learning: under and over fitting

39

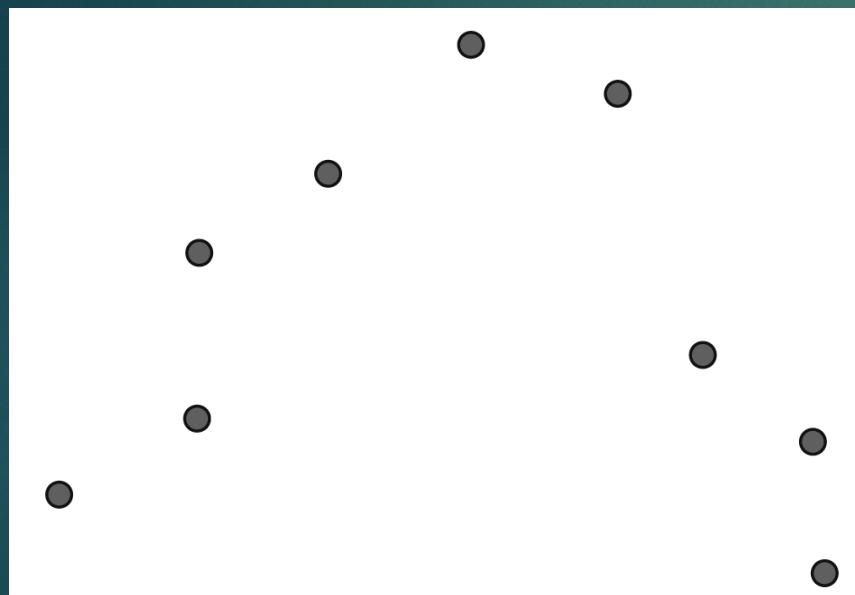
- ▶ **underfitting** occurs when a **too simple** model is used
- ▶ **overfitting** occurs when a **too complex** model is used



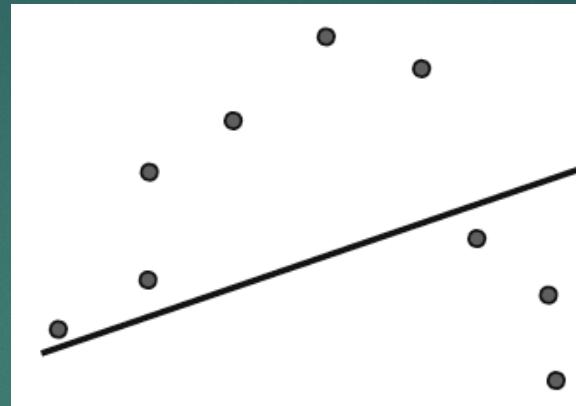
[Serrano2021]

Model Learning: under and over fitting

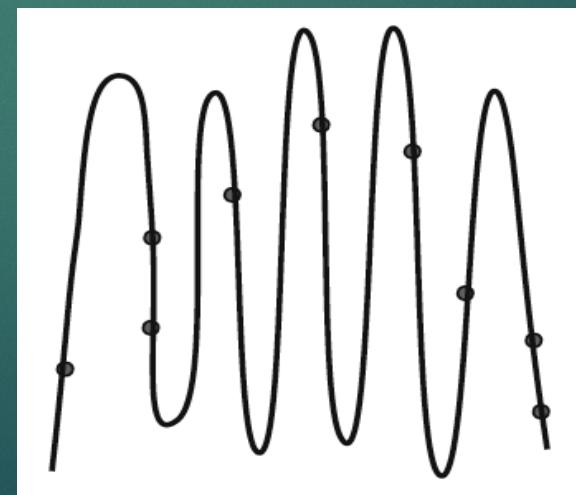
40



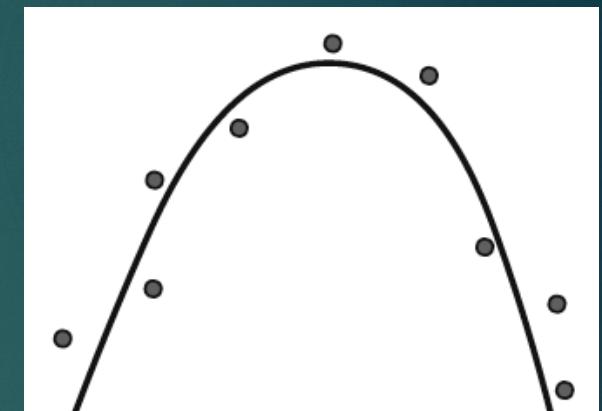
[Serrano2021]



LINEAR



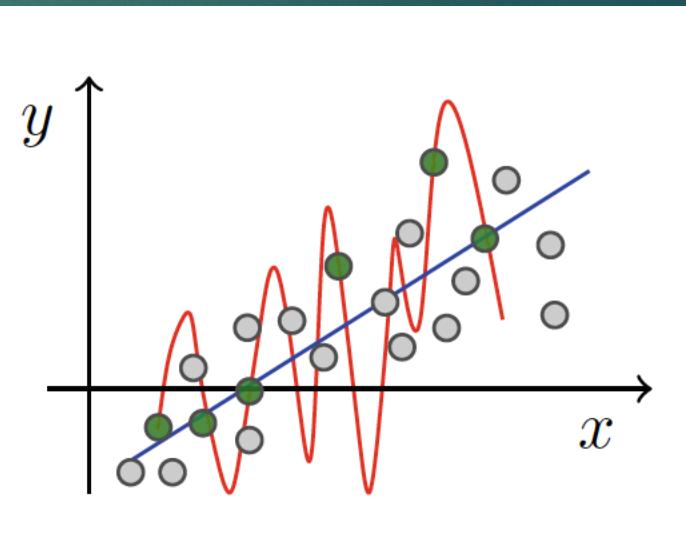
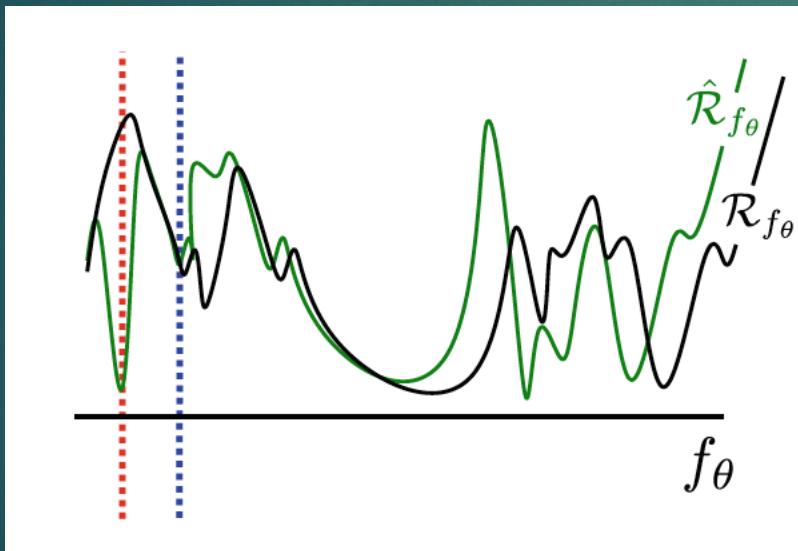
DEGREE 10



DEGREE 2

Overfitting

- ▶ What if empirical risk minimization does not lead to actual risk minimization?
- ▶ What if the trained model does really well on the data it was trained with, but performs poorly on new data?



