

# ESSENTIAL MACHINE LEARNING ALGORITHMS



GRAPHCORE

# MACHINE LEARNING

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning



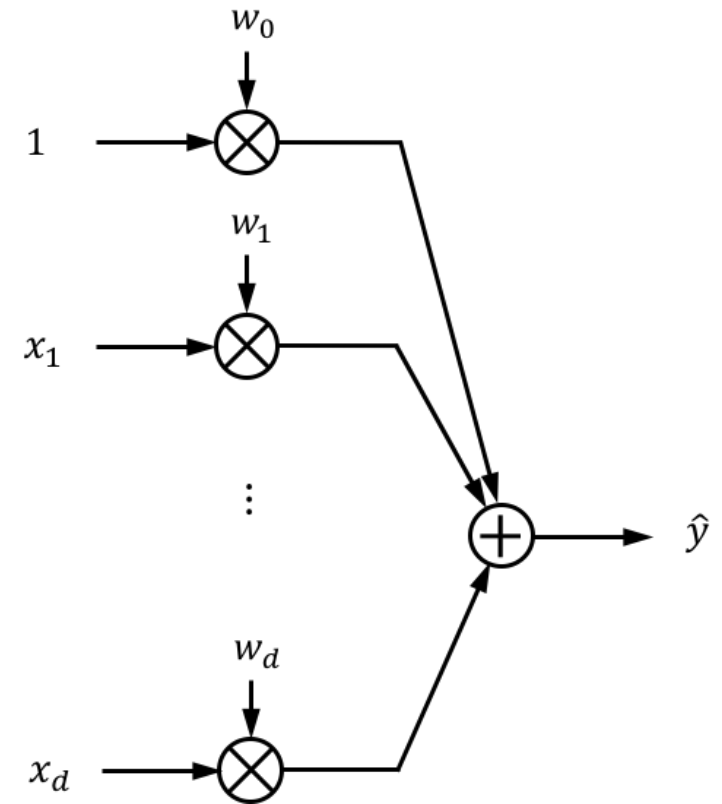
# LINEAR REGRESSION

Data set  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$

Model  $y = w_0 + w_1 x_1 + \dots + w_d x_d + \epsilon$   
 $= \mathbf{w}^T \mathbf{x} + \epsilon$

Prediction

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$



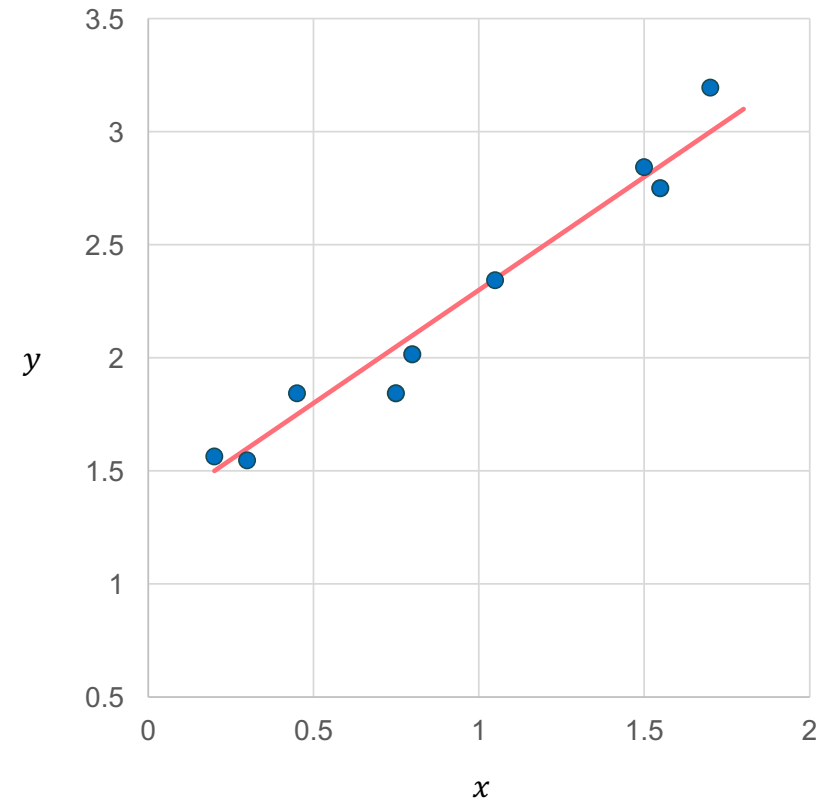
# LINEAR REGRESSION

Data set  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$

Model  $y = w_0 + w_1 x_1 + \dots + w_d x_d + \epsilon$   
 $= \mathbf{w}^T \mathbf{x} + \epsilon$

Prediction

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$



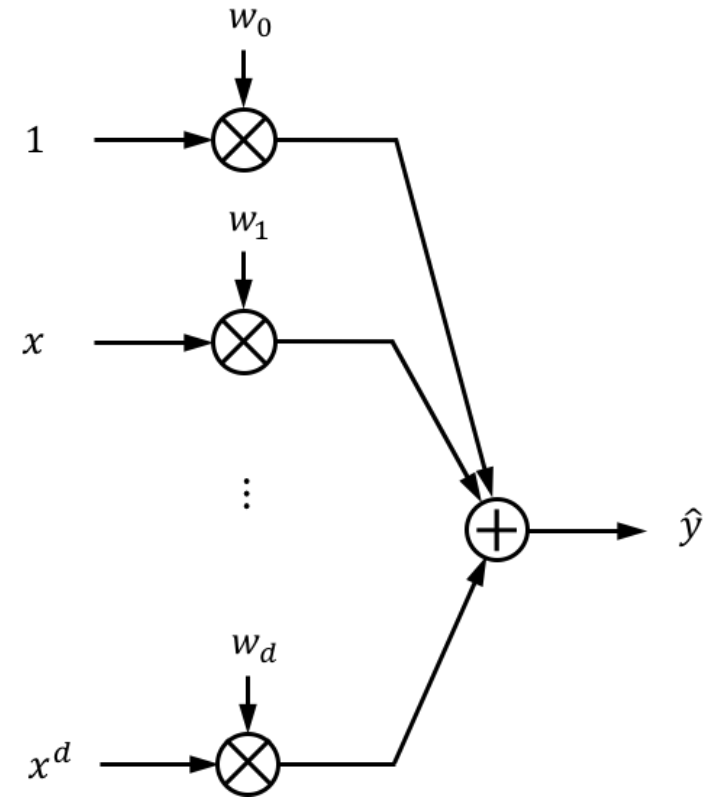
# LINEAR REGRESSION

Data set  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$

Prediction  $\hat{y} = \mathbf{w}^T \phi(\mathbf{x})$

e.g.:

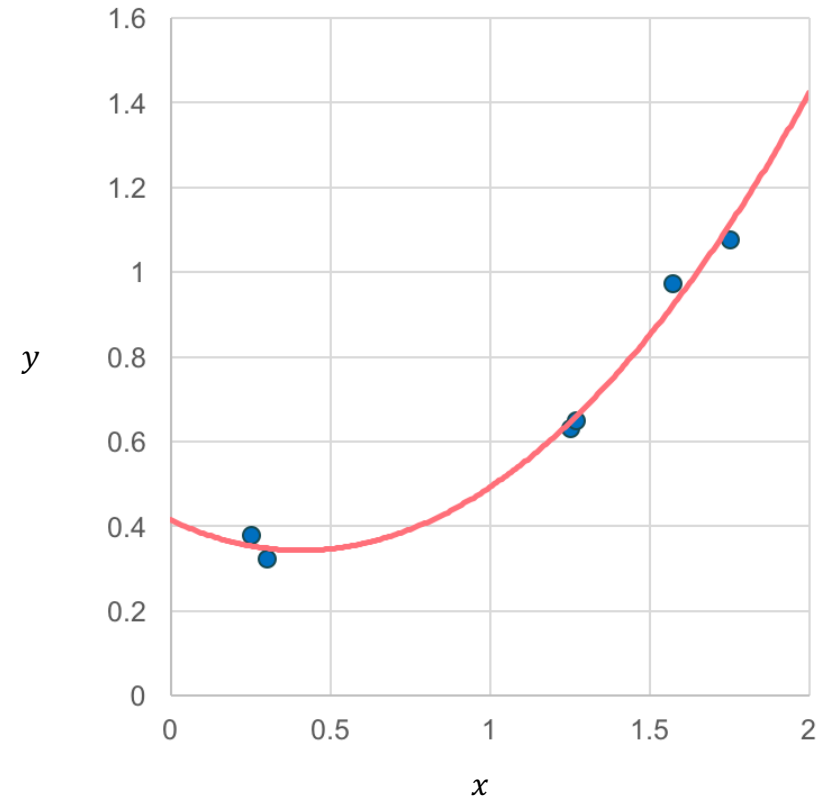
$$\hat{y} = w_0 + w_1 x + w_2 x^2 \dots + w_d x^d$$



# OVERFITTING

- *Model capacity*
- *Size of training data set,  $N$*

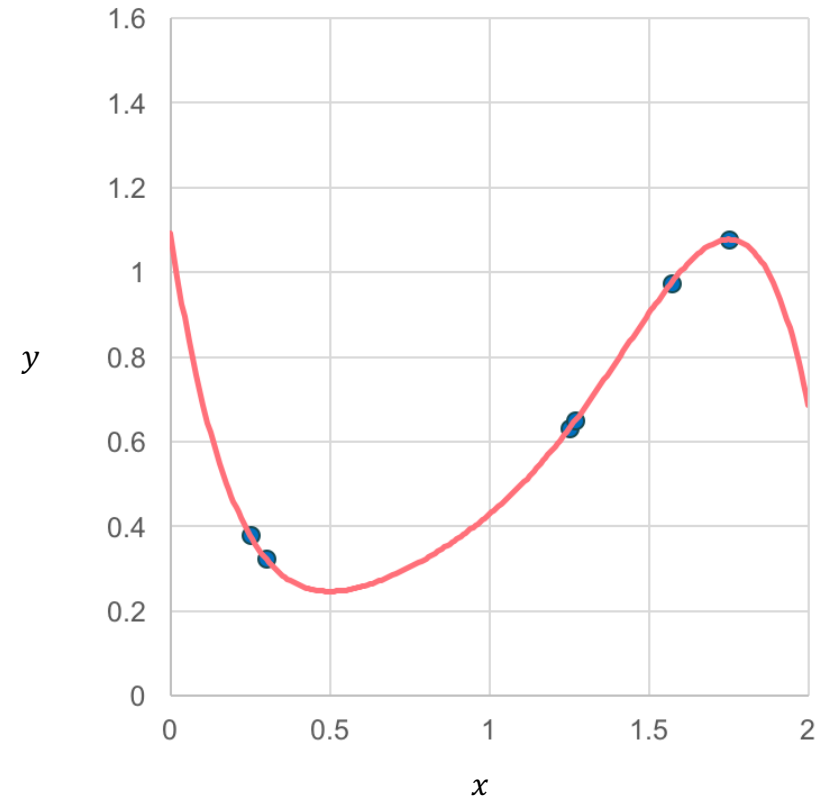
Model capacity and overfitting



# OVERFITTING

- *Model capacity*
- *Size of training data set,  $N$*

Model capacity and overfitting



# REGULARIZATION

Loss function:

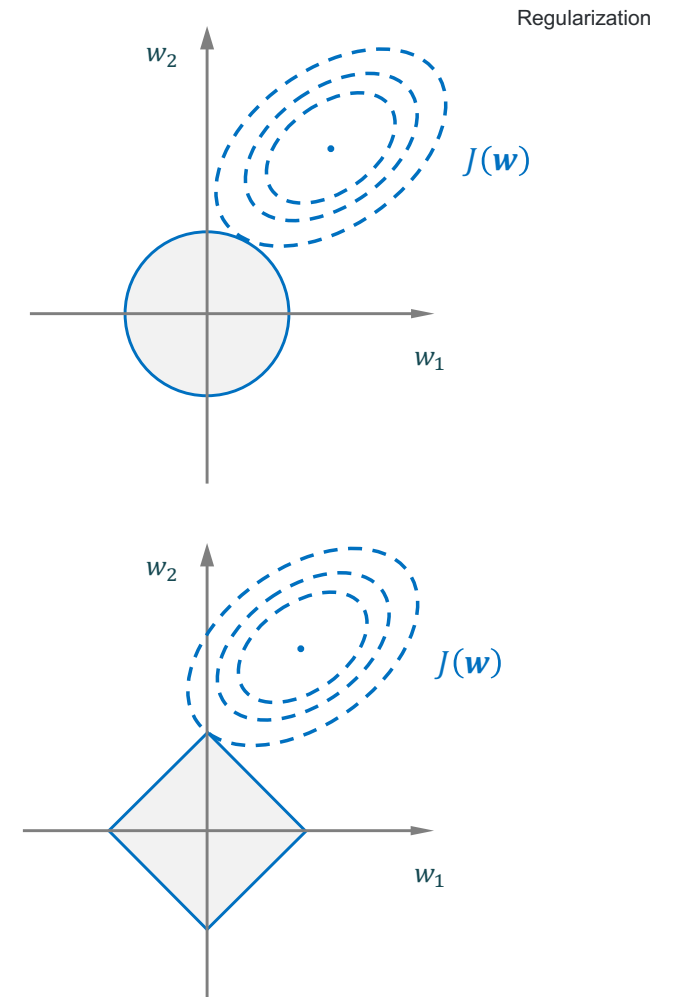
$$J(\mathbf{w}) + \text{norm penalty}$$

$L^2$  regularization (weight decay)

$$\text{norm penalty} = \frac{\alpha}{2} \|\mathbf{w}\|_2^2$$

$L^1$  regularization

$$\text{norm penalty} = \alpha \|\mathbf{w}\|_1$$





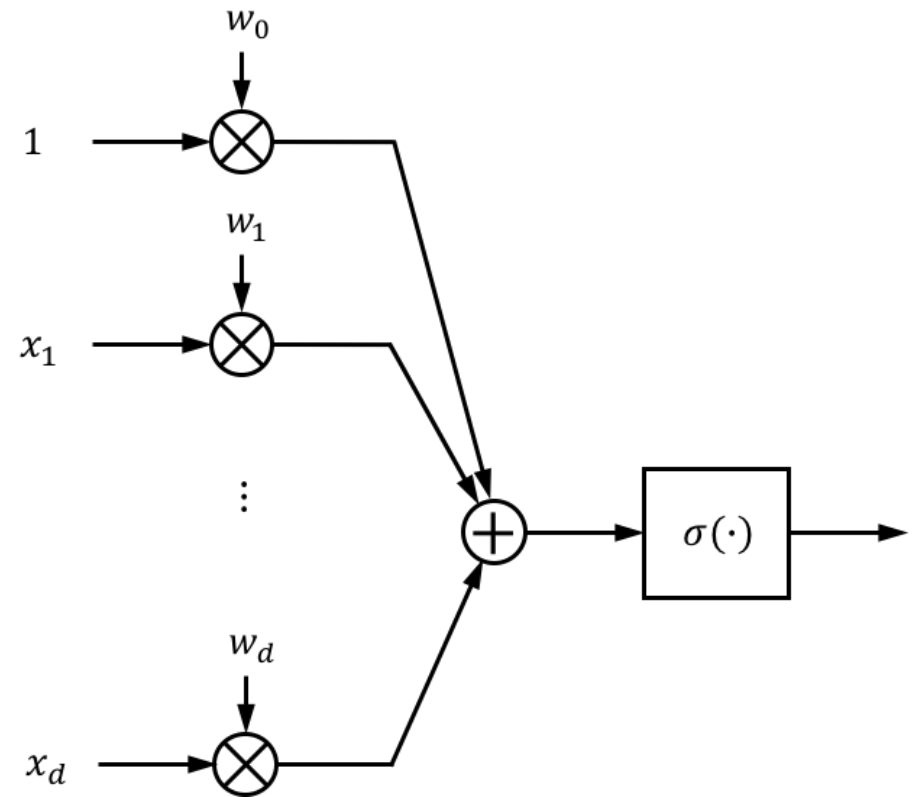
# LOGISTIC REGRESSION

## Binary Classification

$$\mu(\mathbf{x}) = p(y = 1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

where  $\sigma(\cdot)$  is the sigmoid function

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



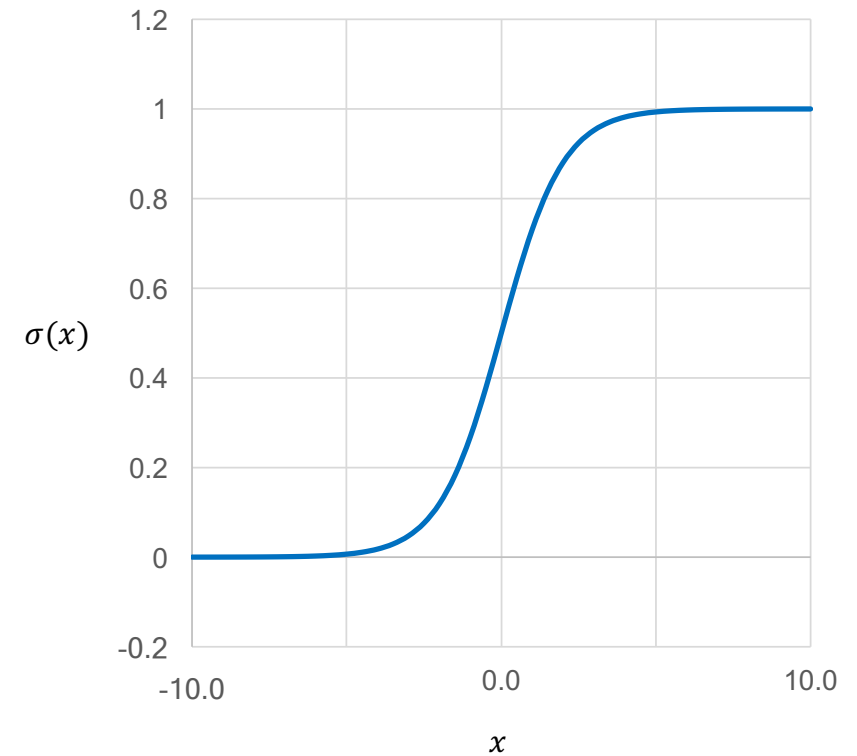
# LOGISTIC REGRESSION

## *Binary Classification*

$$\mu(\mathbf{x}) = p(y = 1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

where  $\sigma(\cdot)$  is the sigmoid function

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



# MAXIMAL MARGIN CLASSIFIER

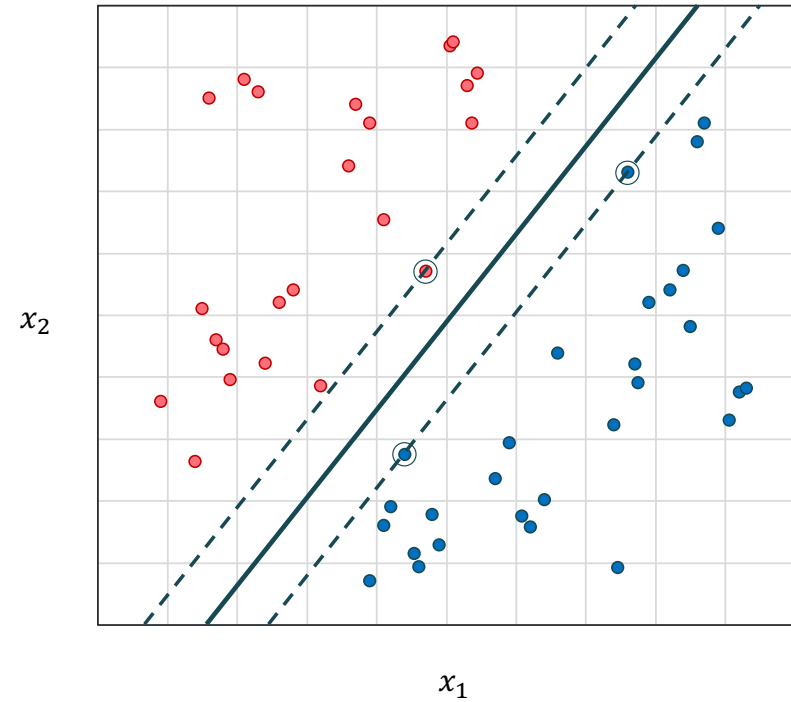
Optimization problem:

$$\max_w M$$

subject to

$$y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} > M, \quad i = 1, \dots, N$$

$$\|\mathbf{w}\|_2^2 = 1$$



# SUPPORT VECTOR CLASSIFIER

Case of not linearly separable classes:

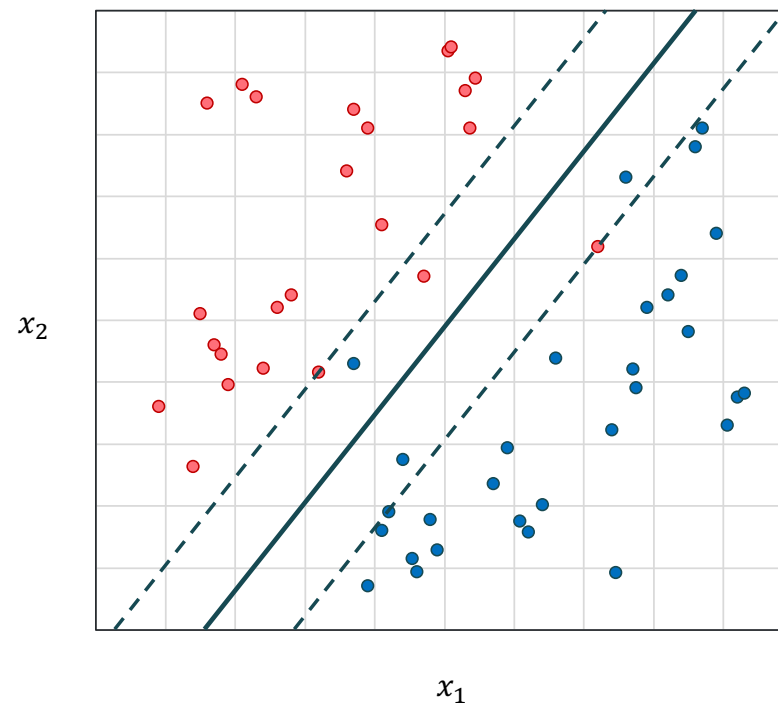
$$\max_w M$$

subject to

$$y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} > M(1 - \epsilon^{(i)}), \quad i = 1, \dots, N$$

$$\|\mathbf{w}\|_2^2 = 1$$

with  $\epsilon^{(i)} \geq 0, \quad \sum_{i=1}^N \epsilon^{(i)} \leq C$



# SUPPORT VECTOR MACHINE (SVM)

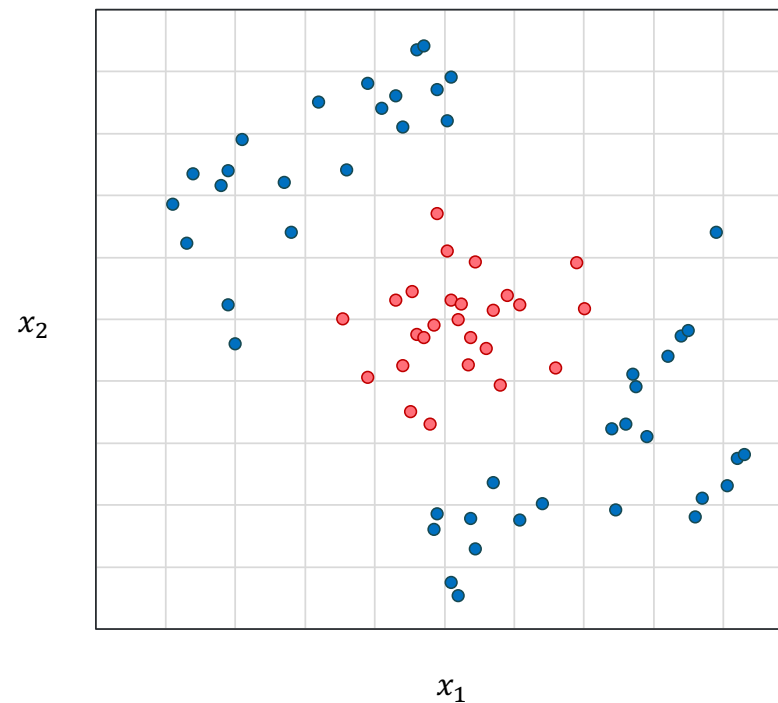
*Kernel trick:*

$$\hat{y} = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^N \alpha_i \mathbf{x}^T \mathbf{x}^{(i)}$$

