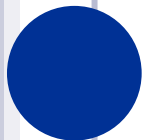


# DEEP LEARNING

Elisabetta Fersini

[elisabetta.fersini@unimib.it](mailto:elisabetta.fersini@unimib.it)



# MAIN TOPICS OF THE COURSE

- Background concepts
- Training Deep Networks:
  - Objective functions
  - Activation Functions
  - Regularization
  - Gradient-based optimization
- Deep Learning in Practice
  - Tips and Tricks
- Focus on Deep Architectures (hands on):
  - Variational Models
  - Transformers



# PEOPLE

- Prof. Elisabetta Fersini
  - Associate professor of CS
  - Research field:
    - Machine Learning
    - Natural Language Processing
  
- Dr. Cristiano De Nobili
  - Theoretical Physicist
  - Lead AI Scientist @ Pi School
  - Speaker



# SCHEDULE

## ○ LECTURES

- Friday, 08/07 (14:00 – 18:00), Building U14, Room T023
- Wednesday 13/07 (14:00– 18:00), Building U14, Room T023
- Monday 18/07 (14:00– 18:00), Building U14, Room T024
- Tuesday 28/07 (10:00 – 12:00, 14:00 – 16:00), Building U14, Room T023
- Friday 29/07 (10:00 – 12:00, 14:00 – 16:00), Building U24, Room C01

## ○ FINAL EXAM

- Friday 30/09 (room and time schedule to be announced)



# FINAL EXAM

- 15 minutes presentation related to:
  - An in depth analysis of a **theoretical** paper



# WHAT IS DEEP LEARNING?

## ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



## MACHINE LEARNING

Ability to learn without explicitly being programmed



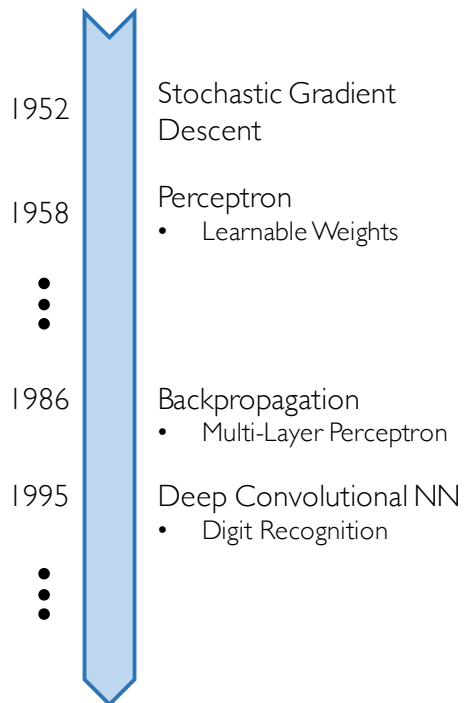
## DEEP LEARNING

Learn underlying features in data using neural networks

3 1 3 4 7 2  
1 7 4 2 3 5



# WHY NOW?



Neural Networks date back decades, so why the resurgence?

## 1. Big Data

- Larger Datasets
- Easier Collection & Storage

IMAGENET



WIKIPEDIA  
The Free Encyclopedia



## 2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



## 3. Software

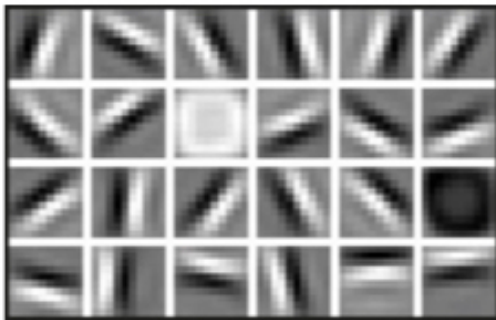
- Improved Techniques
- New Models
- Toolboxes



# WHY DEEP LEARNING?

- Hand engineered features are time consuming, brittle and not scalable in practice
  - Can we learn the **underlying features** directly from data?

Low Level Features



Lines & Edges

Mid Level Features



Eyes & Nose & Ears

High Level Features



Facial Structure





# DOMAINS OF INTEREST FOR DEEP LEARNING

## Image Recognition

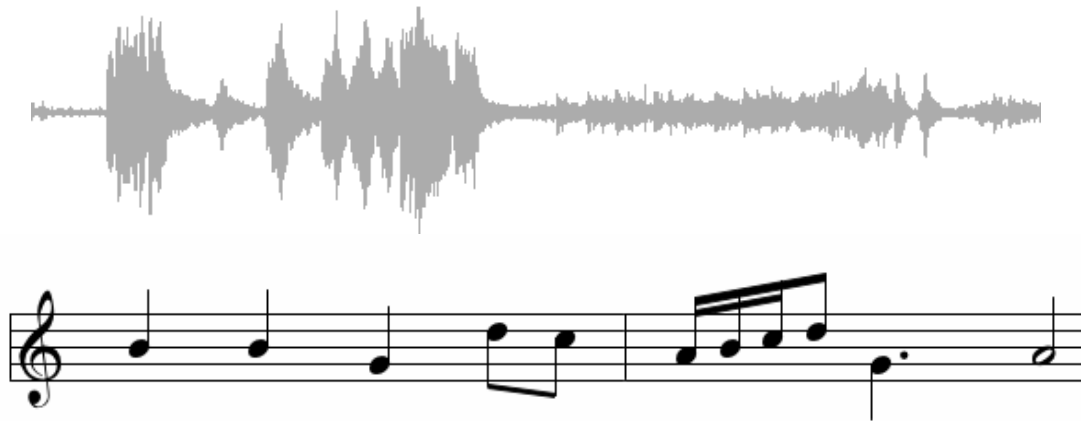


mite	container ship	motor scooter	leopard
mite	container ship	motor scooter	leopard
black widow	lifeboat	go-kart	jaguar
cockroach	amphibian	moped	cheetah
tick	fireboat	bumper car	snow leopard
starfish	drilling platform	golfcart	Egyptian cat