Regularization

- ▶ Strategies to avoid overfitting
 - ▶ Choose less flexible families of models
 - ▶ Stop training prematurely
 - ▶ Add a regularization term: $C(\theta) = \hat{R}_{f_\theta}(\mathcal{D}) + g(\theta)$ and minimize $C(\theta)$

Regularizers impose constraints on the model:

- ▶ $g_{\ell 1}(\theta) = \sum_{i=1}^T |\theta_i|$ - the $\ell 1$ regularizer leads to sparser parameters' vectors
- ▶ $g_{\ell 2}(\theta) = \sum_{i=1}^T \theta_i^2$ - the $\ell 2$ regularizer leads to shorter parameters' vectors,
i.e. parameters with smaller absolute values

► Cost Function

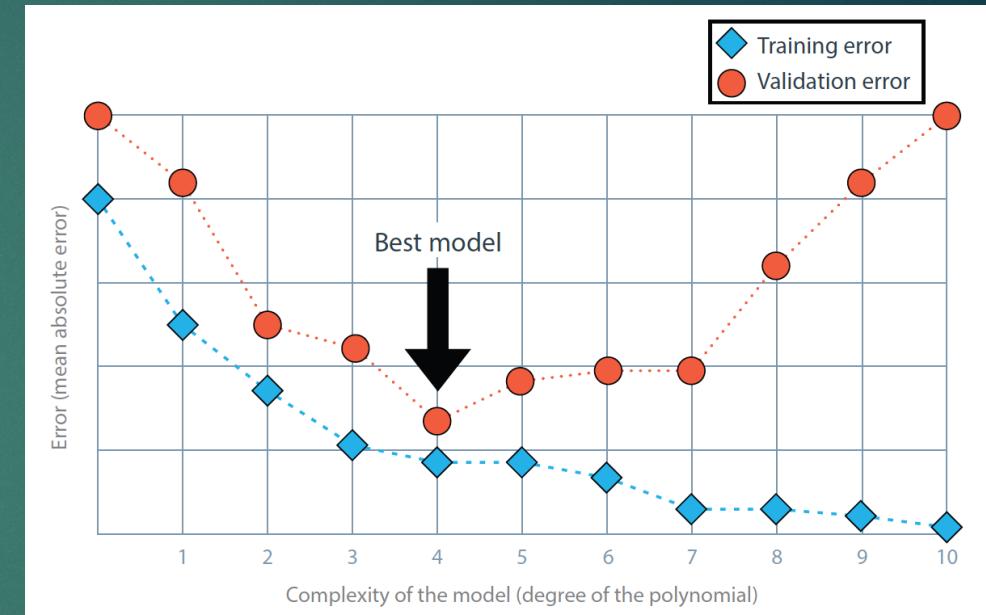
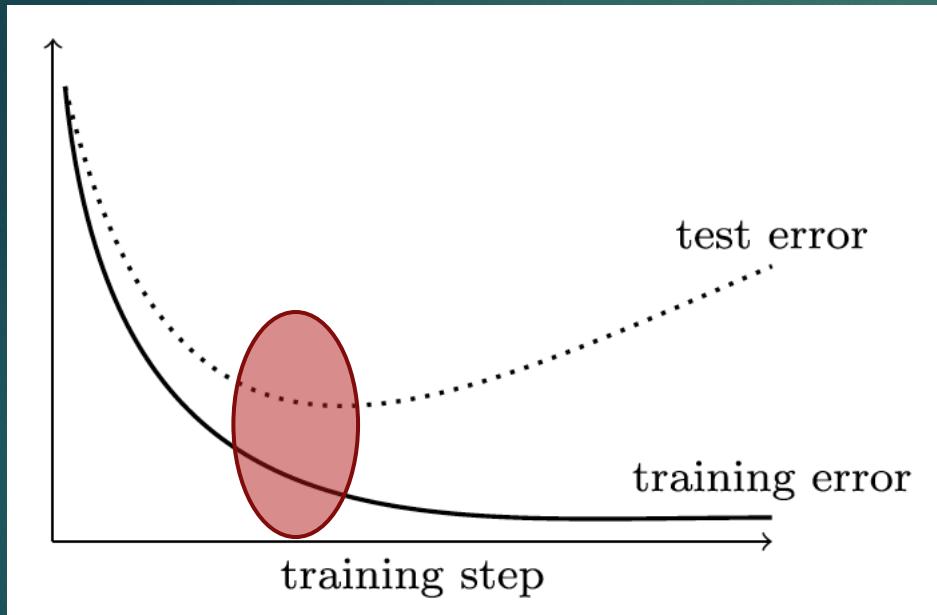
$$C(\theta) = \hat{R}_{f_\theta}(\mathcal{D}) + g(\theta) = \frac{1}{M} \sum_{j=1}^M \left(\theta_0 + \sum_{i=1}^n (\theta_i * x_i^{(j)}) - y^{(j)} \right)^2 + \sum_{i=0}^{T-1} \theta_i^2$$

► Gradient

$$\nabla \hat{R}_{f_\theta}(\mathcal{D}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_{T-1}} \end{pmatrix} = \begin{pmatrix} 2 \sum_{j=1}^M (f_\theta(x^{(j)}) - y^{(j)}) + 2\theta_0 \\ 2 \sum_{j=1}^M (f_\theta(x^{(j)}) - y^{(j)}) * x_1^{(j)} + 2\theta_1 \\ \vdots \\ 2 \sum_{j=1}^M (f_\theta(x^{(j)}) - y^{(j)}) * x_{T-1}^{(j)} + 2\theta_{T-1} \end{pmatrix}$$

Generalization: training vs. test/validation errors

- ▶ the input data \mathcal{D} is divided into a training set and a test set



- ▶ the model is overfitting the training data when the training error keeps decreasing, while the test error starts increasing.
- ▶ The optimal point to stop training is near the test error inflection point OR choose the model to use based on the model complexity graph

Classification

Supervised Learning: Classification

46

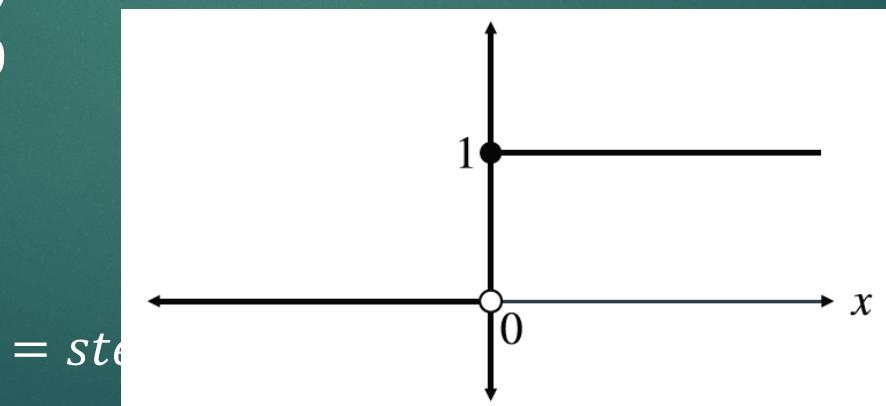
- ▶ Problem:
given a labelled dataset $\mathcal{D} \in \mathcal{X} \otimes \mathcal{Y}$, taken from an unknown function $f^*: \mathcal{X} \rightarrow \mathcal{Y}$, **learn** a function $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$ which approximates f^* and allows predicting $\hat{y} = \hat{f}(x)$ for previously unseen (x, y) , i.e. $(x, y) \notin \mathcal{D}$
- ▶ **Classification** entails that \mathcal{Y} is a discrete set of **labels**
- ▶ **Binary Classification** entails that \mathcal{Y} is a discrete set of **two labels**
Without loss of generality $\mathcal{Y} = \{0, 1\}$