replacing  $\mathbf{x}^T \mathbf{x}'$  by  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$

$$\hat{y} = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)})$$

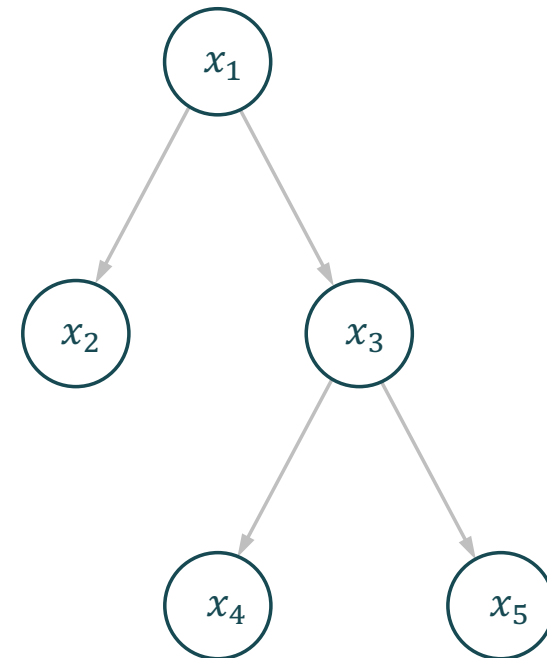
Computation of  $k(\mathbf{x}, \mathbf{x}')$  without working in the enlarged feature space



# DIRECTED GRAPHICAL MODELS (BAYES NETS)

Representation of joint distribution by defining *conditional independence* assumptions

$$P(x_1, x_2, \dots, x_n) = P(x_1) \cdot P(x_2 | x_1) \cdot \\ \cdot P(x_3 | x_2, x_1) \cdots P(x_n | x_{n-1}, \dots, x_2, x_1)$$



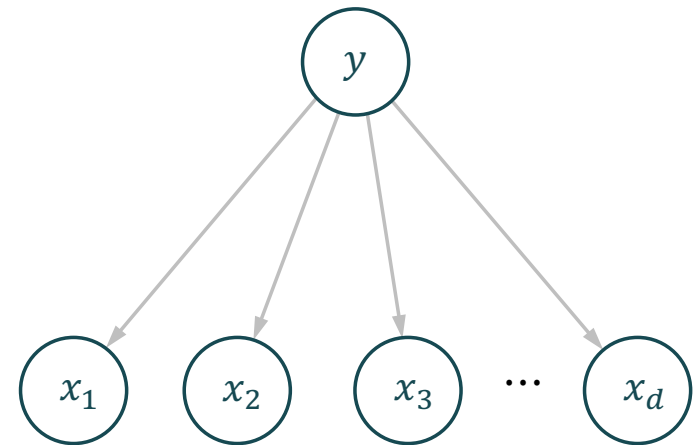
# NAÏVE BAYES CLASSIFIER

$$\hat{y} = \arg \max_y P(y|\mathbf{x}) = \arg \max_y \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

*Conditional independence* assumption:

$$\begin{aligned} P(\mathbf{x}|y) &= P(x_1, x_2, \dots, x_d|y) \\ &= P(x_1|y) \cdot P(x_2|y) \cdots P(x_d|y) \end{aligned}$$

$$\Rightarrow \hat{y} = \arg \max_y P(y) \prod_{i=1}^d P(x_i|y)$$



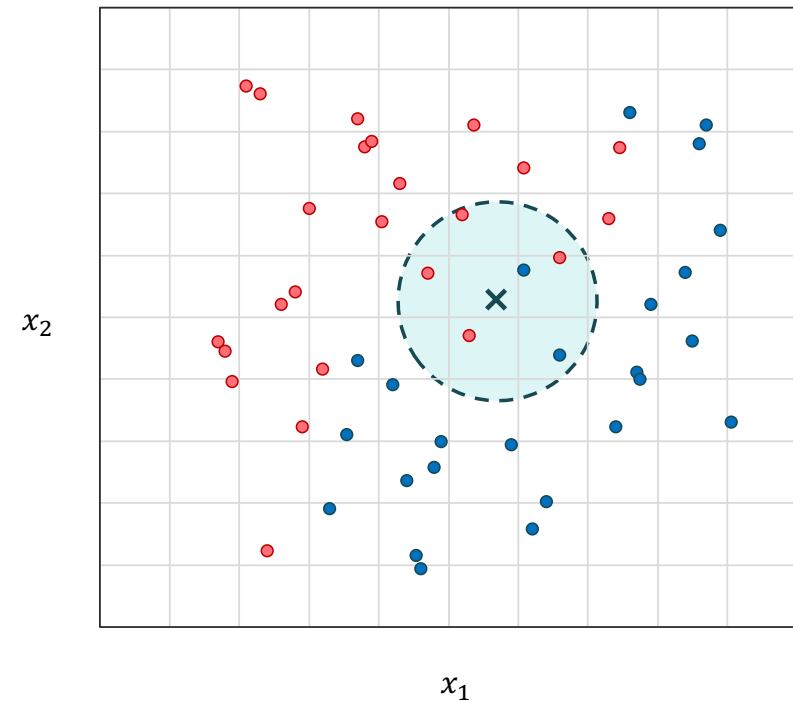
# K-NEAREST NEIGHBOURS (K-NN)

Data set  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$

Procedure for a test observation  $\mathbf{x}$ :

- identify the  $K$  training points closest to  $\mathbf{x}$ , represented by the set  $N_K(\mathbf{x})$
- select the class  $j$  with largest probability

$$P(y = j|\mathbf{x}) = \frac{1}{K} \sum_{i \in N_K(\mathbf{x})} \mathbb{I}(y^{(i)} = j)$$

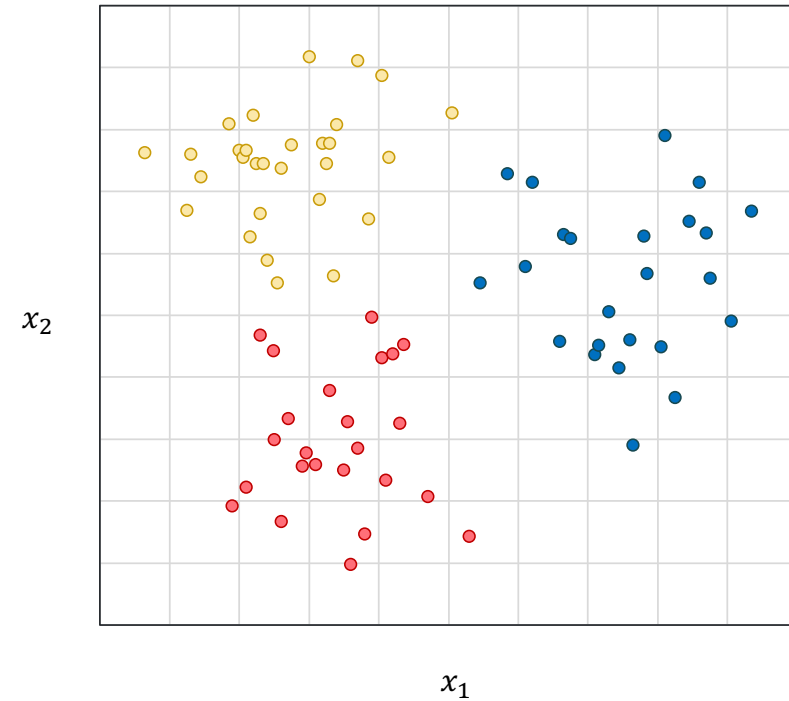




# K-MEANS CLUSTERING

Procedure      after taking  $K$  points as the initial *clusters centres*  $\mu_k$ , repeat the steps:

- assign each data point to the cluster that has the closest centre
- update the position of each cluster centre, by re-computing the mean  $\mu_k$



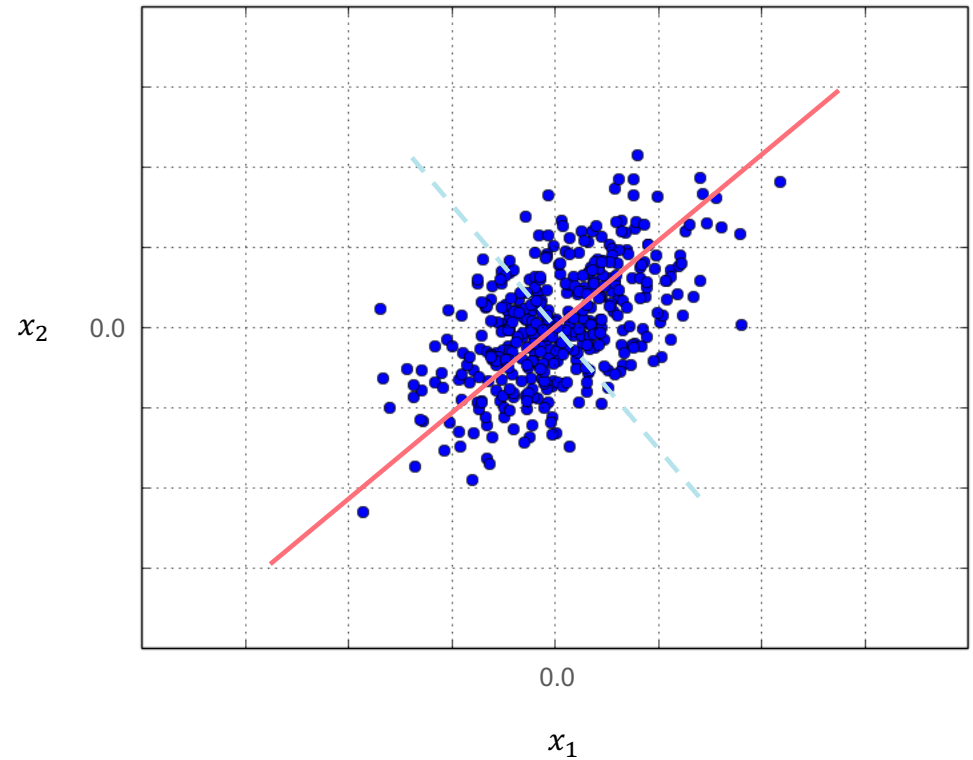
# PRINCIPAL COMPONENT ANALYSIS (PCA)