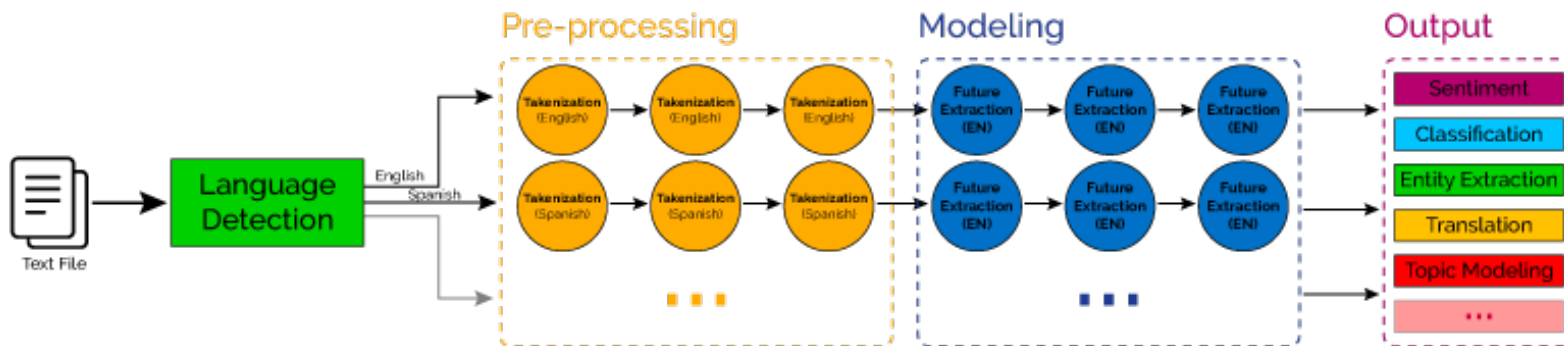
# DOMAINS OF INTEREST FOR DEEP LEARNING

Music Generation

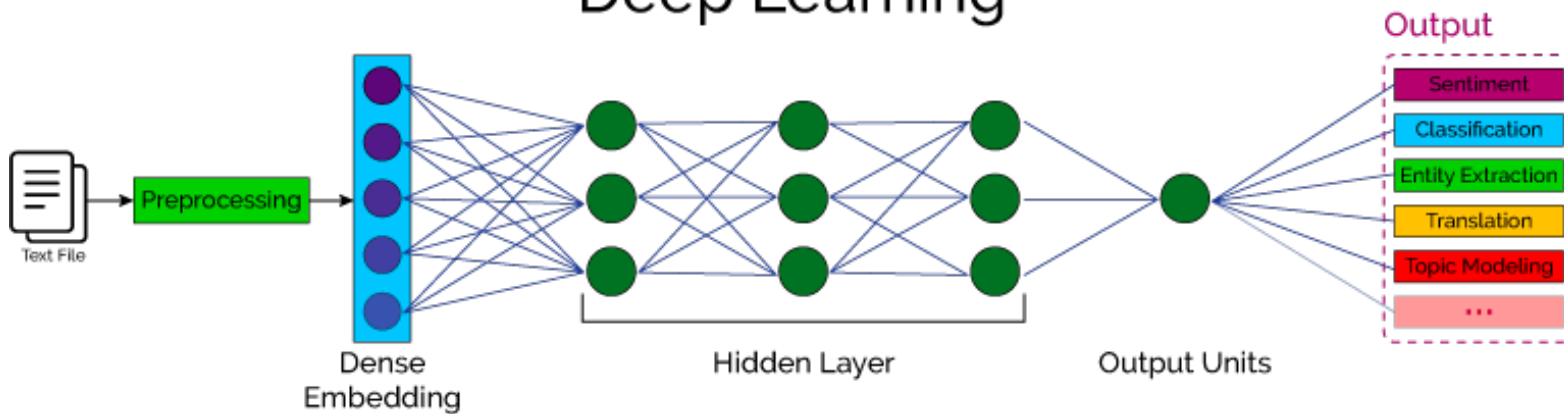


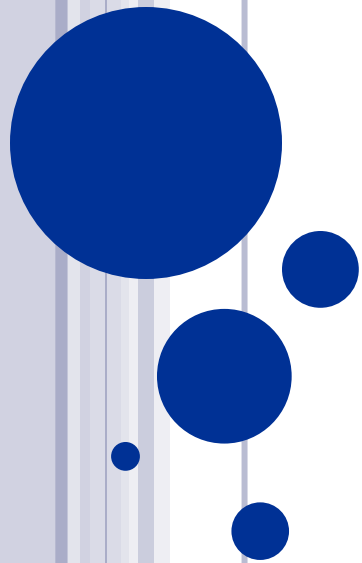
# DOMAINS OF INTEREST FOR DEEP LEARNING

## Classical NLP



## Deep Learning





## BACKGROUND CONCEPTS

# INTRODUCTION

- Deep Learning is grounded on (Artificial) Neural Networks.
- Neural Networks are a class of supervised machine learning algorithms
- We will have a quick introduction to:
  - supervised machine learning terminology and practices
  - linear and log-linear models for binary and multi-class classification