Perceptron Classifier

- ▶ The perceptron classification *line* is expressed as an identity to 0 (small difference compared to the regression in the previous section)
- ▶ If there are n features x_i per data point x and $T = n + 1$ parameters θ_t then the classifier produces a score:

$$\text{score}(x) = \theta_0 + \sum_{i=1}^n \theta_i x_i$$

$$\blacktriangleright \hat{y}(x) = \begin{cases} 0 & \leftarrow \text{score}(x) < 0 \\ 1 & \leftarrow \text{score}(x) \geq 0 \end{cases}$$



$$\blacktriangleright \hat{y}(x) = \text{step}(\text{score}(x))$$

Activation function:

$$\text{step}(x) = \begin{cases} 0 & \leftarrow x < 0 \\ 1 & \leftarrow x \geq 0 \end{cases}$$

The "alien" dataset

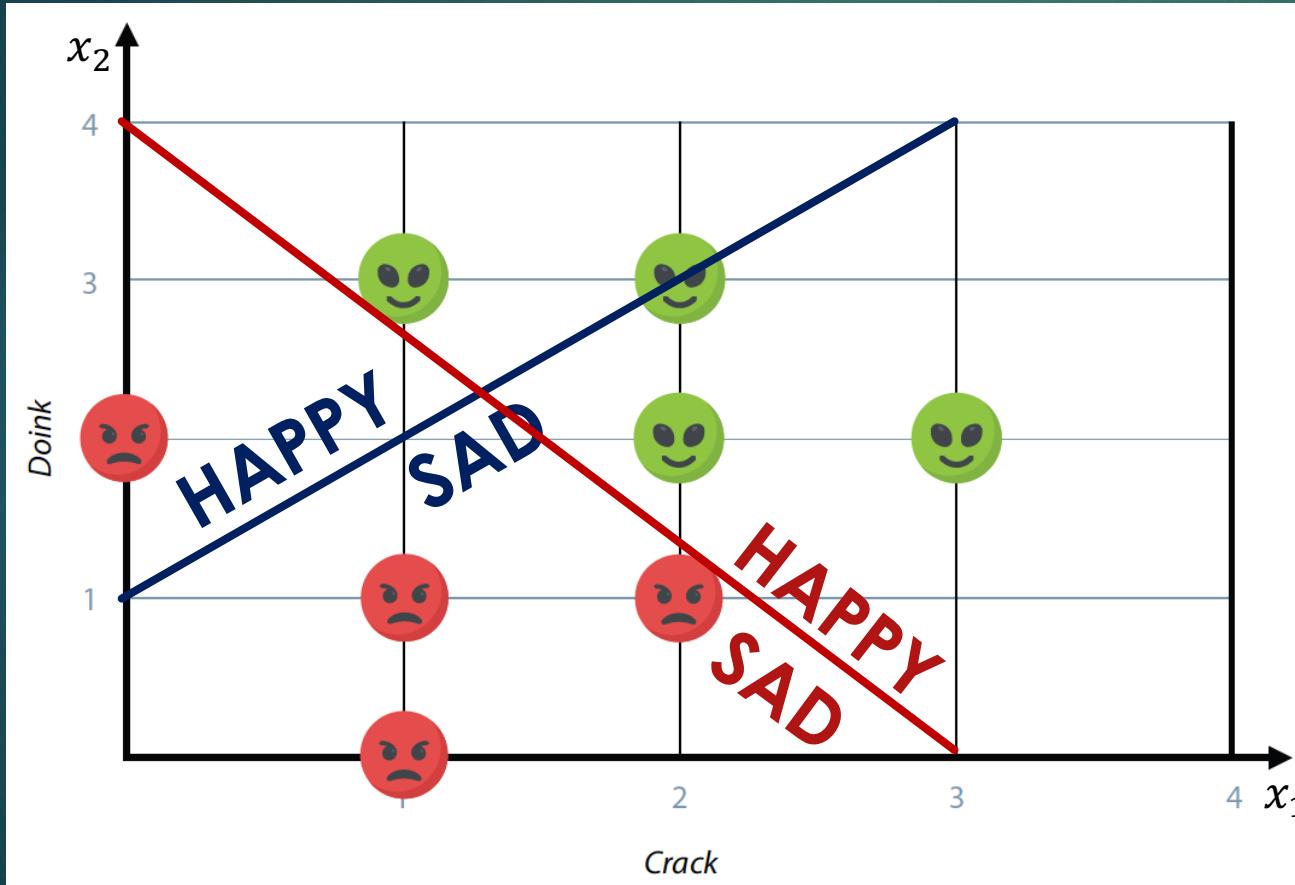
[Serrano'2021-cap 5]

Sentence	Crack	Doink	Mood
Crack!	1	0	Sad
Doink doink!	0	2	Sad
Crack doink!	1	1	Sad
Crack doink crack!	2	1	Sad
Doink crack doink doink!	1	3	Happy
Crack doink doink crack!	2	2	Happy
Doink doink crack crack crack!	3	2	Happy
Crack doink doink crack doink!	2	3	Happy



The "alien" dataset

[Serrano'2021-cap 5]



$$\theta^{(0)} = [1., 0.5, 0.5]$$

$$score_{\theta^{(0)}}(x) = 1. + \frac{x_1}{2} + \frac{x_2}{2}$$

$$\hat{y}_{\theta^{(0)}}(x) = step\left(1. + \frac{x_1}{2} + \frac{x_2}{2}\right)$$