Learns a lower dimensional representation

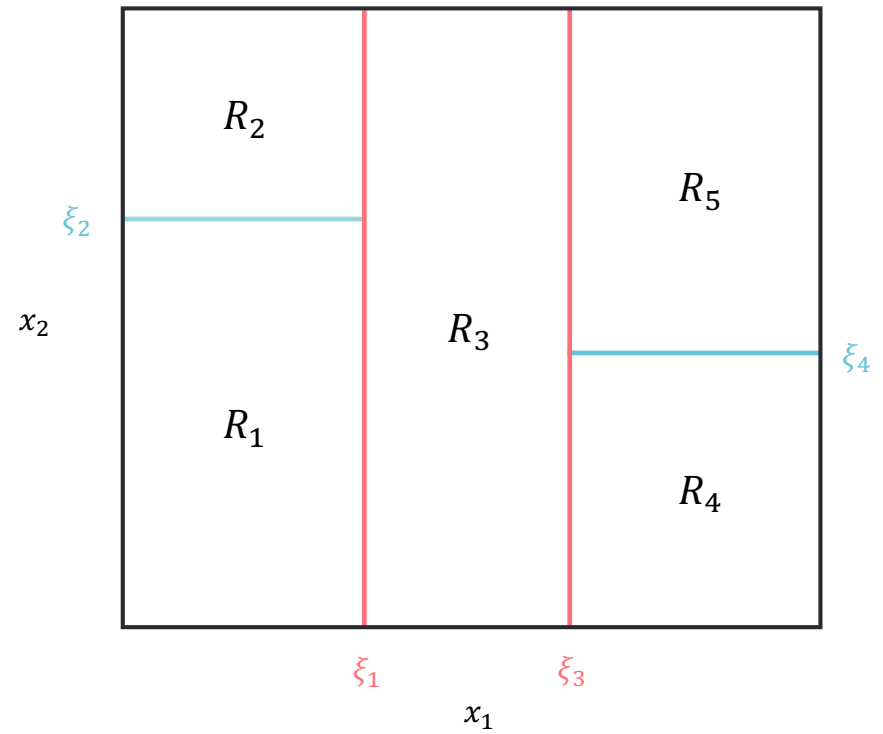
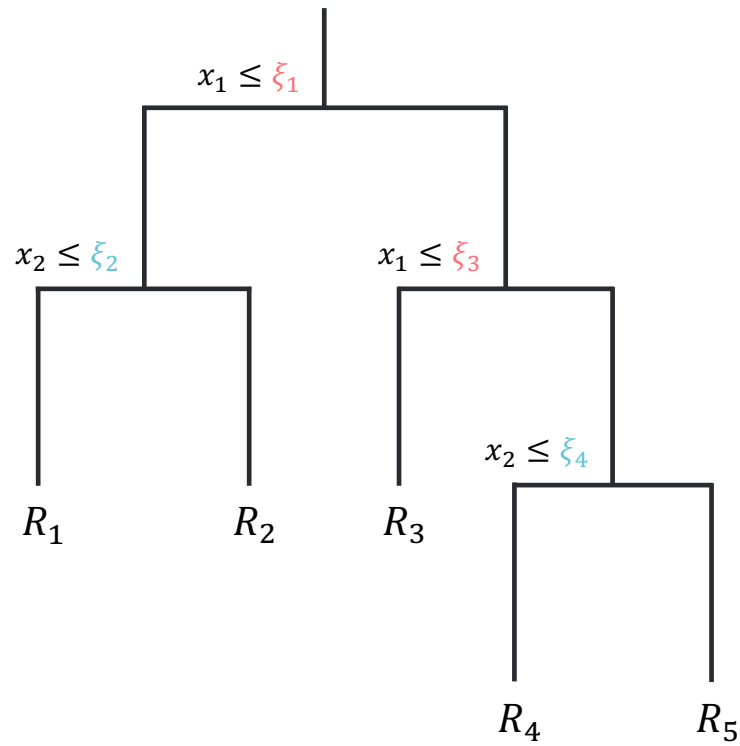
Orthogonal linear transformation  $\mathbf{z} = \mathbf{W}\mathbf{x}$ ,  
where  $\mathbf{W}$  are eigenvectors of  $\mathbf{X}^T\mathbf{X}$  (right  
singular vectors of  $\mathbf{X}$ ):

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{W}^T \Rightarrow \mathbf{Z}^T\mathbf{Z} = \mathbf{\Sigma}^2$$

- (1) The matrix  $\mathbf{X}$  is obtained by stacking the vectors  $\mathbf{x}^{(i)T}$ ,  
 $i = 1, \dots, N$



# DECISION TREES



# CLASSIFICATION AND REGRESSION TREES (CART)

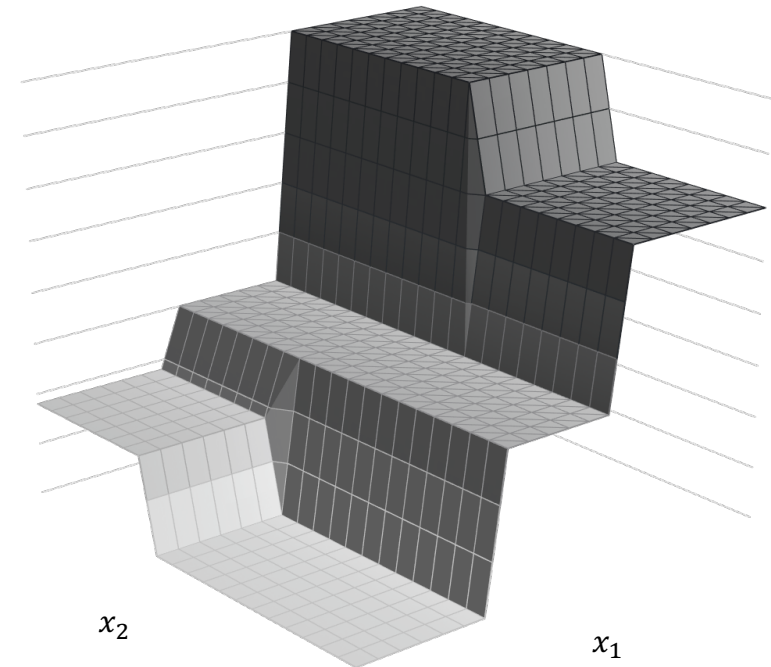
Data set  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N$

Recursive binary splitting:

*Regression trees – RSS cost*

$$\sum_{j=1}^J \sum_{i \in R_j} (y^{(i)} - \hat{y}_{R_j})^2$$

where  $\hat{y}_{R_j}$  is the mean response for the training observations within the  $j$ -th region



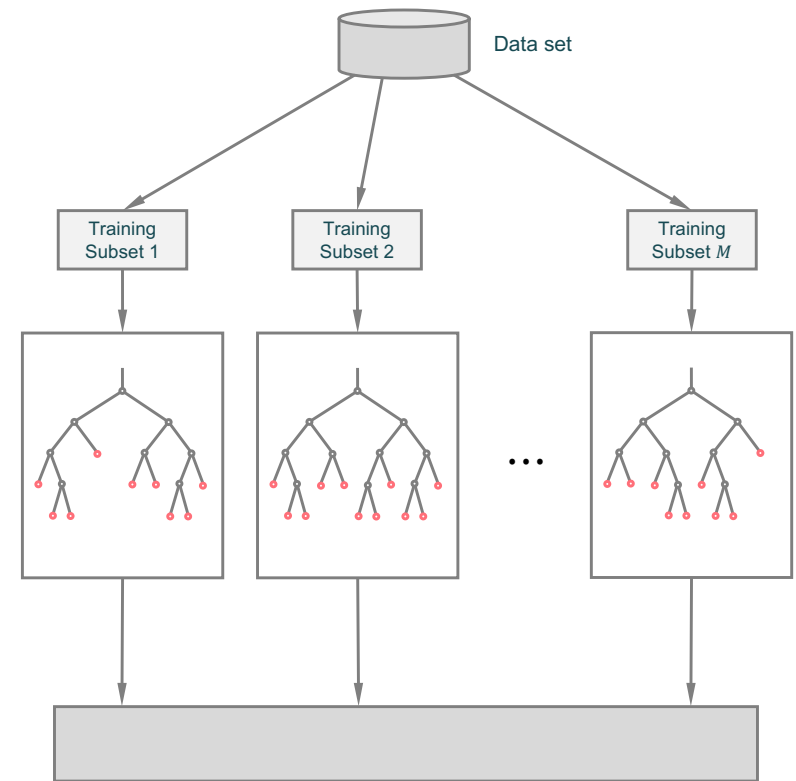


# BAGGING (BOOTSTRAP AGGREGATING)

Improves performance by averaging the estimates produced by  $M$  different trees on different *subsets of the data*, drawn randomly with replacement (*Bootstrap*):

$$f(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x})$$

where  $f_m$  denotes the  $m$ -th tree

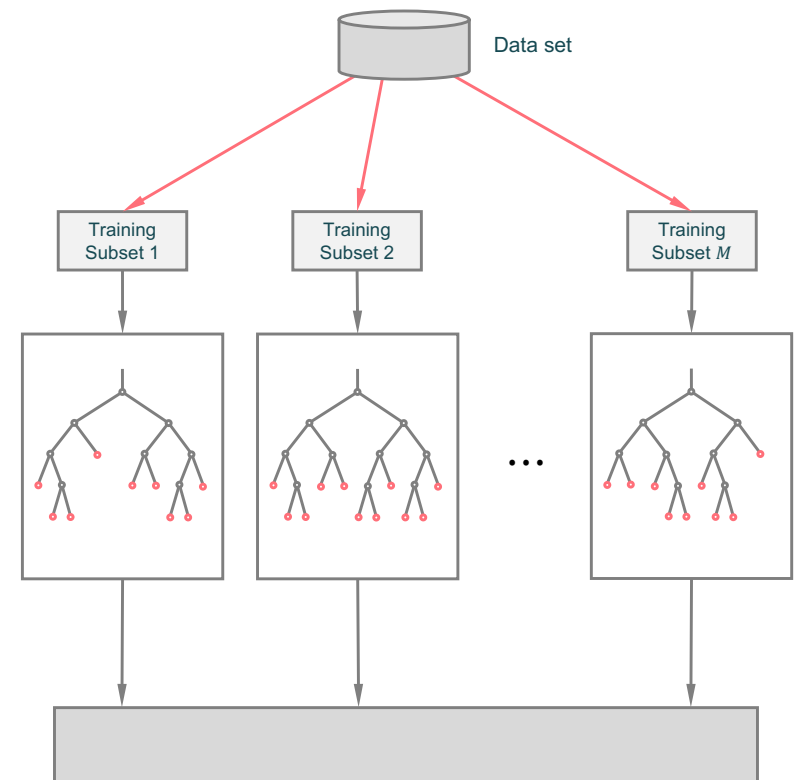


# RANDOM FORESTS

Tries to *decorrelate* the base learners, learning the  $M$  trees based on both

- a different *subsets of the data*, drawn randomly with replacement (*Bootstrap*), and
- a different, randomly chosen *subset of the predictors*  $x_1, \dots, x_d$  (typically  $\sqrt{d}$  predictors) at each split

Reduced variance and reduced test error

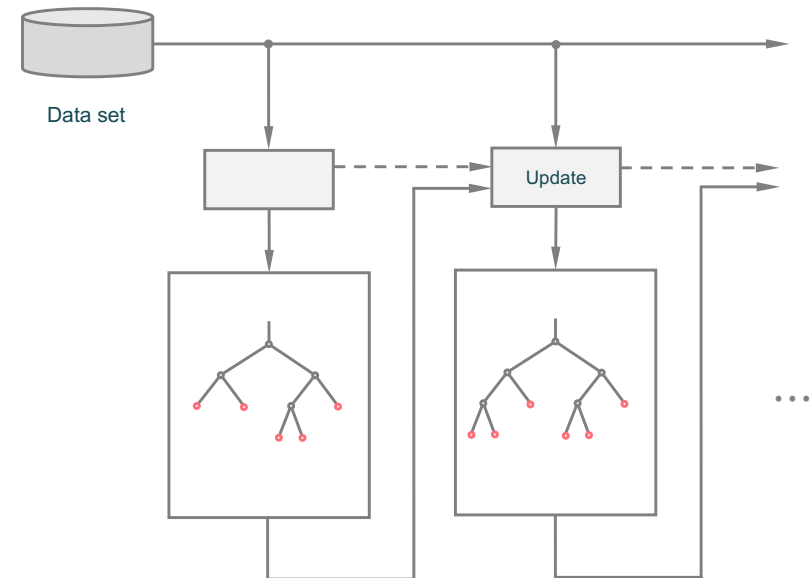


# BOOSTING

Builds an ensemble combining *weak learners*

Similar to Bagging, except that

- trees are grown *sequentially*:  
each tree is grown using information from previously grown trees
- does not involve Bootstrap resampling:  
each tree is fit on a modified version of the original data set