# SUPERVISED LEARNING AND PARAMETERIZED FUNCTIONS

- The essence of supervised machine learning is the creation of *functions* that can look at examples (instances) and produce generalizations (good prediction on unseen data).
- As searching over the set of all possible functions is a very hard problem, we often restrict ourselves to search over specific families of functions called *hypothesis classes*



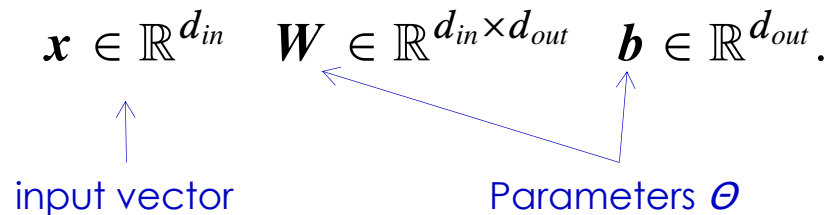
# SUPERVISED LEARNING AND PARAMETERIZED FUNCTIONS

- A common hypothesis class is represented by high-dimensional **linear function**, i.e. functions of the form:

$$f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

$$\mathbf{x} \in \mathbb{R}^{d_{in}} \quad \mathbf{W} \in \mathbb{R}^{d_{in} \times d_{out}} \quad \mathbf{b} \in \mathbb{R}^{d_{out}} .$$

input vector Parameters  $\theta$



- The goal of the learner is to find the values of the parameters  $\mathbf{W}$  and  $\mathbf{b}$  such that the function behaves as intended on a collection of input instances  $x_1, \dots, x_k$  and the corresponding output labels  $y_1, \dots, y_k$ .



# LINEAR MODELS

- Binary classification:

$$f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$$

Diagram illustrating the components of the linear function  $f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{w} + b$ :

- $\mathbf{x}$ : input vector
- $\mathbf{w}$ : weight vector
- $b$ : scalar

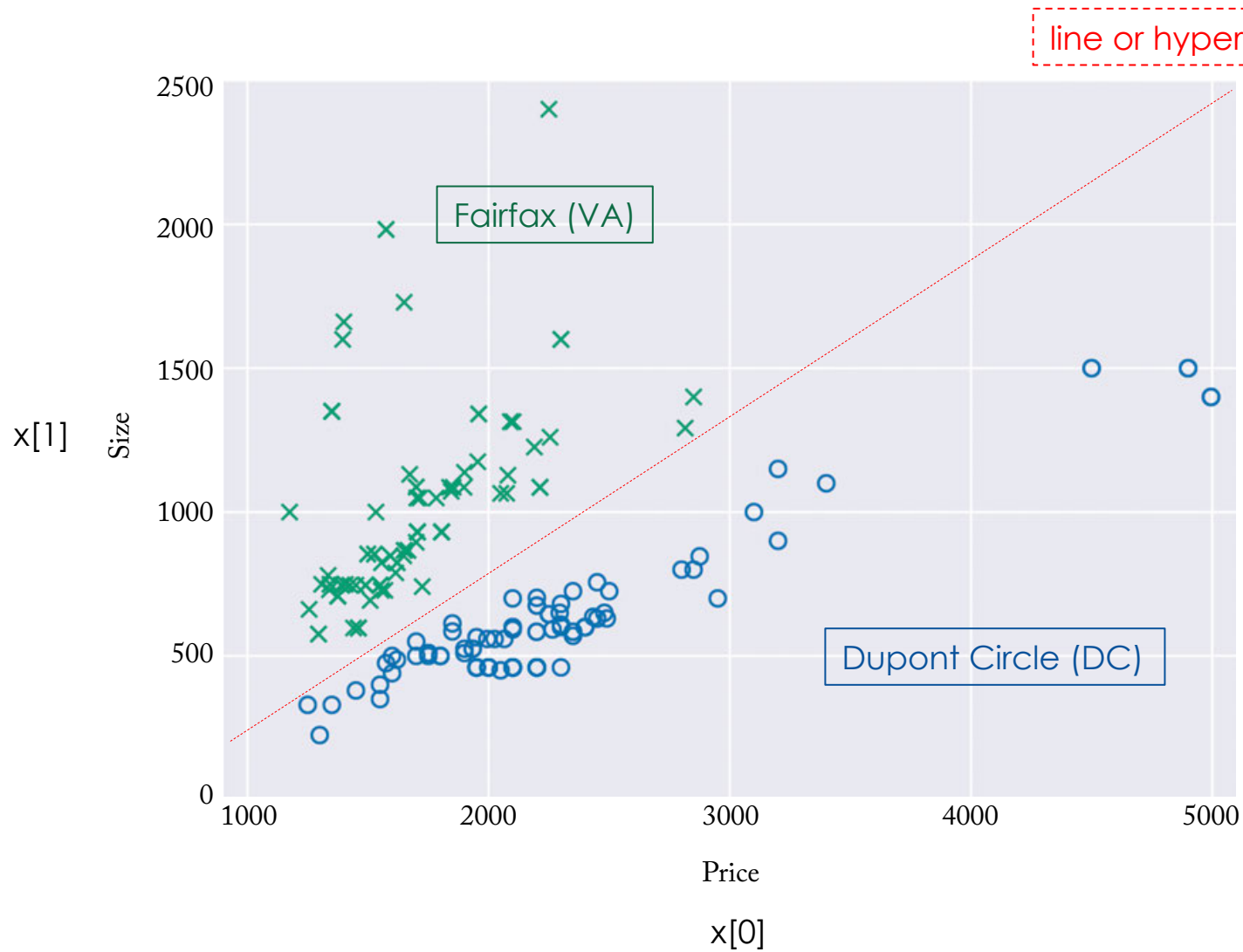
- The range of the linear function is  $[-\infty; +\infty]$ 
  - We need to map the in order to have only values +1 and -1
    - *sign function*

$$\hat{y} = \text{sign}(f(\mathbf{x})) = \text{sign}(\mathbf{x} \cdot \mathbf{w} + b)$$