$$\theta^{(t)} = [-4., 4/3, 1.]$$

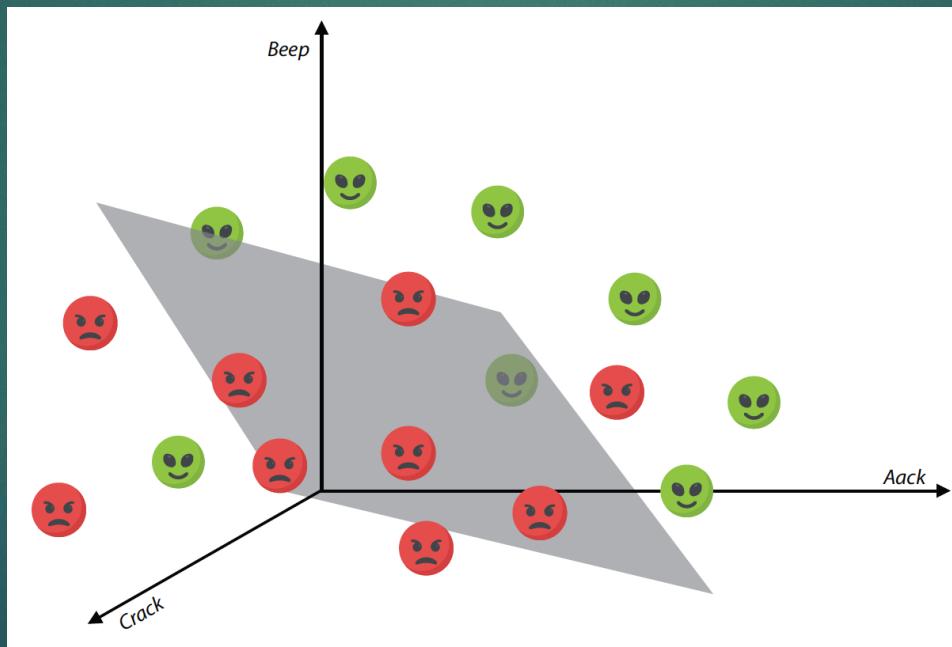
$$score_{\theta^{(t)}}(x) = -4. + \frac{4x_1}{3} + x_2$$

$$\hat{y}_{\theta^{(t)}}(x) = step\left(-4. + \frac{4x_1}{3} + x_2\right)$$

Hyper-Planes

$$score(x) = \theta_0 + \sum_{i=1}^n \theta_i x_i$$

- ▶ The classifier is represented by an hyperplane with $(n - 1)$ -dimensionality



Perceptron Error Function

- ▶ Option 1: count the number of errors

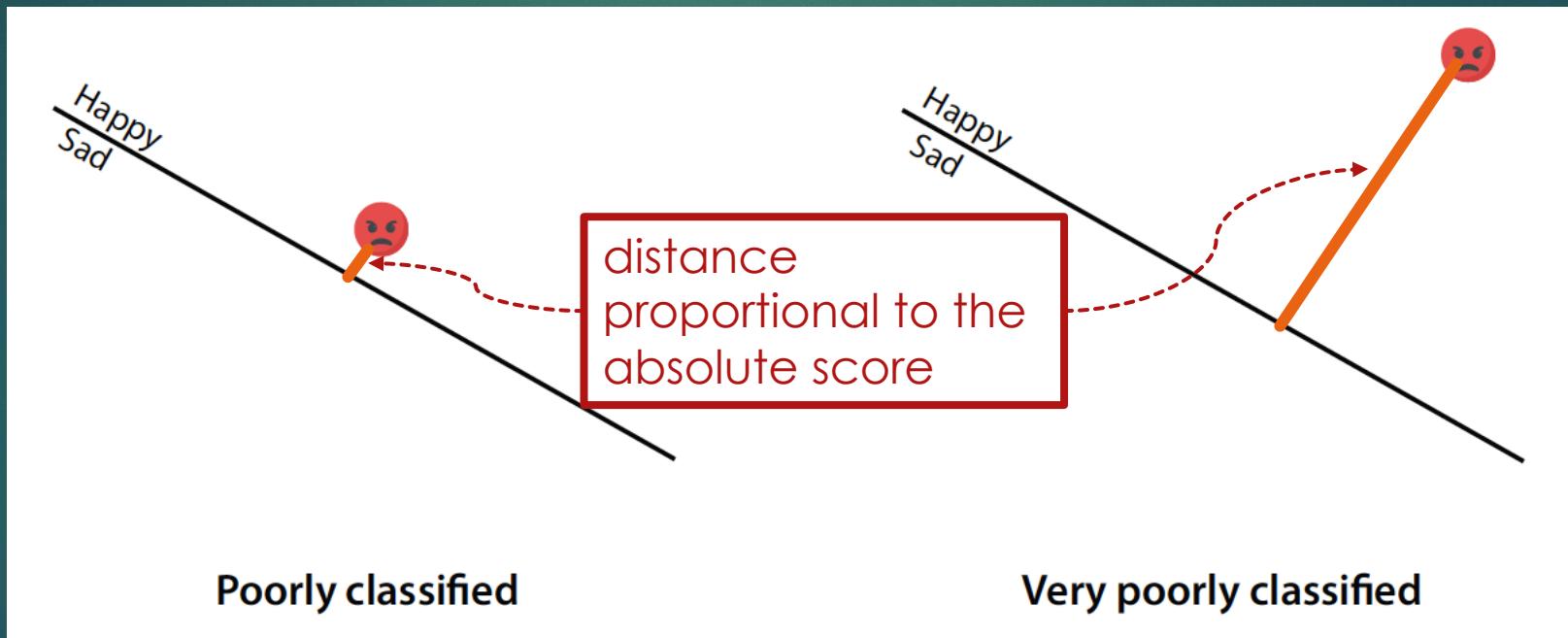


Does not capture **how misclassified** each point was!

Perceptron Error Function

- ▶ Option 1: count the number of errors

Does not capture **how misclassified** each point was!



Perceptron Error Function

- ▶ Option 2: absolute score of misclassified points

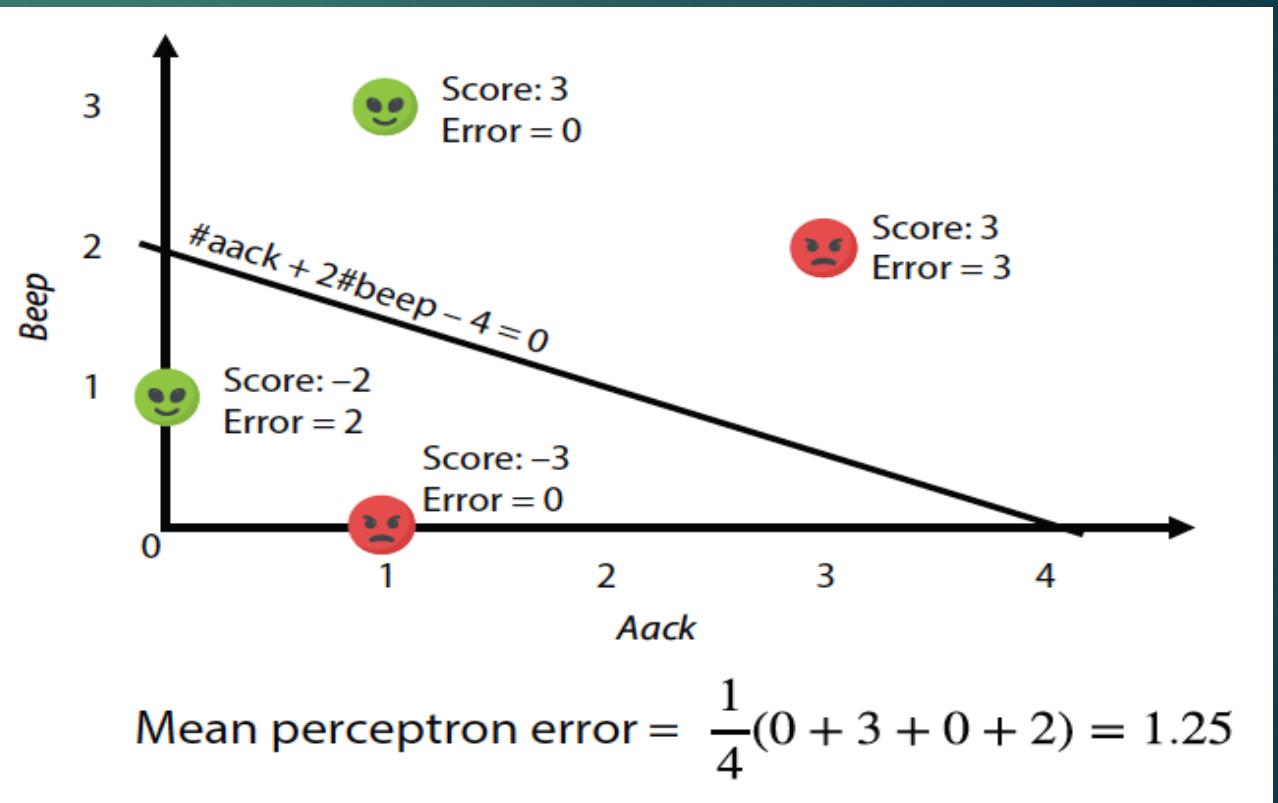
$$L(\hat{y}(x), y) = \begin{cases} 0 & \leftarrow \hat{y}(x) = y \\ \left| \theta_0 + \sum_{i=1}^n \theta_i x_i \right| & \leftarrow \hat{y}(x) \neq y \end{cases}$$