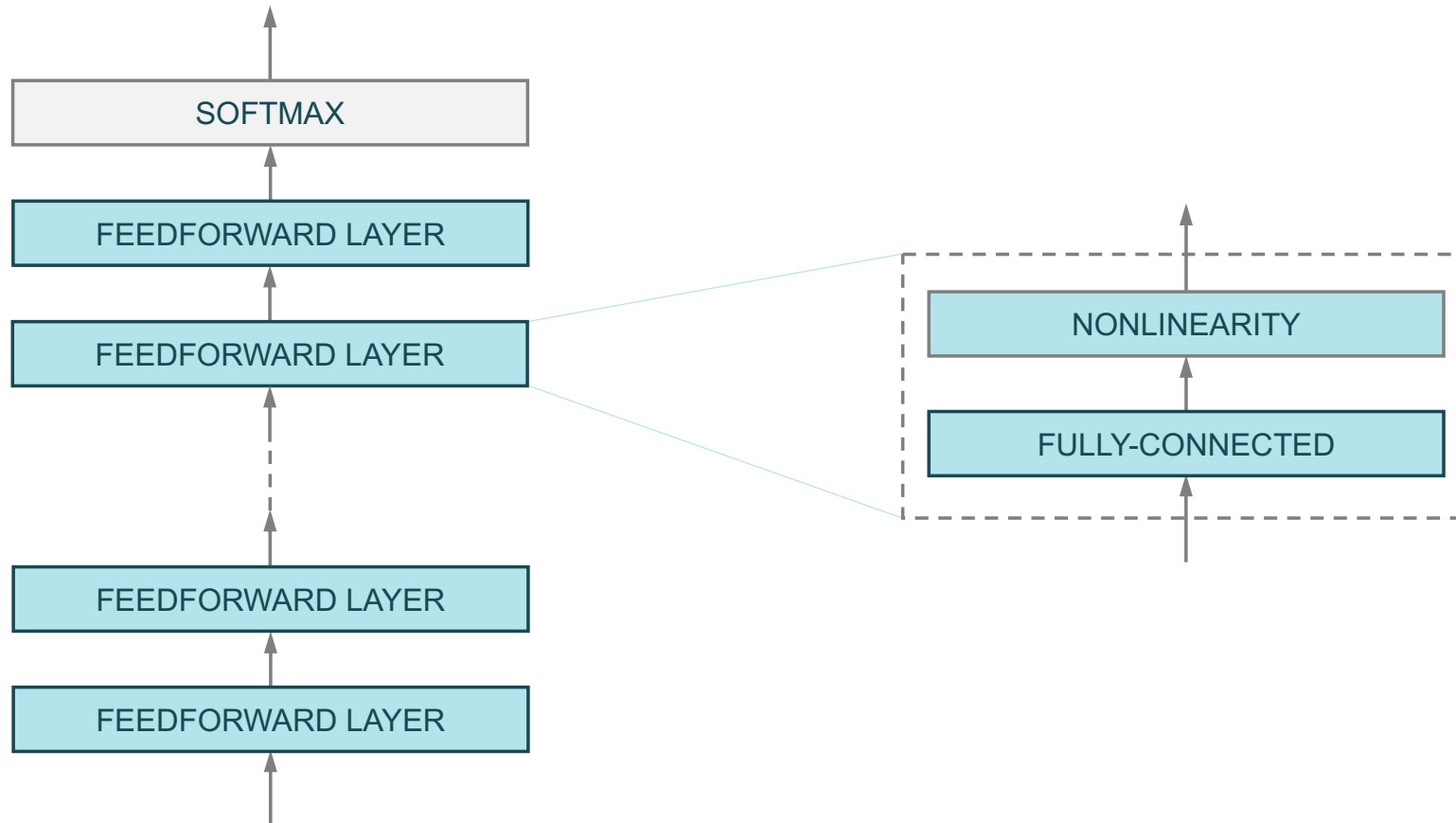


# NEURAL NETWORKS

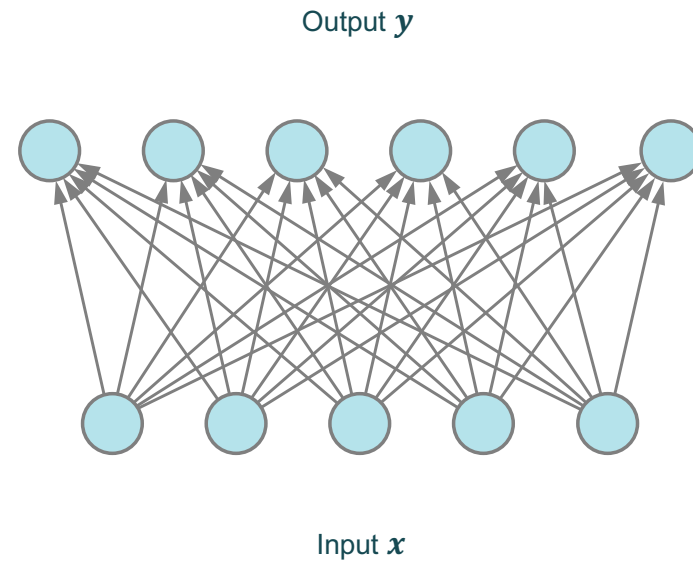
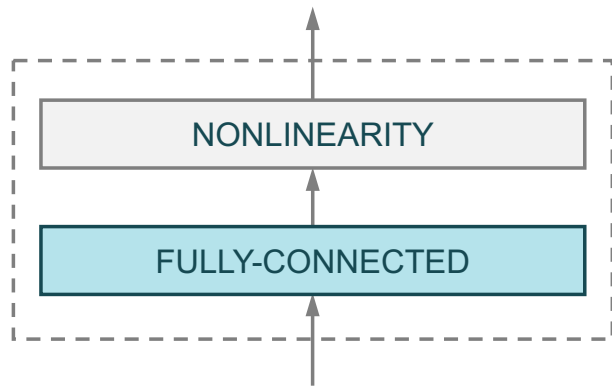
- **Feedforward Neural Networks (FNN)**
- **Convolutional Neural Networks (CNN)**  
*image processing applications (e.g., face recognition, sign and pedestrian detection for autonomous cars)*
- **Recurrent Neural Networks (RNN)**  
*natural language processing applications (e.g., speech transcription, machine translation)*



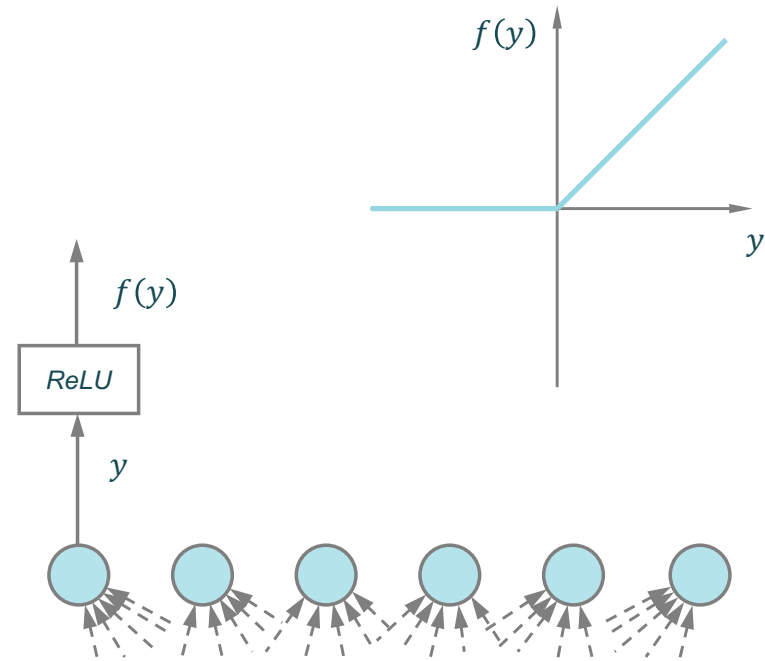
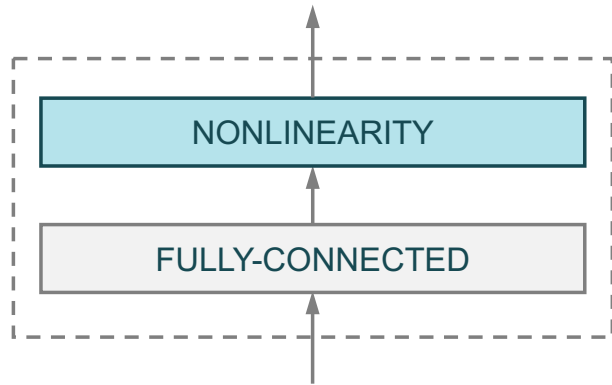
# FEEDFORWARD NEURAL NETWORKS



# FEEDFORWARD NEURAL NETWORKS



# FEEDFORWARD NEURAL NETWORKS



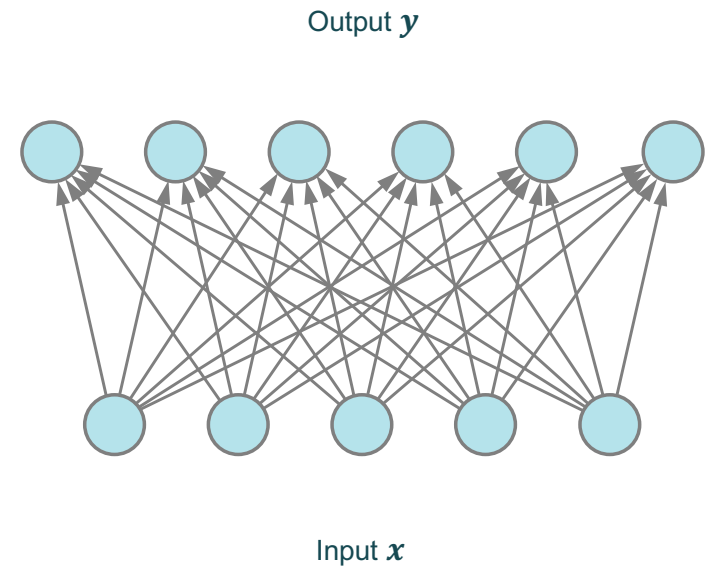
# FEEDFORWARD NEURAL NETWORKS

Fully-connected *linear layer*:

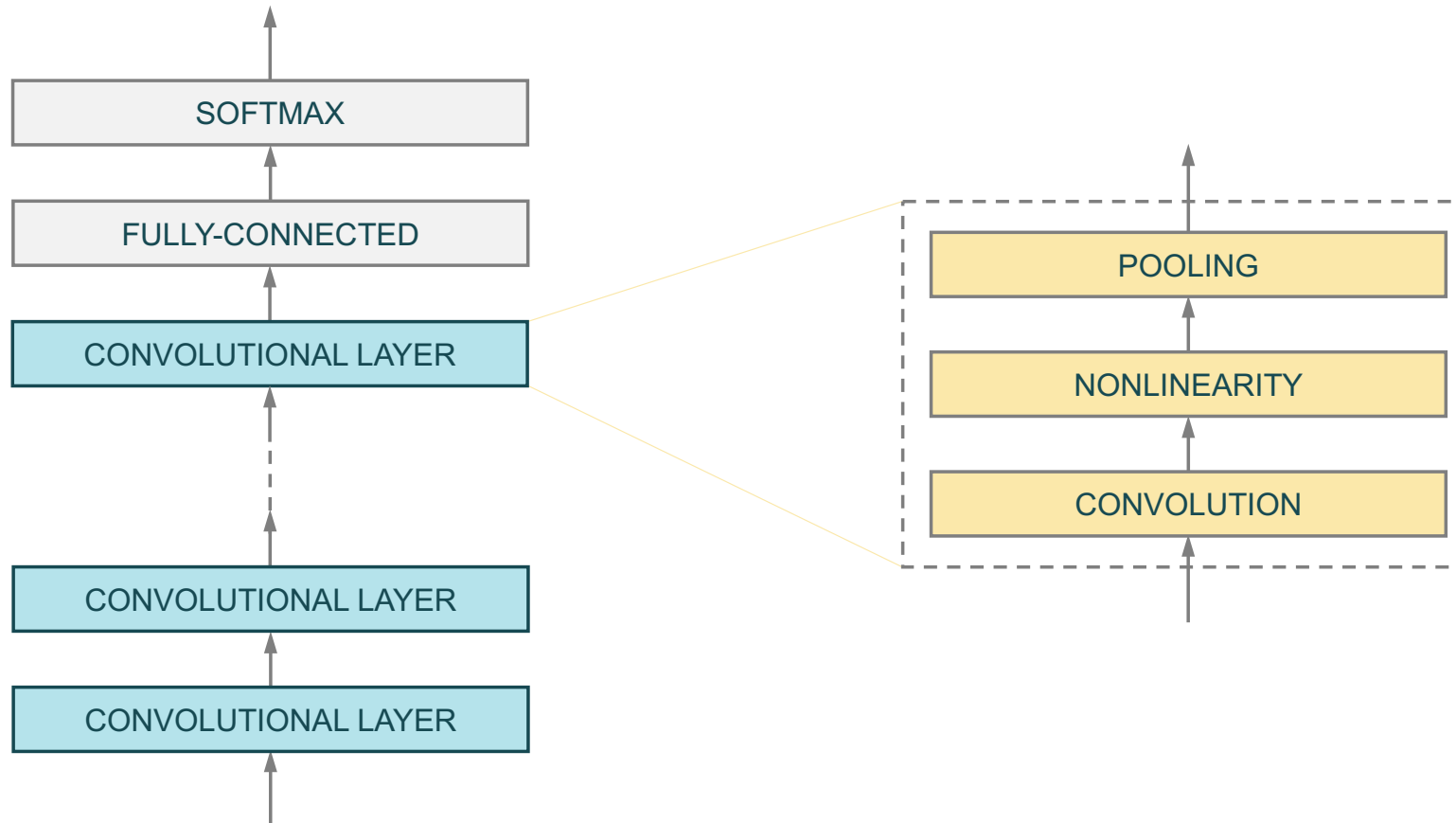
$$y = W^T x + b$$

Nonlinear *activation function*:

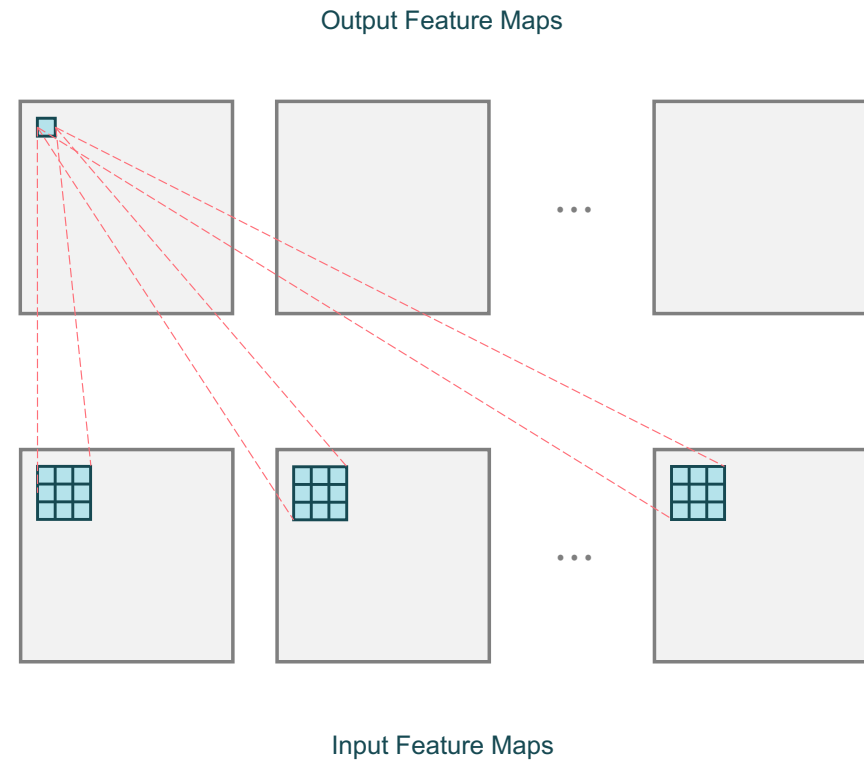
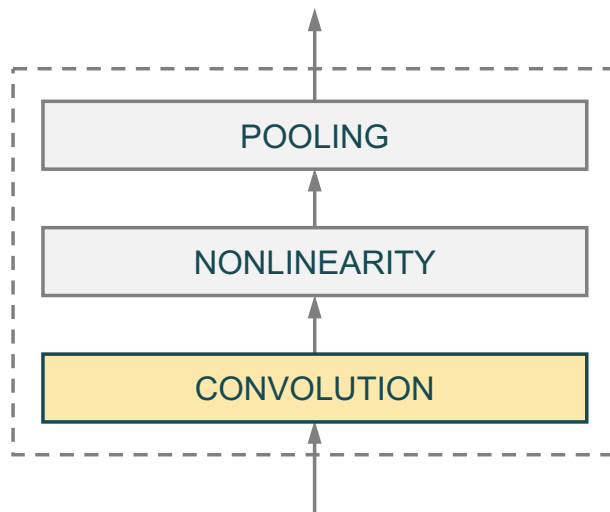
$$z = f(y)$$



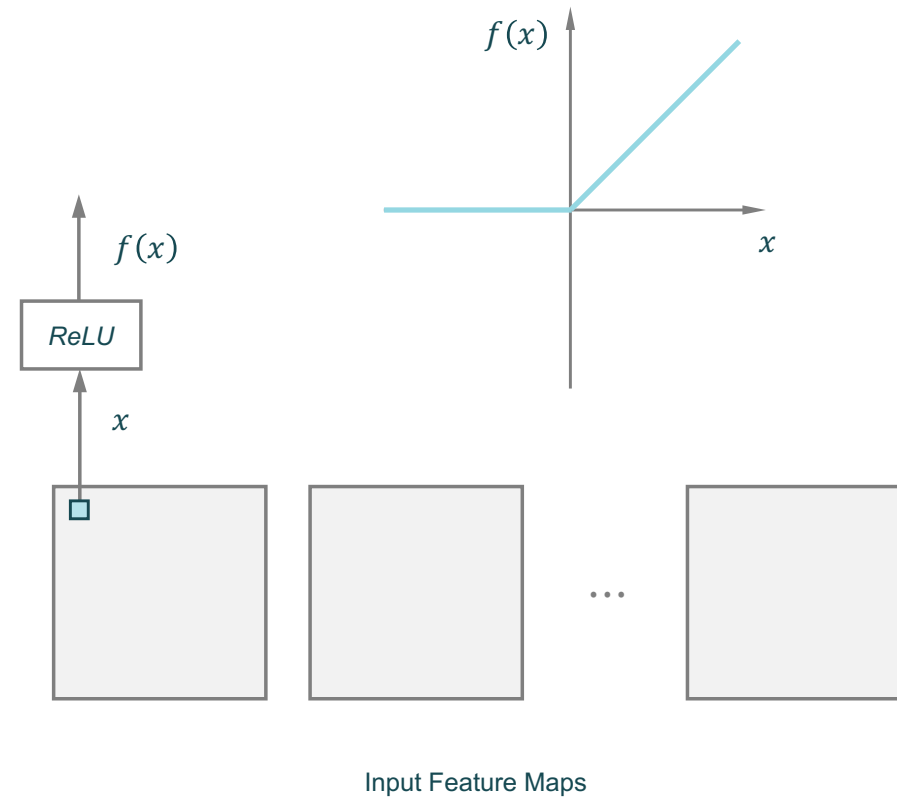
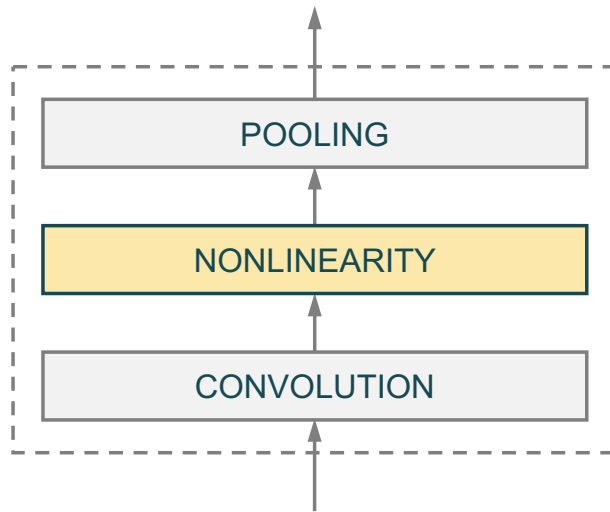
# CONVOLUTIONAL NEURAL NETWORKS



# CONVOLUTIONAL NEURAL NETWORKS

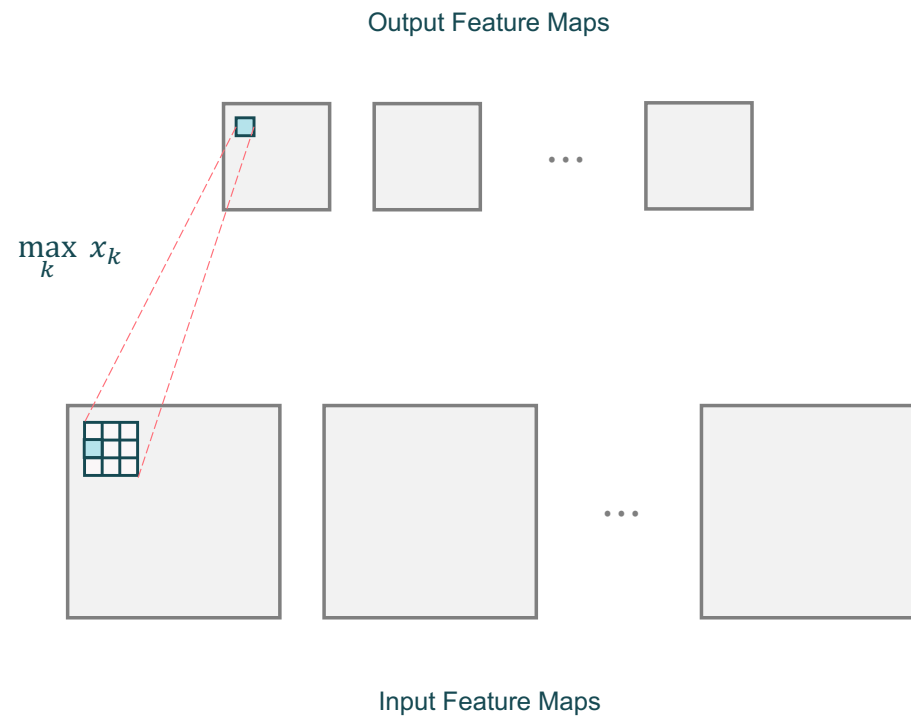
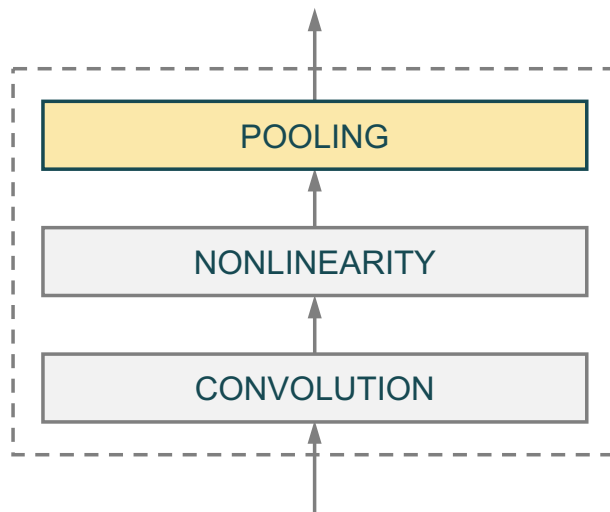