# LINEAR MODELS



# LINEAR MODELS

- In our example, the linear model is as follows:

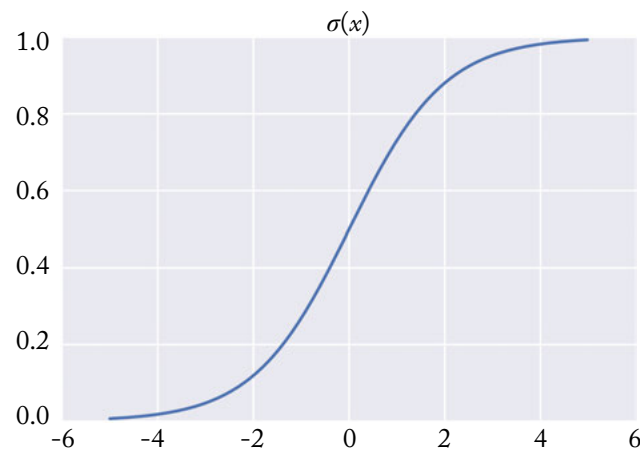
$$\begin{aligned}\hat{y} &= \text{sign}(f(\mathbf{x})) = \text{sign}(\mathbf{x} \cdot \mathbf{w} + b) \\ &= \text{sign}(\text{size} \times w_1 + \text{price} \times w_2 + b)\end{aligned}$$

- where the prediction is **Fairfax** if  $\hat{y} \geq 0$  and **Dupont Circle** if  $\hat{y} \leq 0$ .
- The goal of learning is setting the values of  $w_1$ ,  $w_2$  and  $b$  such that all  $\hat{y}$  are correct for all data points we observe.
  - datasets are not always linearly separable, and therefore they are beyond the hypothesis class of linear classifiers
    - Possible solutions: move to a higher dimension, move to a richer hypothesis class, etc..



# LINEAR MODELS

- **Log-Linear Binary Classification**: the need of a confidence of the decision
  - Instead of mapping to  $\{-1,+1\}$ , we map to the range  $[0,1]$  by using a squashing function such as the *sigmoid function*



- Our prediction becomes

$$\hat{y} = \sigma(f(\mathbf{x})) = \frac{1}{1 + e^{-(\mathbf{x} \cdot \mathbf{w} + b)}}$$



# LINEAR MODELS

- **Multi-class classification:**

- we want to predict an instance as belonging to one of  $k$  different classes

- A possible solution is to consider  $k$  weight vectors  $w^1, w^2, \dots, w^k$  and biases

- And predict the class with the highest score:

$$\hat{y} = f(x) = \operatorname{argmax}_{k \in K} \mathbf{x} \cdot \mathbf{w}^k + b^k$$

- The sets of parameters  $w^k \in R^{in}$  and  $b^k$  can be arranged as a matrix  $W \in R^{in \times k}$  and vector  $b \in R^k$ , leading to:

$$\hat{\mathbf{y}} = f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W} + \mathbf{b}$$

$$\text{prediction} = \hat{y} = \operatorname{argmax}_k \hat{\mathbf{y}}_{[k]}$$

- where  $\hat{\mathbf{y}} \in R^k$  is a vector of the scores assigned by the model to each class



# LINEAR MODELS

- Log-linear multi-class classification: transform the score vector into a probability estimate through the *softmax function*
- Our prediction becomes:

$$\hat{y} = \text{softmax}(\mathbf{x}\mathbf{W} + \mathbf{b})$$
$$\hat{y}_{[k]} = \frac{e^{(\mathbf{x}\mathbf{W} + \mathbf{b})_{[k]}}}{\sum_j e^{(\mathbf{x}\mathbf{W} + \mathbf{b})_{[j]}}}$$

- forcing the values in  $\hat{y}$  to be positive and sum to 1, making them interpretable as a probability distribution.



# TRAINING AS OPTIMIZATION

- The goal of a **training algorithm** is to return a function  $f()$  that accurately maps the input vectors to the corresponding labels
  - i.e. a function  $f()$  such that  $\hat{y} = f(x)$  over the training data is accurate
- In order to measure how a function  $f()$  is accurate, we need to introduce the concept of **loss function**
  - Formally, a loss function  $L(\hat{y}, y)$  assign a **numerical score** to a predicted output  $\hat{y}$  given the true expected output  $y$
  - The parameters of the learned function ( $W$  and  $b$ ) are then set to **minimize** the **loss**  $L$  over the training samples



# TRAINING AS OPTIMIZATION

- Given a **labelled training set**  $(x_{1:n}, y_{1:n})$ , a per-instance **loss function**  $\mathcal{L}$  and a **parameterized function**  $f(x; \Theta)$ , we define the global loss as follows:

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n L(f(x_i; \Theta), y_i)$$

- The main goal of the training algorithm is to set the values of the parameters  $\Theta$  such that the value of  $\mathcal{L}$  is **minimized**:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \mathcal{L}(\Theta) = \underset{\Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n L(f(x_i; \Theta), y_i)$$

- Here we are minimizing the loss at all costs, which may result in **overfitting**



# TRAINING AS OPTIMIZATION

- To overcome this issue, we could pose a soft restriction on the form of the solution through a **regularization function**:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \left( \overbrace{\frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)}^{\text{loss}} + \overbrace{\lambda R(\Theta)}^{\text{regularization}} \right)$$

- Where:
  - $R(\Theta)$  corresponds to a scalar that reflects the parameter “complexity” (that we want to keep low for efficiency reasons)
  - $\lambda$  is an hyperparameter that control the trade-off between bias and variance





# BIAS-VARIANCE TRADE-OFF

Who knows the bias-variance trade-off?