Mean Perceptron Error Function (Risk)

- ▶ average over the dataset of the absolute scores of misclassified points

$$\hat{\mathcal{R}}_{\theta}(\mathcal{D}) = \frac{1}{M} \sum_{x^{(j)} \in \mathcal{D}}$$



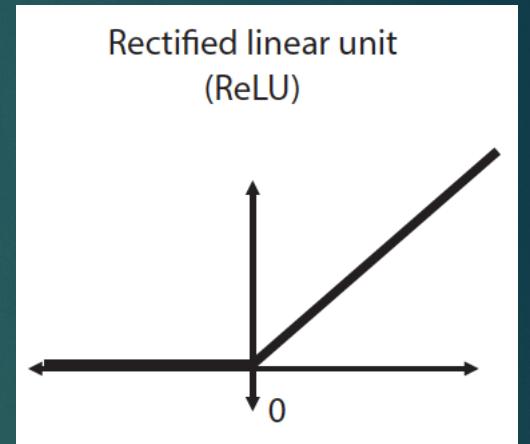
Gradient Descent

$$\hat{\mathcal{R}}_{\theta}(\mathcal{D}) = \frac{1}{M} \sum_{j=1}^M L(\hat{y}(x^{(j)}), y^{(j)}) \quad L(\hat{y}(x), y) = \begin{cases} 0 & \leftarrow \hat{y}(x) = y \\ |\theta_0 + \sum_{i=1}^n \theta_i x_i| & \leftarrow \hat{y}(x) \neq y \end{cases}$$

$$\frac{\partial \text{ReLU}(f(x))}{\partial x} = \text{step}(f(x)) \frac{\partial f(x)}{\partial x} \quad \text{ReLU}(x) = \begin{cases} 0 & \leftarrow x < 0 \\ x & \leftarrow x \geq 0 \end{cases}$$

$$L(\hat{y}(x), y) = y * \text{ReLU}(-\theta_0 - \sum_{i=1}^n \theta_i x_i) + (1 - y) * \text{ReLU}(\theta_0 + \sum_{i=1}^n \theta_i x_i)$$

$$\frac{\partial L(\hat{y}(x), y)}{\partial \theta_{>0}} = y * \text{step}(-\text{score}(x)) * (-x_i) + (1 - y) * \text{step}(\text{score}(x)) * x_i$$



Gradient Descent

$$L(\hat{y}(x), y) = y * \text{ReLU}(-\theta_0 - \sum_{i=1}^n \theta_i x_i) + (1 - y) * \text{ReLU}(\theta_0 + \sum_{i=1}^n \theta_i x_i)$$

$$\begin{aligned}\frac{\partial L(\hat{y}(x), y)}{\partial \theta_{>0}} &= y * \text{step}(-\text{score}(x)) * (-x_i) + (1 - y) * \text{step}(\text{score}(x)) * x_i = \\ &= \left(-y * \text{step}(-\text{score}(x)) + (1 - y) * \text{step}(\text{score}(x)) \right) x_i\end{aligned}$$