# CONVOLUTIONAL NEURAL NETWORKS

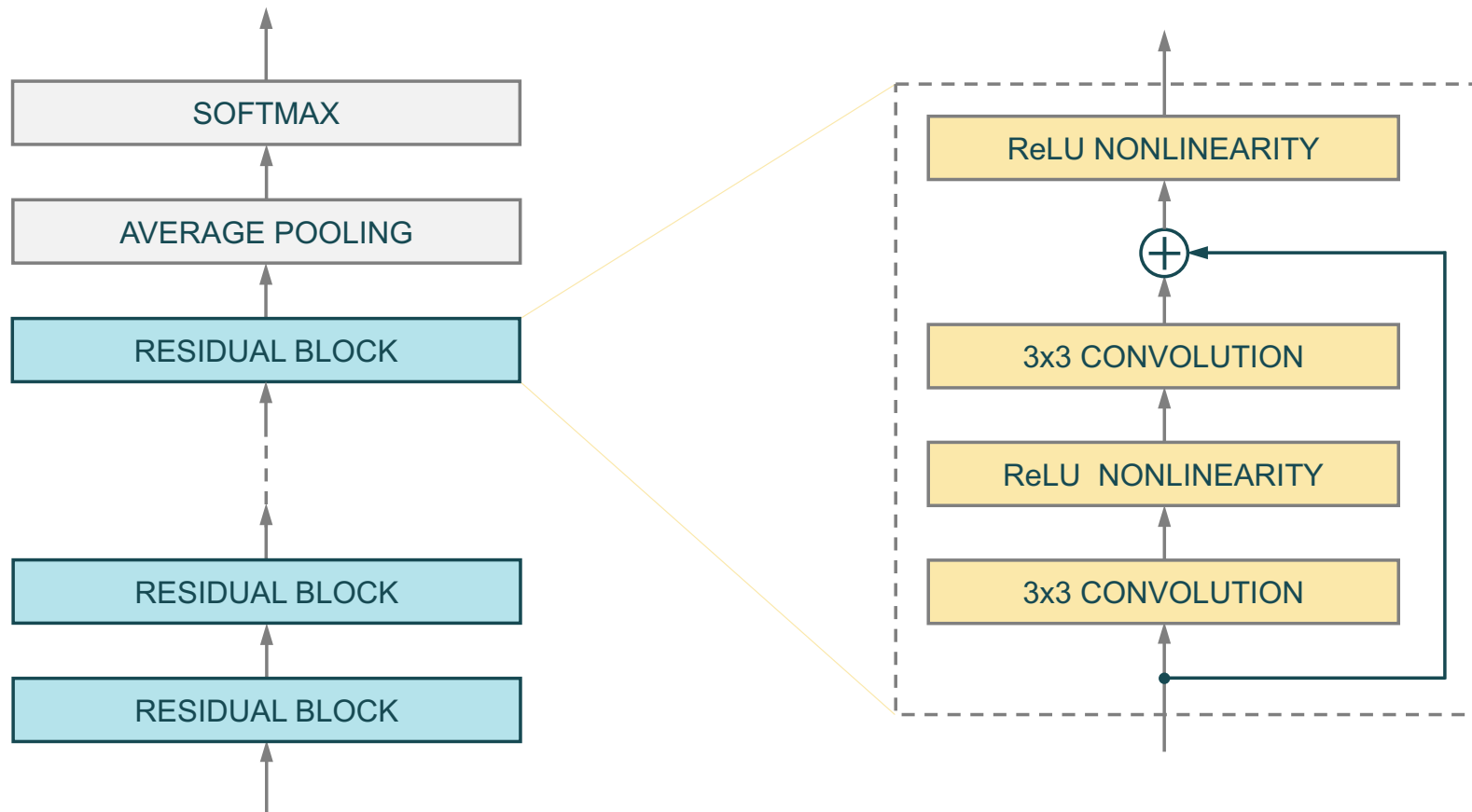




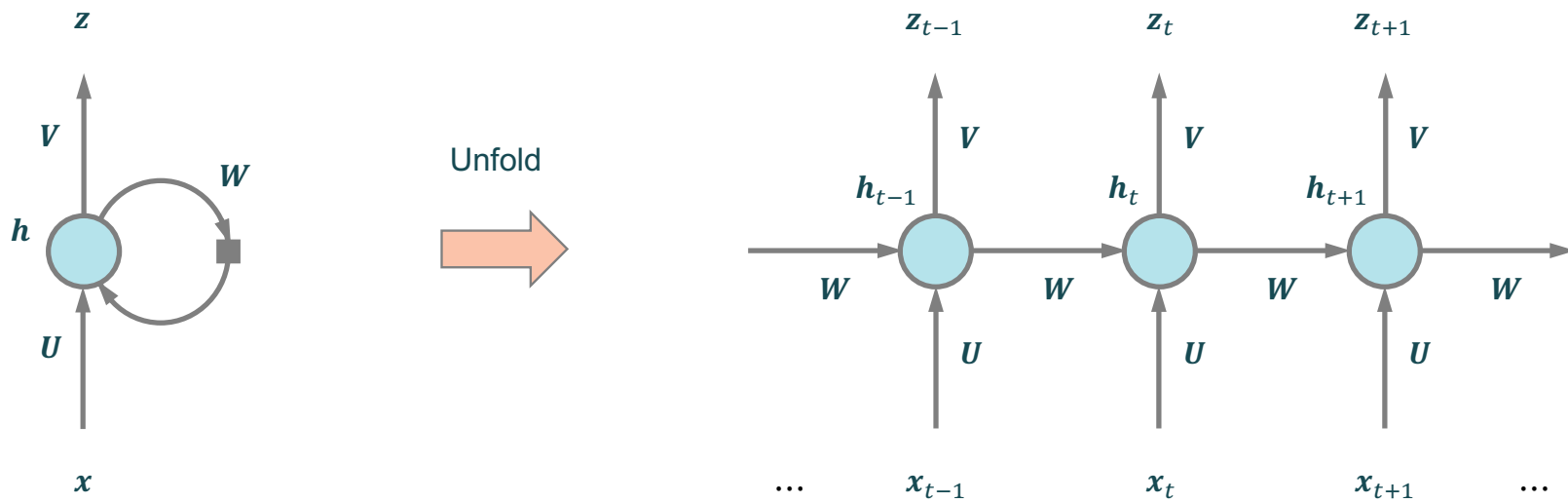
# CONVOLUTIONAL NEURAL NETWORKS



# RESNET ARCHITECTURE



# RECURRENT NEURAL NETWORKS



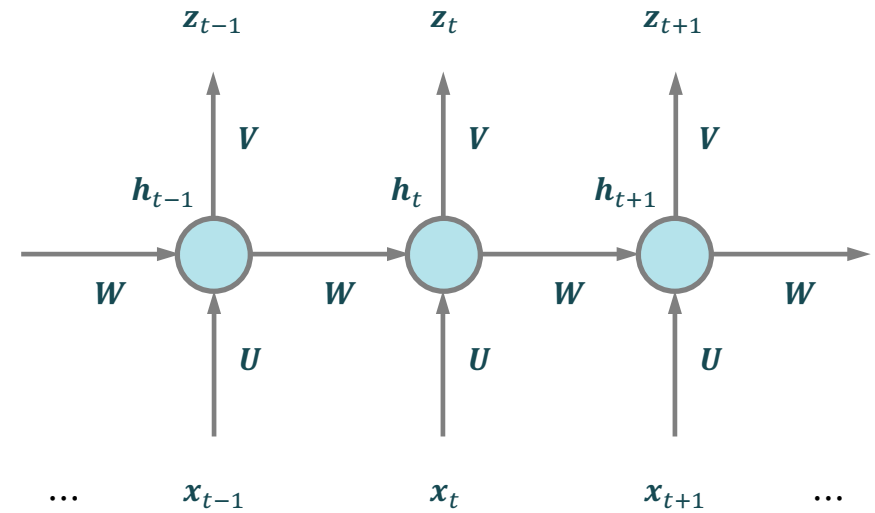
# RECURRENT NEURAL NETWORKS

Efficient parametrization: recurrent application of the same function

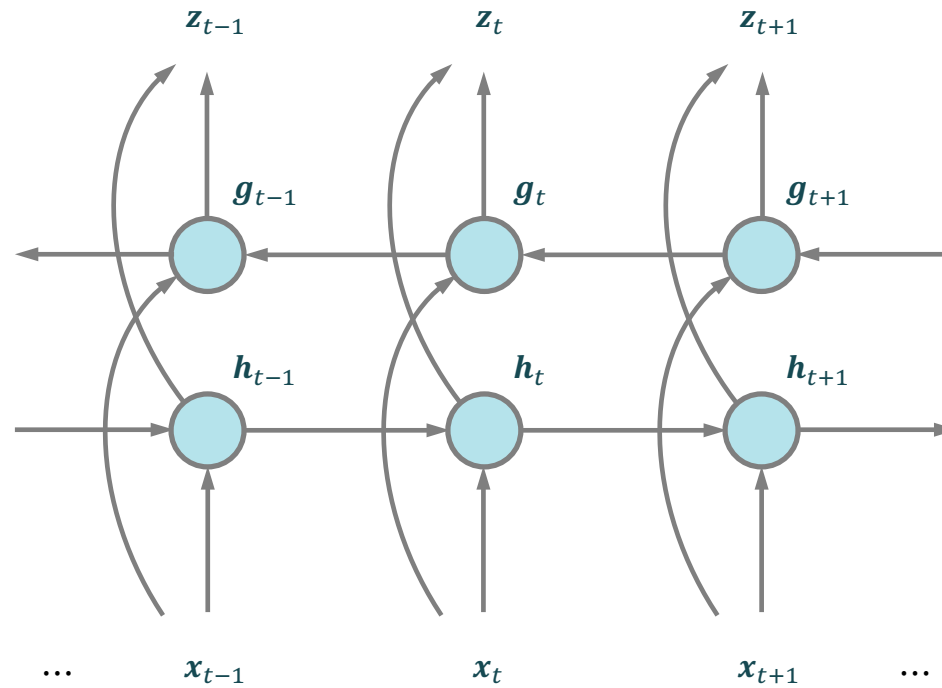
$$y_t = Ux_t + Wh_t + b$$

$$h_t = \tanh(y_t)$$

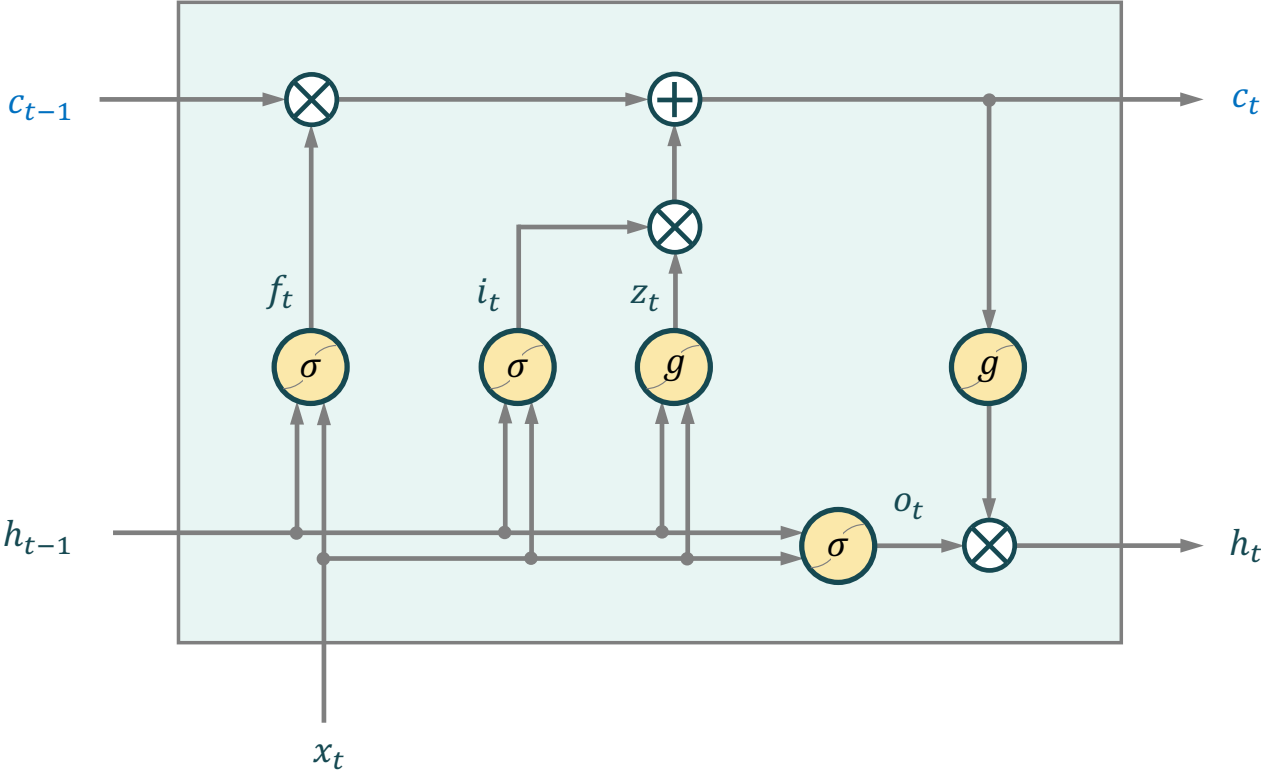
$$z_t = Vh_t + c$$



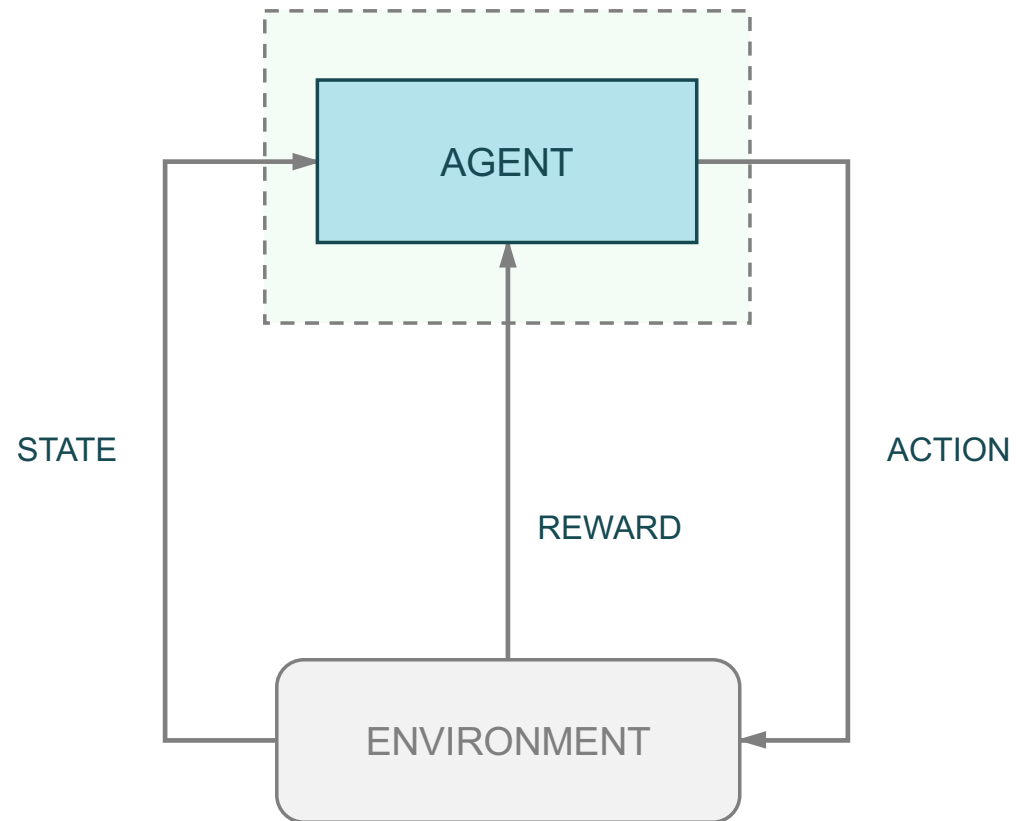
# RECURRENT NEURAL NETWORKS



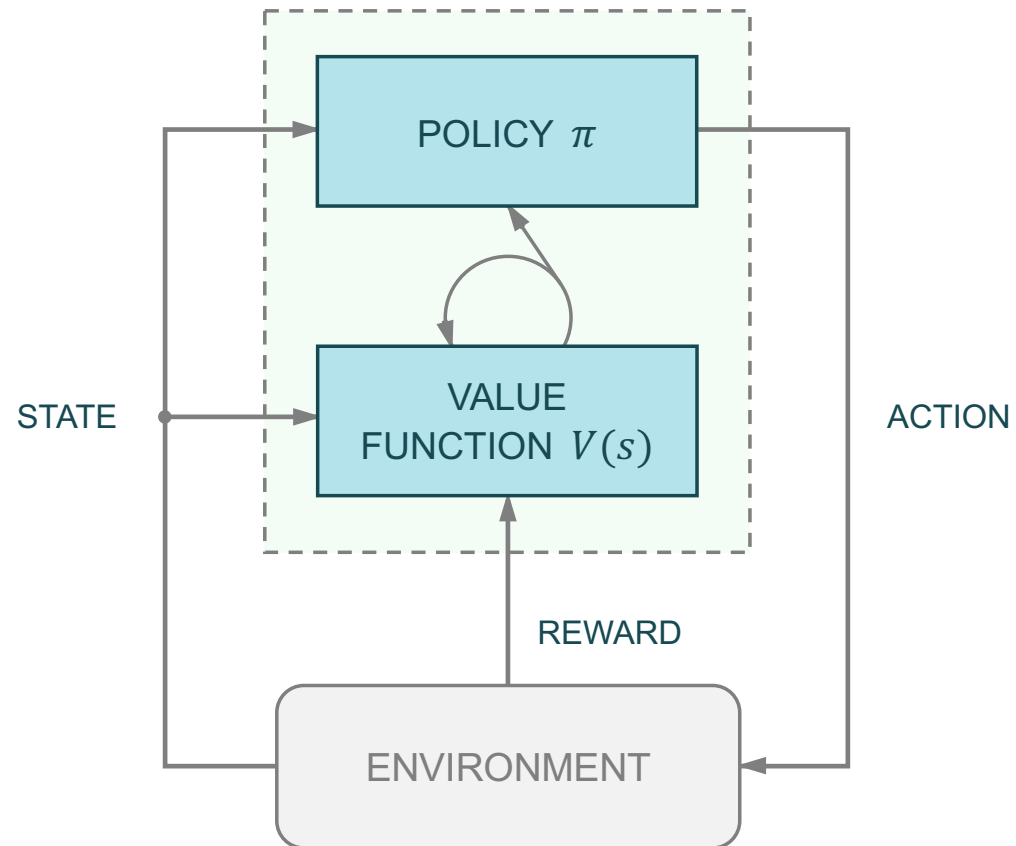