## Bias - Variance Tradeoff

$$\text{Error}(x) = \left( \underset{\substack{\downarrow \\ \text{predicted}}}{E[\hat{f}(x)]} - \underset{\substack{\downarrow \\ \text{true}}}{f(x)} \right)^2 + E \left[ \underset{\substack{\downarrow \\ \text{predicted}}}{\hat{f}(x)} - \underset{\substack{\downarrow \\ \text{average predicted value}}}{E[\hat{f}(x)]} \right]^2 + \underset{\substack{\downarrow \\ \text{irreducible error}}}{\sigma_e^2}$$

Bias<sup>2</sup>

How much predicted values differ from true values.

Variance

How predictions made on the same value vary on different realizations of the model

BY CHRIS ALBON



# BIAS-VARIANCE TRADE-OFF: MATHEMATICALLY

- Let the variable we are trying to predict as  $Y$  and other covariates as  $X$ .

$$Y = f(X) + e$$

- The **expected squared error** is

$$Err(x) = E \left[ (Y - \hat{f}(x))^2 \right]$$

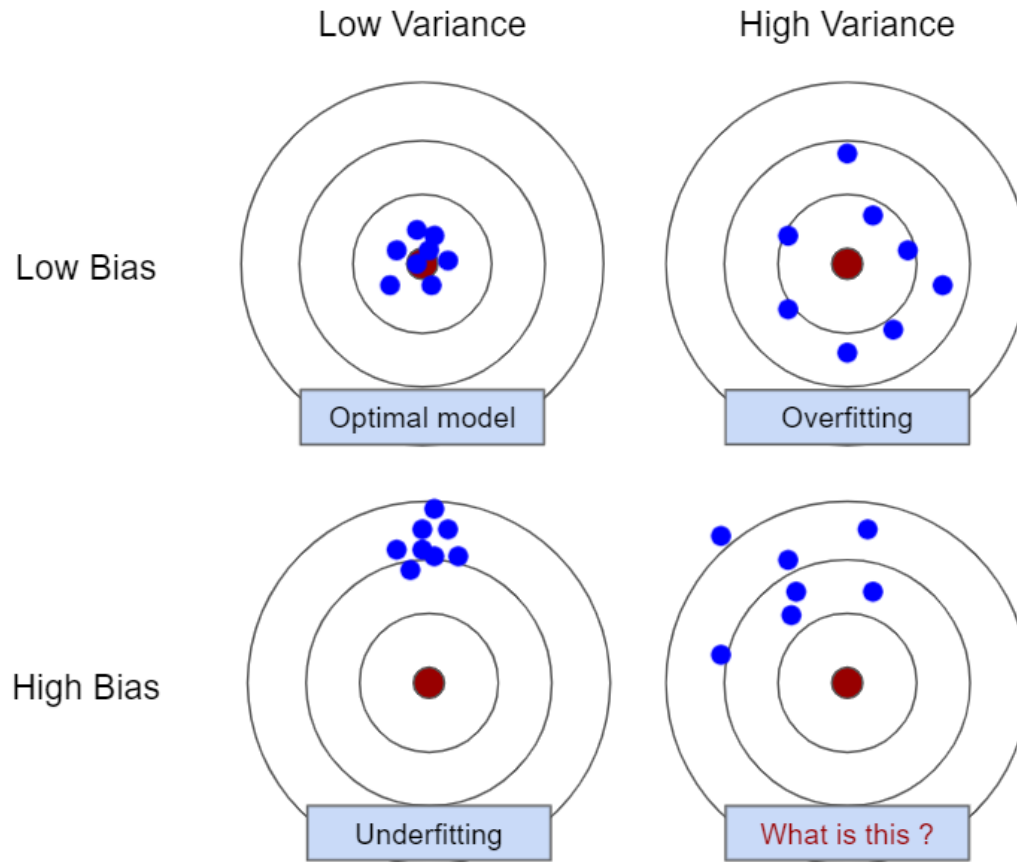
- Which can be further decomposed as

$$Err(x) = \left( E[\hat{f}(x)] - f(x) \right)^2 + E \left[ \left( \hat{f}(x) - E[\hat{f}(x)] \right)^2 \right] + \sigma_e^2$$