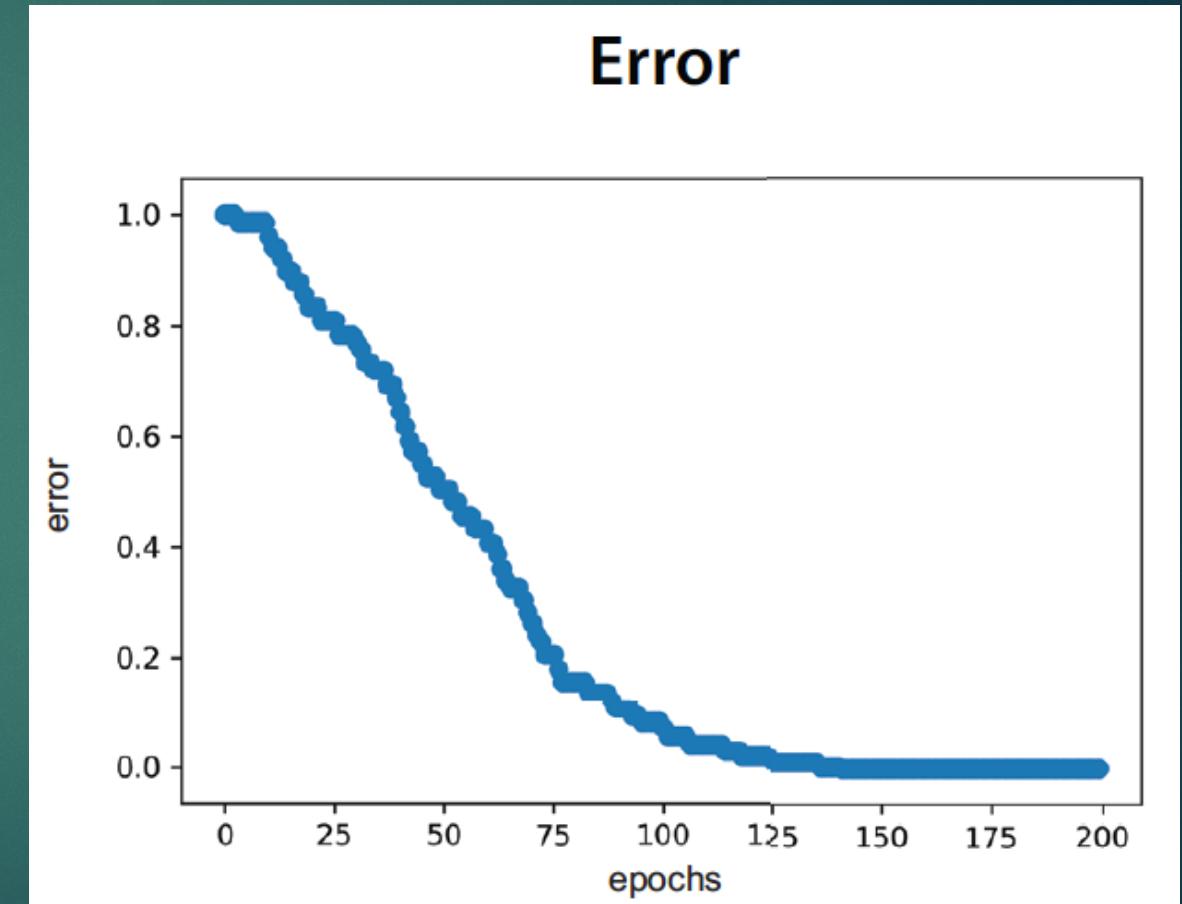
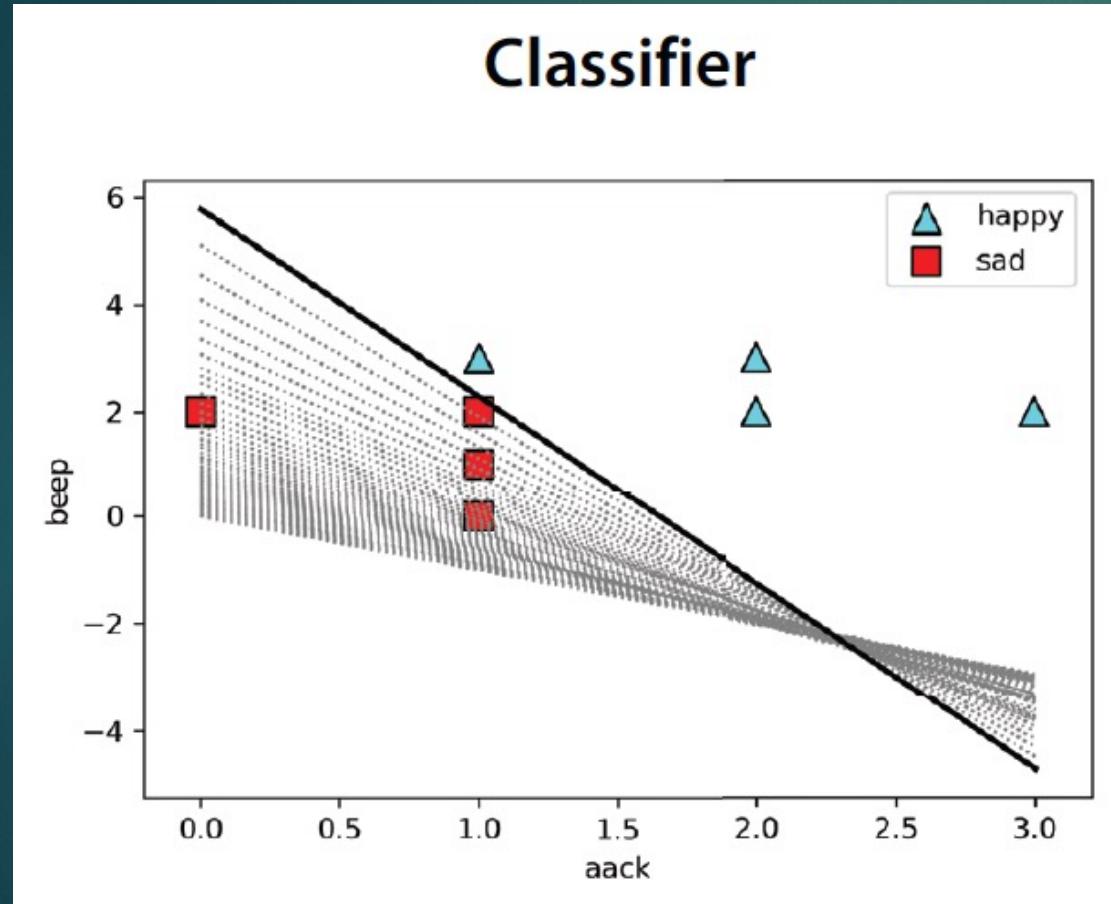
$$\frac{\partial L(\hat{y}(x), y)}{\partial \theta_0} = -y * \text{step}(-\text{score}(x)) + (1 - y) * \text{step}(\text{score}(x))$$

$$\hat{\mathcal{R}}_\theta(\mathcal{D}) = \frac{1}{M} \sum_{j=1}^M L(\hat{y}(x^{(j)}), y^{(j)})$$

$$\nabla \hat{\mathcal{R}}_{f_\theta}(\mathcal{D}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_{T-1}} \end{pmatrix} = \begin{pmatrix} \frac{1}{M} \sum_{j=1}^M [-y^{(j)} * \text{step}(-\text{score}(x^{(j)})) + (1 - y^{(j)}) * \text{step}(\text{score}(x^{(j)}))] \\ \frac{1}{M} \sum_{j=1}^M [-y^{(j)} * \text{step}(-\text{score}(x^{(j)})) + (1 - y^{(j)}) * \text{step}(\text{score}(x^{(j)}))] * x_1^{(j)} \\ \vdots \\ \frac{1}{M} \sum_{j=1}^M [-y^{(j)} * \text{step}(-\text{score}(x^{(j)})) + (1 - y^{(j)}) * \text{step}(\text{score}(x^{(j)}))] * x_{T-1}^{(j)} \end{pmatrix}$$