# LONG SHORT-TERM MEMORY (LSTM) CELL



# REINFORCEMENT LEARNING

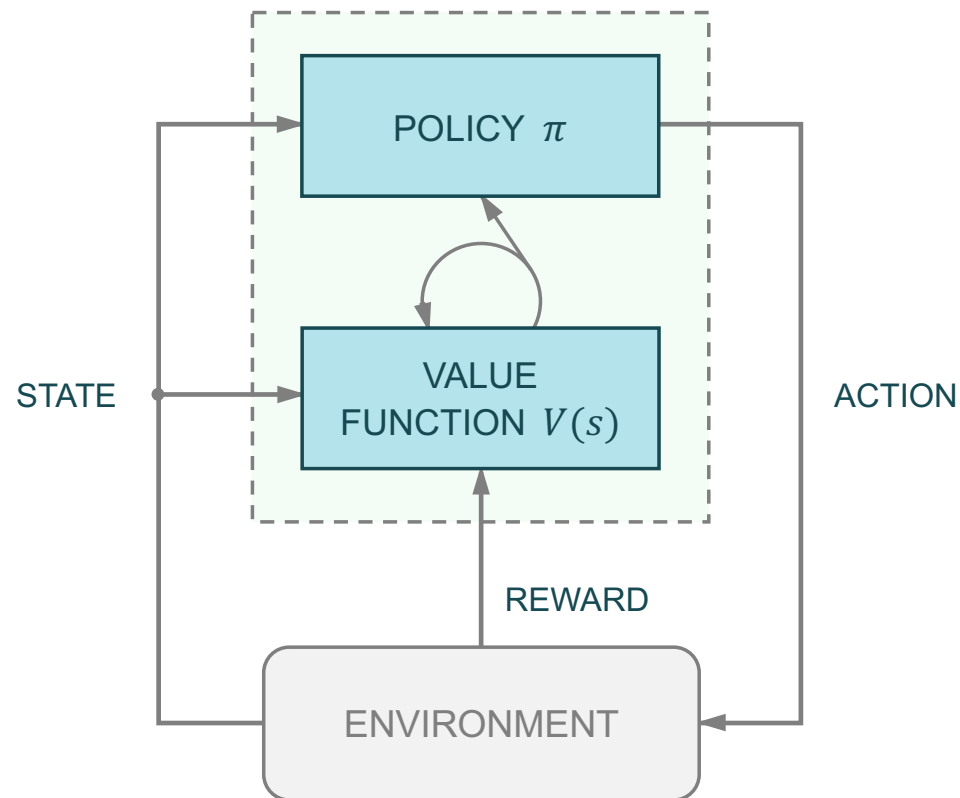


# REINFORCEMENT LEARNING





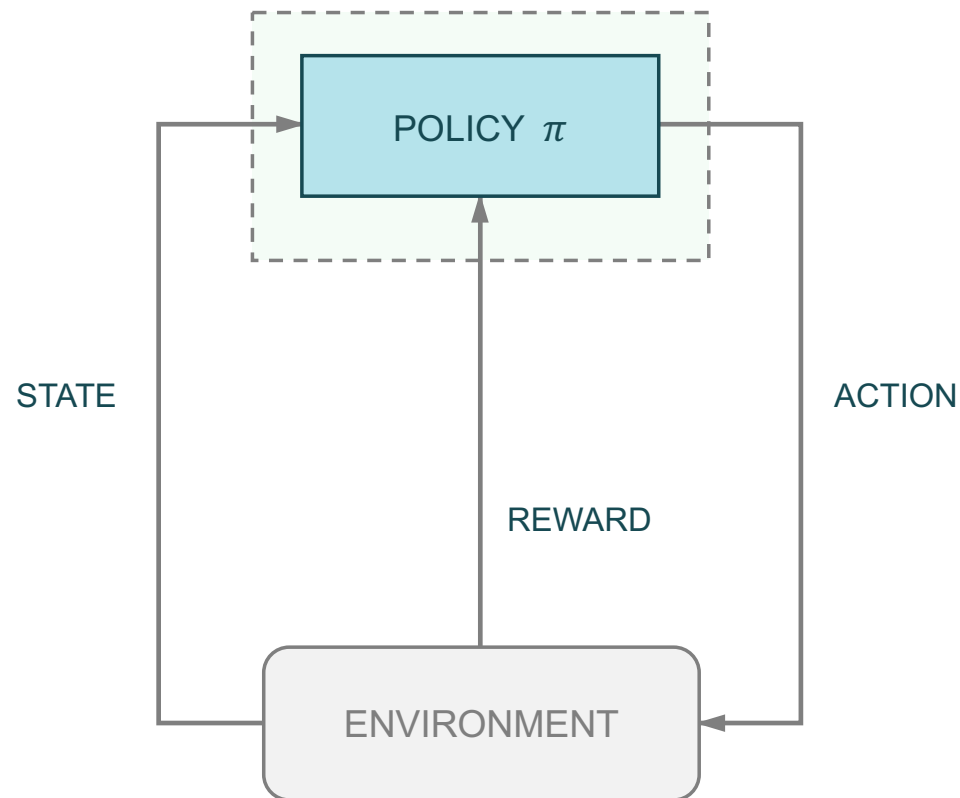
# REINFORCEMENT LEARNING TASKS



## *Reinforcement Learning (RL):*

- Maximization of *expected reward*
- Gradient based optimization by *backpropagation*
- Exploration by *trial-and-error on actions*

# REINFORCEMENT LEARNING TASKS



## *Evolution Strategies (ES):*

- Maximization of *expected reward*
- *Black-box optimization* using *finite differences*
- Exploration by *trial-and-error* on *parameters*

# THANK YOU

**Carlo Luschi**

carlo@graphcore.ai