$$Err(x) = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

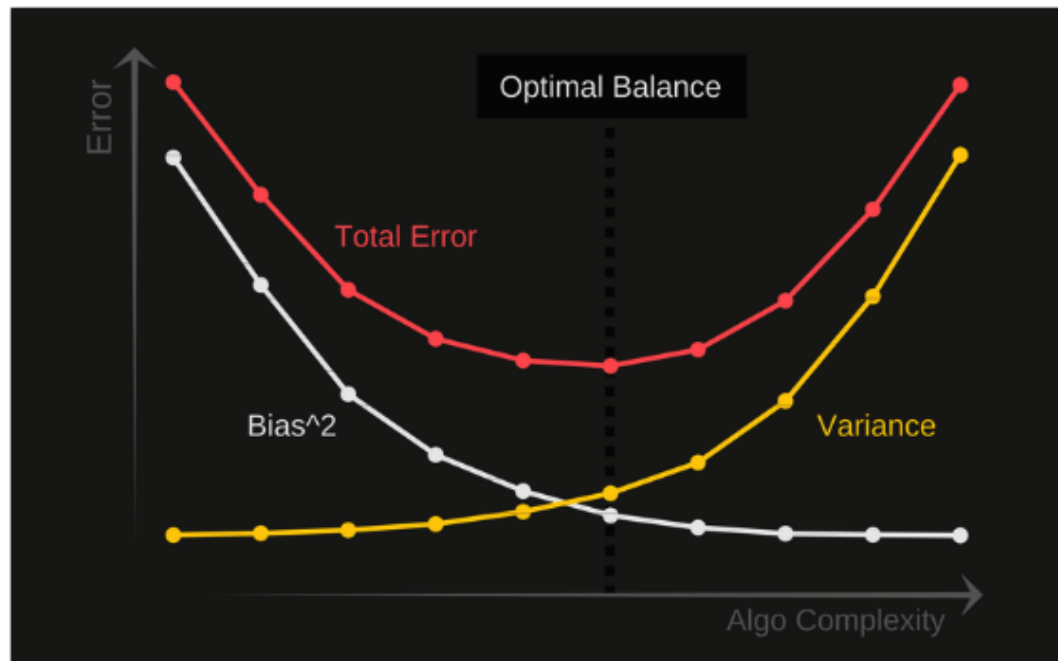


# BIAS-VARIANCE TRADE-OFF: BULLS-EYE DIAGRAM

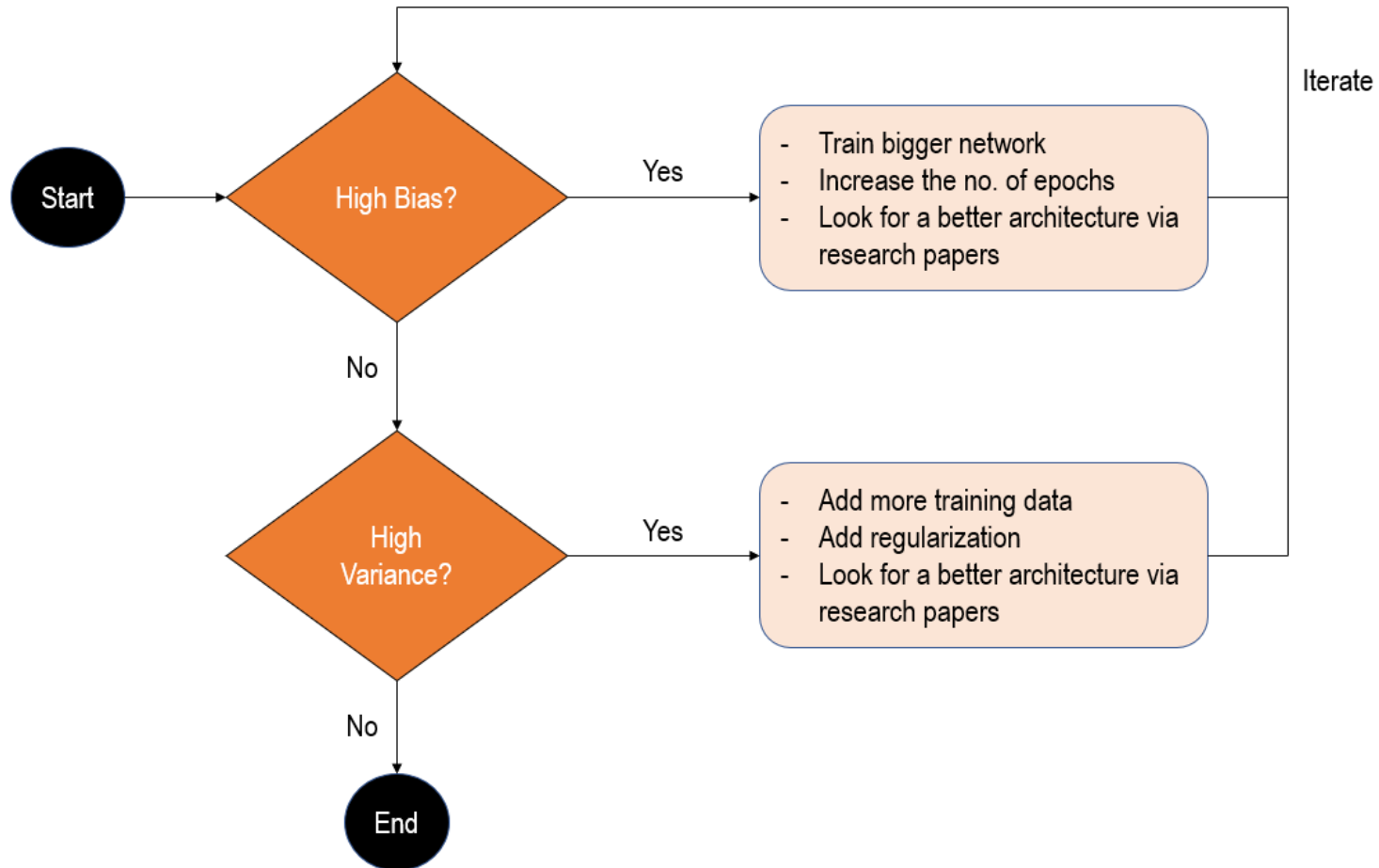


# WHY IS BIAS VARIANCE TRADEOFF?

- Tradeoff in complexity is why there is a tradeoff between bias and variance.