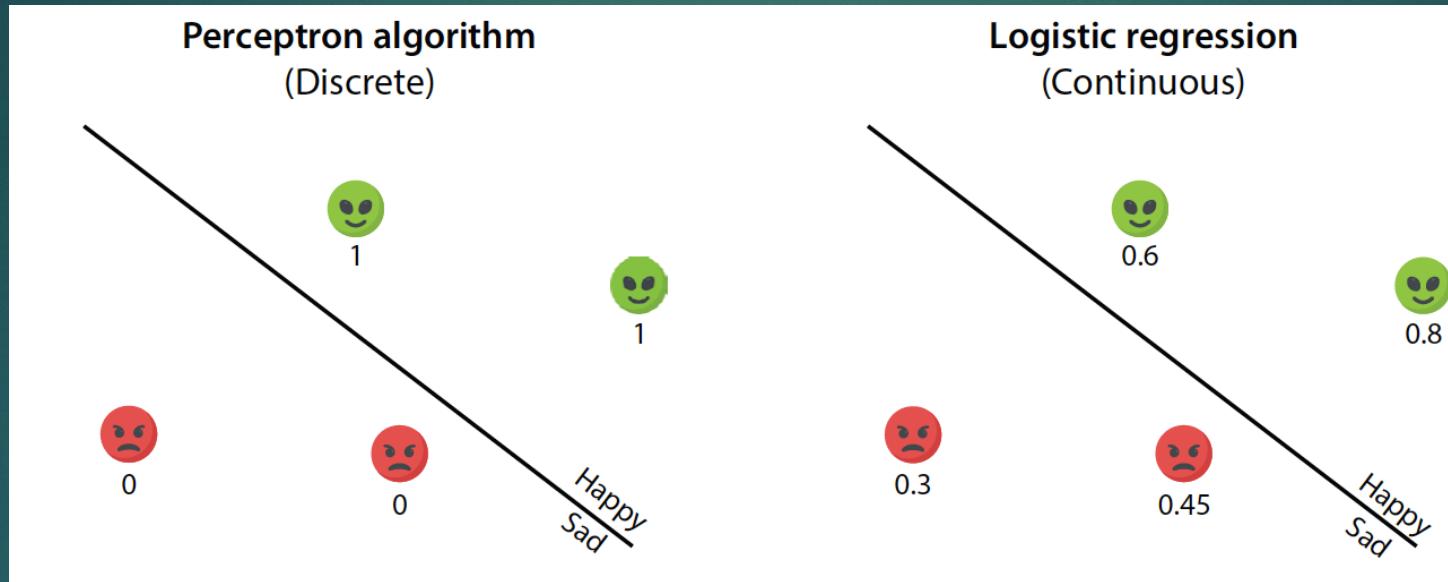
$\hat{y}(x)$	y	$\text{score}(x)$	$\frac{\partial L(\hat{y}(x), y)}{\partial \theta_{>0}}$
0	0	< 0	0
0	1	< 0	$-x_i$
1	0	> 0	x_i
1	1	> 0	0

Perceptron Classifier



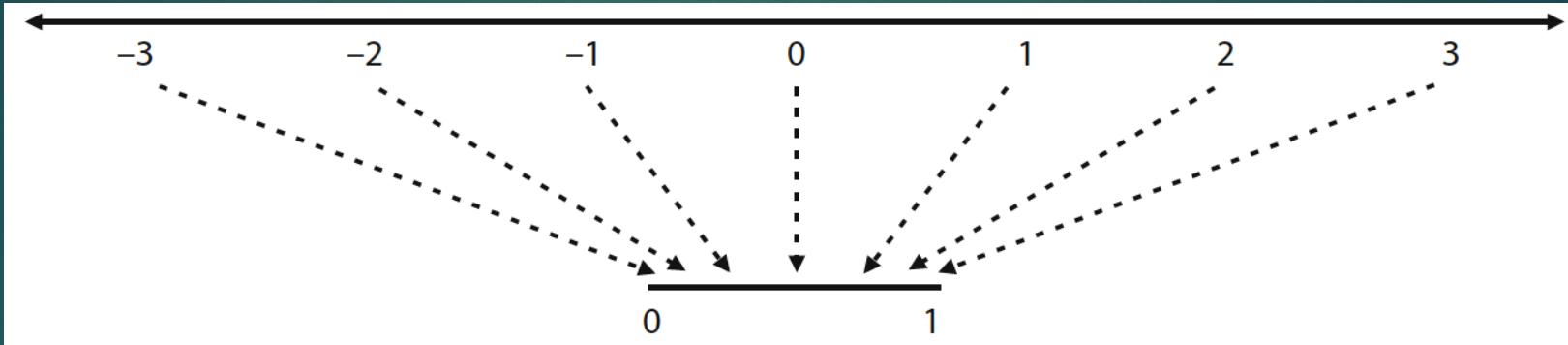
Logistic Classifiers

- ▶ Binary classifiers which return a continuous value in [0..1]



- ▶ The number assigned to each data point is the probability that its label is $y=1$: $p_1(x)$
- ▶ The probability of being $y=0$ is $p_0(x)=1-p_1(x)$

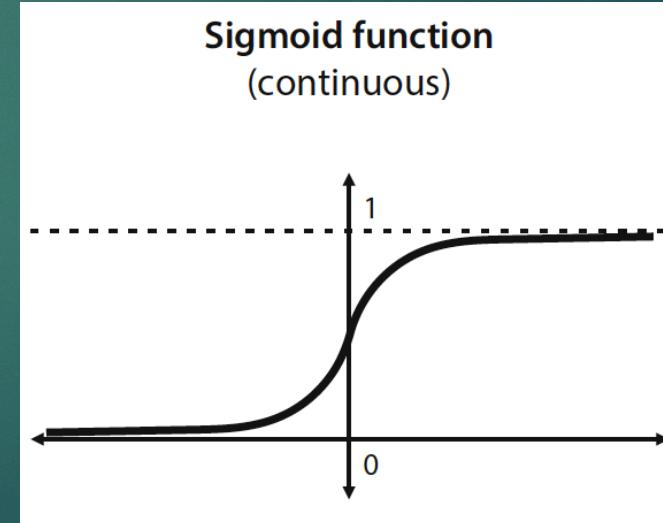
Logistic Classifiers



Sigmoid function: the score $\in \mathbb{R}$ is compressed into $[0..1]$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\hat{y}(x) = \sigma(score(x)) \\ = \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i x_i\right)$$



Logistic Classifier Error: log loss

- If $y = 1$ then $L(\hat{y}(x), y) = -\ln(\hat{y}(x))$
- If $y = 0$ then $L(\hat{y}(x), y) = -\ln(1 - \hat{y}(x))$

$$L(\hat{y}(x), y) = -y * \ln(\hat{y}(x)) - (1 - y) * \ln(1 - \hat{y}(x))$$

$$\hat{\mathcal{R}}_{\theta}(\mathcal{D}) = \sum_{j=1}^M L(\hat{y}(x^{(j)}), y^{(j)})$$

Point	Label	Predicted label	Probability of being its label	Log loss
1	0 (Sad)	0.953	0.047	$-\ln(0.047) = 3.049$
2	1 (Happy)	0.731	0.731	$-\ln(0.731) = 0.313$
3	1 (Happy)	0.119	0.119	$-\ln(0.119) = 2.127$
4	0 (Sad)	0.119	0.881	$-\ln(0.881) = 0.127$

Gradient Descent

$$\hat{y}(x) = \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i x_i\right)$$

$$L(\hat{y}(x), y) = -y * \ln(\hat{y}(x)) - (1 - y) * \ln(1 - \hat{y}(x))$$

$$\frac{\partial \sigma(f(\theta))}{\partial \theta} = \sigma(f(\theta)) \left(1 - \sigma(f(\theta))\right) \frac{\partial f(\theta)}{\partial \theta}$$

$$\frac{\partial \ln(f(z))}{\partial z} = \frac{1}{f(z)} \frac{\partial f(z)}{\partial z}$$

$$\frac{\partial \hat{y}(x)}{\partial \theta_{>0}} = \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i x_i\right) \left(1 - \sigma\left(\theta_0 + \sum_{i=1}^n \theta_i x_i\right)\right) x_i = \hat{y}(x)(1 - \hat{y}(x))x_i$$

$$\frac{\partial L(\hat{y}(x), y)}{\partial \theta_{>0}} = -y * \frac{1}{\hat{y}(x)} \frac{\partial \hat{y}(x)}{\partial \theta_{>0}} - (1 - y) * \frac{-1}{1 - \hat{y}(x)} \frac{\partial \hat{y}(x)}{\partial \theta_{>0}} = \dots = (\hat{y} - y) x_i$$