# GOOD TRADE-OFF: TIPS



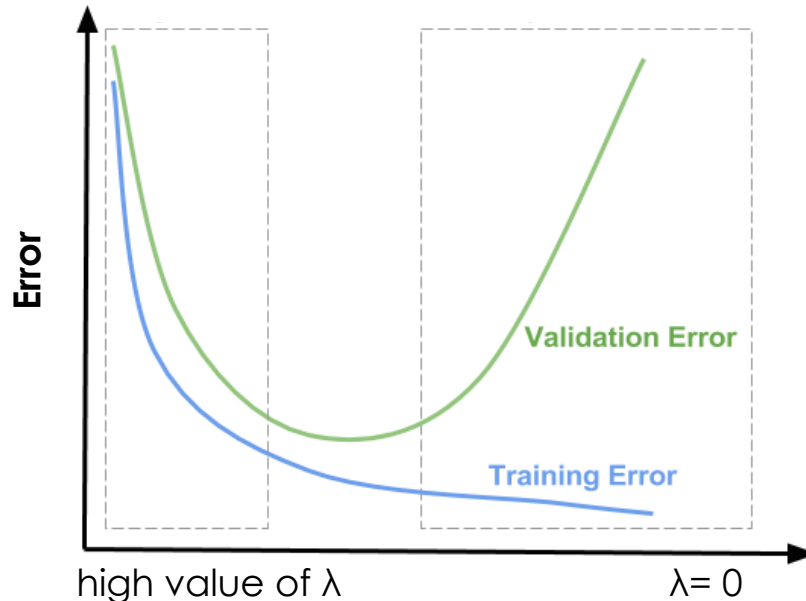
# READING

Yang, Z., Yu, Y., You, C., Steinhardt, J., & Ma, Y. (2020, November).  
**Rethinking bias-variance trade-off for generalization of neural networks.**  
In *International Conference on Machine Learning* (pp. 10767-10777).



# IF WE GO BACK TO THE REGULARIZATION

- The optimal value for  $\lambda$  will be probably somewhere around the minimum of the Cross Validation loss function



# LOSS FUNCTIONS

- The loss can be an arbitrary function mapping two vectors to a scalar
- **HINGE** (binary) **LOSS**:
  - the classifier's output is a single scalar  $\tilde{y}$  and the true label is in  $\{+1, -1\}$ ;
  - the classification rule is  $\hat{y} = \text{sign}(\tilde{y})$  and a prediction is correct if

$$y \cdot \tilde{y} \geq 0 \quad \rightarrow \quad \text{they share the same sign}$$

- The Hinge Loss is defined as:

$$L_{\text{hinge(binary)}}(\tilde{y}, y) = \max(0, 1 - y \cdot \tilde{y})$$





# LOSS FUNCTIONS

- **HINGE** (multi-class) **LOSS**:

- Let  $\hat{\mathbf{y}} = \hat{y}_{[1]}, \dots, \hat{y}_{[k]}$  be the classifier's output vector, and  $\mathbf{y}$  be the one-hot vector for the true label.
- The classification rule is defined as selecting the class with the highest score

$$\text{prediction} = \underset{i}{\operatorname{argmax}} \hat{y}_{[i]}$$

- The Hinge Loss is defined as:

$$L_{\text{hinge(multi-class)}}(\hat{\mathbf{y}}, \mathbf{y}) = \max(0, 1 - (\hat{y}_{[t]} - \hat{y}_{[j]}))$$

**TIPS:** In general HINGE LOSSES are useful when we require a hard decision rule, do not considering the class membership probability.



# LOSS FUNCTIONS

- LOGISTIC LOSS (BINARY CROSS ENTROPY):

- It is used in **binary classification** with **conditional probability outputs**
- We assume a set of two target classes labelled 0 and 1, with correct label  $y \in \{0, 1\}$
- The classifier's output  $\tilde{y}$  is transformed using the **sigmoid** (aka logistic) function  $\sigma(x) = 1/(1 + e^{-x})$  to the range  $[0, 1]$ 
  - It is interpreted as the conditional probability  $\hat{y} = \sigma(\tilde{y}) = P(y = 1|\mathbf{x})$
- The **prediction rule** is:

$$\text{prediction} = \begin{cases} 0 & \hat{y} < 0.5 \\ 1 & \hat{y} \geq 0.5 \end{cases}$$

- The training algorithm tries to maximize the log conditional probability  $P(y=1 | \mathbf{x})$  for each training example  $(\mathbf{x}, y)$ :

$$L_{\text{logistic}}(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$



# LOSS FUNCTIONS

## ○ NEGATIVE LOG LIKELIHOOD (CATEGORICAL CROSS-ENTROPY LOSS):

- It is used when a **probabilistic distribution** of scores is desired
- Let  $\mathbf{y} = y_{[1]}, \dots, y_{[k]}$  be a vector denoting the **true multinomial distribution** over the labels  $1, \dots, k$
- Let  $\hat{\mathbf{y}} = \hat{y}_{[1]}, \dots, \hat{y}_{[k]}$  be the linear classifier's predictions (transformed by the softmax function) that represent the **class membership conditional probability distribution**  $\hat{y}_{[i]} = P(y = i | x)$
- The negative log likelihood loss measures the dissimilarity between the true label distribution and the predicted one:

$$L_{\text{negative log likelihood}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_i y_{[i]} \log(\hat{y}_{[i]})$$

$$L_{\text{cross-entropy(hard classification)}}(\hat{\mathbf{y}}, \mathbf{y}) = -\log(\hat{y}_{[t]}),$$



# GRADIENT-BASED OPTIMIZATION

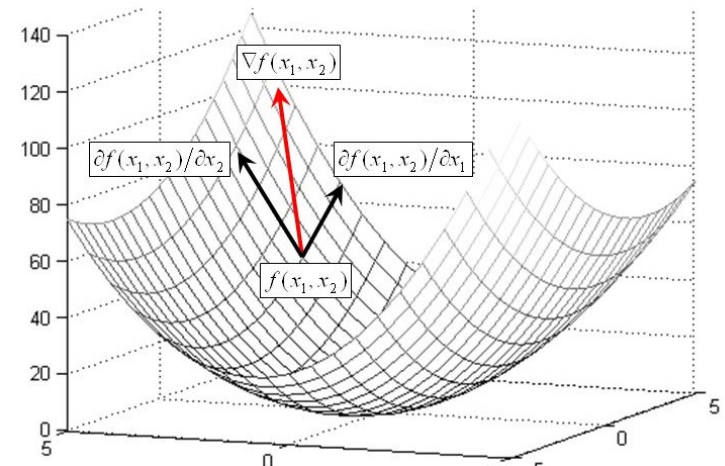
- In order to **train** the model, we need to solve the optimization problem:

$$\begin{aligned}\hat{\Theta} &= \underset{\Theta}{\operatorname{argmin}} \mathcal{L}(\Theta) + \lambda R(\Theta) \\ &= \underset{\Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n L(f(x_i; \Theta), y_i) + \lambda R(\Theta)\end{aligned}$$

- Solution: **gradient-based methods**

- Estimates the loss  $\mathcal{L}$  over the training set
- Computes the gradients  $\Theta$  w.r.t. the loss  $\mathcal{L}$
- Moves the parameters in the directions of the gradient

- When optimizing the loss  $\mathcal{L}$ , the parameters are considered as input, while the training examples are treated as constants



# STOCHASTIC GRADIENT DESCENT

- Stochastic Gradient Descent (SGD)
  - It is a general optimization algorithm
  - It receives a function  $f$  parameterized by  $\Theta$ , a loss function  $\mathcal{L}$ , input-output pairs  $x_{1:n}, y_{1:n}$
  - The main goal is to set the parameters  $\Theta$  such that the cumulative loss of  $f$  on the training examples is small

$$\mathcal{L}(\Theta) = \sum_{i=1}^n L(f(x_i; \theta), y_i)$$



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

It starts by sampling a training example



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

It computes the loss



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

It computes the gradient w.r.t  $\Theta$





# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

It computes the gradient w.r.t  $\Theta$

Input and predictions are assumed to be fixed and the loss is treated as a function of  $\Theta$



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$  ←
  - 6: **return**  $\Theta$
- 

Parameters  $\Theta$  are updated in the opposite direction of the gradient



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:   Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:   Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:    $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:    $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

Parameters  $\Theta$  are updated in the opposite direction of the gradient

scaled by a learning rate  $\eta_t$



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

Parameters  $\Theta$  are updated in the opposite direction of the gradient

scaled by a learning rate  $\eta_t$

**TIPS:** The learning rate can be either fixed or decay as function of the time step  $t$



# STOCHASTIC GRADIENT DESCENT

---

**Algorithm 2.1** Online stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

- 
- 1: **while** stopping criteria not met **do**
  - 2:     Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$
  - 3:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$
  - 4:      $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$
  - 5:      $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
  - 6: **return**  $\Theta$
- 

The error is based on a single training example

It is a rough estimation of the entire training set that we want to minimize

This may result in inaccurate gradients



# STOCHASTIC GRADIENT DESCENT

- A common way to reduce this noise is to estimate the error and the gradients based on a sample of  $m$  examples:
  - **Minibatch Stochastic Gradient Descent**

---

**Algorithm 2.2** Minibatch stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

---

```
1: while stopping criteria not met do
2:   Sample a minibatch of  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$ 
3:    $\hat{\mathbf{g}} \leftarrow 0$ 
4:   for  $i = 1$  to  $m$  do
5:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
6:      $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradients of } \frac{1}{m}L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i) \text{ w.r.t } \Theta$ 
7:    $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
8: return  $\Theta$ 
```

---

Estimate of the training loss based on the minibatch



# STOCHASTIC GRADIENT DESCENT

- A common way to reduce this noise is to estimate the error and the gradients based on a sample of  $m$  examples:
  - **Minibatch Stochastic Gradient Descent**

---

**Algorithm 2.2** Minibatch stochastic gradient descent training.

---

*Input:*

- Function  $f(\mathbf{x}; \Theta)$  parameterized with parameters  $\Theta$ .
- Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and desired outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .
- Loss function  $L$ .

---

```
1: while stopping criteria not met do
2:   Sample a minibatch of  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$ 
3:    $\hat{\mathbf{g}} \leftarrow 0$ 
4:   for  $i = 1$  to  $m$  do
5:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
6:      $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradients of } \frac{1}{m}L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i) \text{ w.r.t } \Theta$ 
7:    $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
8: return  $\Theta$ 
```

---

Final gradient estimate

Parameters are updated toward  
the gradient

# STOCHASTIC GRADIENT DESCENT

- Minibatch Stochastic Gradient Descent

- The minibatch can vary in size from  $m=1$  to  $m=n$ 
  - High values of  $m$  provide better estimates of the training gradients
  - Small values of  $m$  allow more updates and faster convergence
  - It provides also efficient training:
    - for small values of  $m$  some computing architectures (e.g. GPU) allow an efficient parallel implementation of the gradient computation