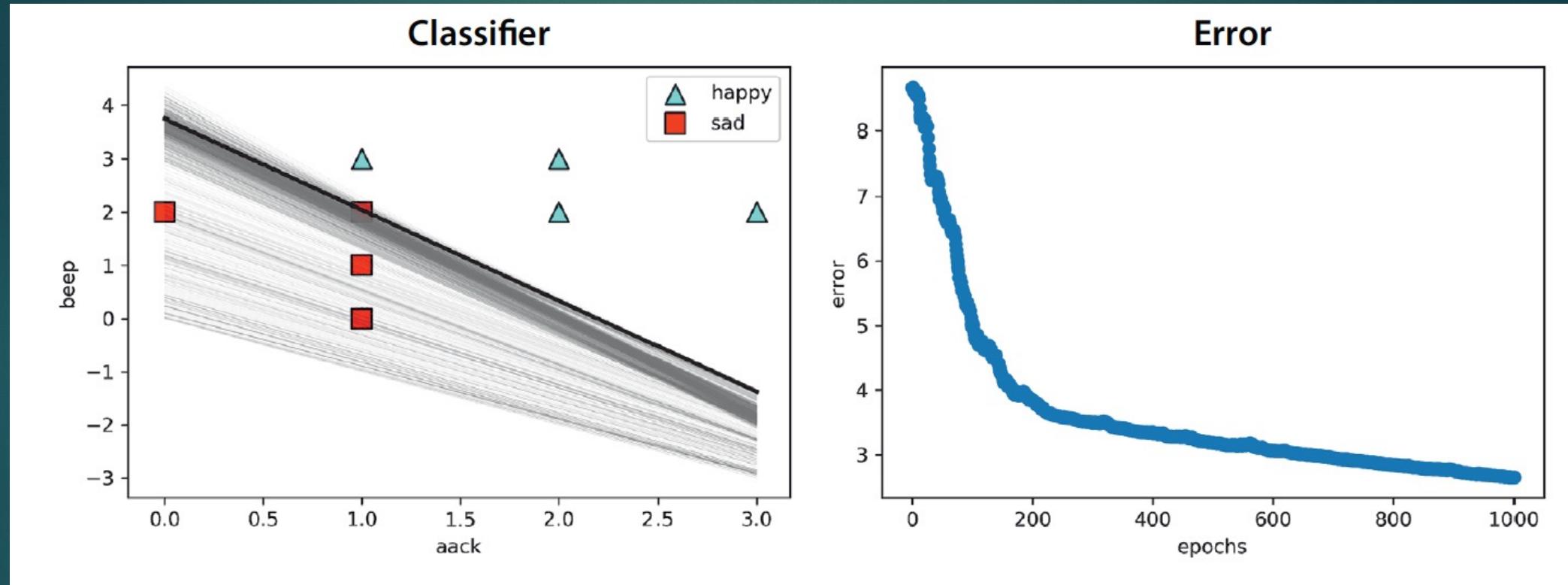
$$\frac{\partial L(\hat{y}(x), y)}{\partial \theta_0} = (\hat{y} - y)$$

Gradient Descent

$$\hat{\mathcal{R}}_{\theta}(\mathcal{D}) = \sum_{j=1}^M L(\hat{y}(x^{(j)}), y^{(j)})$$

$$\nabla \hat{\mathcal{R}}_{f_{\theta}}(\mathcal{D}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_{T-1}} \end{pmatrix} = \begin{pmatrix} \sum_{j=1}^M [\hat{y}(x^{(j)}) - y^{(j)}] \\ \sum_{j=1}^M [\hat{y}(x^{(j)}) - y^{(j)}] * x_1^{(j)} \\ \vdots \\ \sum_{j=1}^M [\hat{y}(x^{(j)}) - y^{(j)}] * x_{T-1}^{(j)} \end{pmatrix}$$

Logistic Classifier

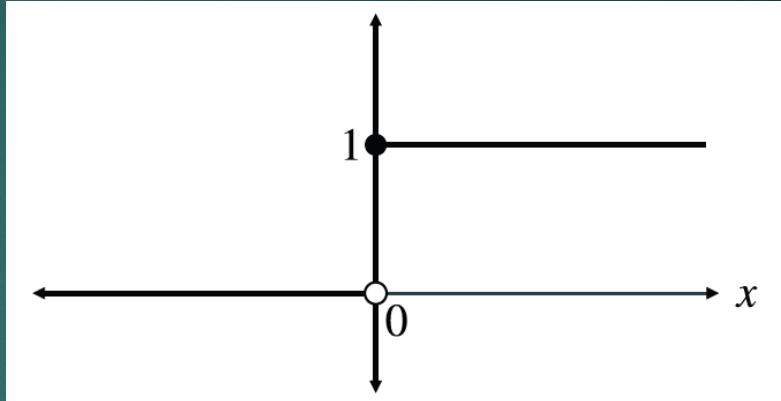


Perceptrons

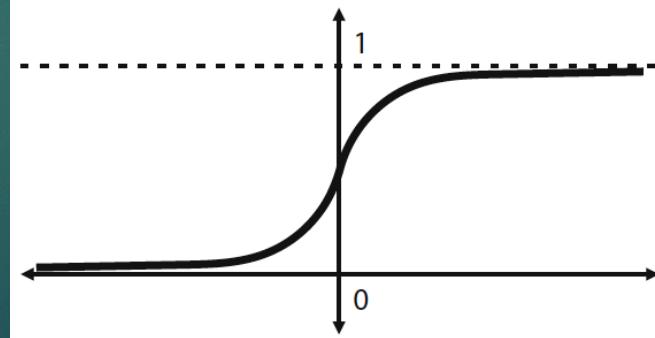
Perceptron Classifier:
 $f_{activation}(\cdot) = \text{step}(\cdot)$

$$\hat{y}(x) = f_{activation}\left(\theta_0 + \sum_{i=1}^n \theta_i x_i\right)$$

Logistic Classifier:
 $f_{activation}(\cdot) = \sigma(\cdot)$



Sigmoid function
 (continuous)



Multi Label Classifier

- ▶ Train one classifier for each class $l_i \in \mathcal{Y} = \{l_0, \dots, l_{L-1}\}, i = 0 \dots L - 1$
- ▶ To classify a data point x obtain the L different scores
- ▶ Compute a probability distribution over \mathcal{Y} by normalizing the scores

Multi Label Classifier: Example

- ▶ $y = \{"dog", "cat", "bird"\}, L = 3$
- ▶ $scores(x) = \{3, 2, -1\}$
- ▶ Get a probability distribution:
 - ▶ normalize $p(\hat{y}(x)) = \{3/5, 2/5, -1/5\}$
- ▶ use softmax $p(\hat{y}_i(x)) = \frac{e^{\hat{y}_i(x)}}{\sum_{l=0}^{L-1} e^{\hat{y}_l(x)}}$
$$p(\hat{y}(x)) = \{0.721, 0.265, 0.013\}$$