

Deep Fool Phenomenon

Smart AI or Clever Hans

Machine Learning

When to apply Machine Learning?

Types of Machine Learning and which one did we focus on?

Learning Algorithm

Loss

Risk

A smart AI genuinely understands and solves problems, while a "Clever Hans" AI appears intelligent but relies on deceptive correlations or unintended cues rather than true learning

Adversarial attack method that iteratively perturbs an input to find the minimal perturbation needed to misclassify it in a deep neural network

- No explicitly known rules (speech recognition)
- Beyond human capabilities (100k rows of data)
- Adaptivity required (handwriting recognition)

A subset of AI that involves training algorithms on data to identify patterns and make predictions or decisions without explicit programming

$A: D^m \rightarrow Y^X$ → set of all (measurable) mappings $h: X \rightarrow Y$

Given a data sample $s \in D^m$, the learning algorithm outputs a learned hypothesis $h = A(s)$

- supervised vs unsupervised
- active vs passive
- helpful vs neutral vs adversarial
- online vs batch

Given a loss function l , an unknown data distribution μ on D we define the (expected) risk of a hypothesis $h: X \rightarrow Y$ by

$$R_\mu(h) := E_\mu[l(h(x,y))] = \int_{X \times Y} l(h(x,y)) \mu(dx dy)$$

To evaluate how well a hypothesis $h: X \rightarrow Y$ fits the given data (x_i, y_i) we use a loss function

$$l: Y^X \times D \rightarrow [0, \infty)$$

Empirical Risk

Empirical Risk Minimization
(ERM)

Biased ERM Rule

Approximation Error

Estimation Error

Idealized Error Decomposition

Optimization Error

Realistic Error Decomposition

Given training data $s = ((x_1, y_1), \dots, (x_m, y_m)) \in \mathcal{D}^m$,
 loss function $\ell: Y^X \times \mathcal{D} \rightarrow [0, \infty)$,

$$\text{ERM}(s) := \underset{h \in Y^X}{\operatorname{argmin}} \mathcal{R}_s(h) = \underset{h \in Y^X}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, y_i))$$

* Also known as **naive** ERM rule

Given a loss function ℓ , a data sample $s = ((x_1, y_1), \dots, (x_m, y_m)) \in \mathcal{D}^m$ the **empirical risk** of a hypothesis $h: X \rightarrow Y$ is given by

$$\mathcal{R}_s(h) := \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, y_i))$$

Given a loss function ℓ , data distribution μ , the **approximation error** of a hypothesis class H is given by

$$\mathcal{E}_{\text{app}}(H) := \min_{h \in H} \mathcal{R}_{\mu}(h)$$

Given training data $s = ((x_1, y_1), \dots, (x_m, y_m)) \in \mathcal{D}^m$,
 Hypothesis class $H \subseteq Y^X$, loss function $\ell: H \times \mathcal{D} \rightarrow [0, \infty)$

$$\text{ERM}_H(s) := \underset{h \in H}{\operatorname{argmin}} \mathcal{R}_s(h) = \underset{h \in H}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \ell(h, (x_i, y_i))$$

$$\mathcal{E}_{\text{gen}}(s, H) = \mathcal{E}_{\text{app}}(H) + \mathcal{E}_{\text{est}}(s)$$

Given a loss function ℓ , data distribution μ ,
 training data $s \in \mathcal{D}^m$, a learning algorithm A ,
 we define the **estimation error** by

$$\mathcal{E}_{\text{est}}(s) := \mathcal{R}_{\mu}(A(s)) - \min_{h \in H} \mathcal{R}_{\mu}(h)$$

$$\mathcal{E}_{\text{gen}}(s, H, k) = \mathcal{E}_{\text{app}}(H) + \mathcal{E}_{\text{est}}(s) + \mathcal{E}_{\text{opt}}(k)$$

We typically use **iterative optimization methods**
 for solving ERM. These recursively compute iterates $h_s^{(k)} \in H$, $k \in \mathbb{N}$, such that (hopefully) $h_s^{(k)} \rightarrow h_s$ as $k \rightarrow \infty$.
 However, we have to stop the iteration at some point, say after $k \in \mathbb{N}$ steps. This yields the **optimization error**

$$\mathcal{E}_{\text{opt}}(k) := \mathcal{R}_{\mu}(h_s^{(k)}) - \mathcal{R}_{\mu}(h_s)$$

Approximation vs Estimation
(Bias - Variance Trade off)

Linear Hypotheses

Halfspaces

Perceptron Algorithm

Convergence of Perceptron

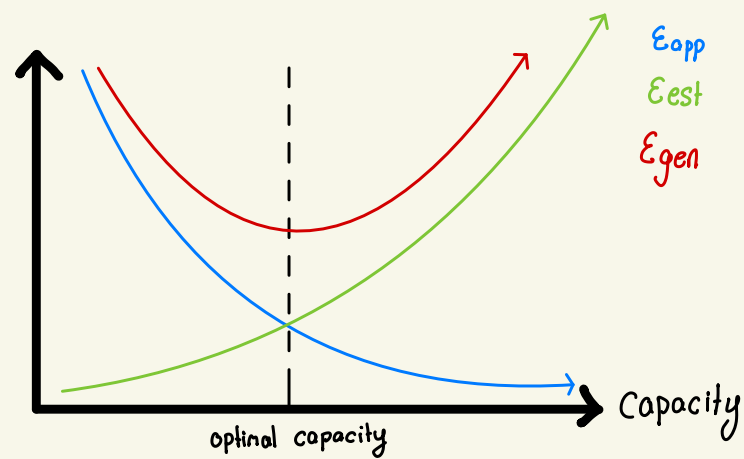
Likelihood

Maximum Likelihood Estimation

Log Loss

A hypothesis $h: X \rightarrow Y$, $X \subseteq \mathbb{R}^d$, is a **linear hypothesis** if there exists an affine-linear function $f_{w,b} \in L_d$ and an **activation function** $\phi: \mathbb{R} \rightarrow Y$ such that $h(x) = \phi \circ f_{w,b}(x) = \phi(w \cdot x + b)$, $x \in X$

The set of all linear hypotheses with activation function ϕ $\mathcal{L}_{d,\phi} := \{h_{w,b}(x) := \phi(wx+b) \mid w \in \mathbb{R}^d, b \in \mathbb{R}\}$



input: Sample s with m data points $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$

start: Set $t=0$ and $w_t' = (w_t, b_t) = 0$

iteration: For $t=1, 2, \dots$ do

 If there exists an $i=1, \dots, m$ with

$$y_i(w_t' \cdot x_i') \leq 0 \quad x_i' := (x_i, 1)$$

 then set

$$w_{t+1}' = w_t' + y_i x_i'$$

 else stop and output w_t'

A linear hypothesis $h_{w,b}(x) = \text{sgn}(wx+b)$, $x \in X$ divides the features space $X = \mathbb{R}^d$ in two halvespaces

$$H_{w,b}^+ := \{x \in \mathbb{R}^d : wx+b \geq 0\},$$

$$H_{w,b}^- := \{x \in \mathbb{R}^d : wx+b < 0\},$$

by a **separating hyperplane**

$$H_{w,b} := \{x \in \mathbb{R}^d : wx+b = 0\}$$

$$IP_{w,b}(Y=y \mid X=x) := \frac{1}{1 + e^{-y(wx+b)}}$$

The recursive rule $w_{t+1}' = w_t' + y_i x_i'$ for i with $y_i(w_t' \cdot x_i') \leq 0$ points the vector w_{t+1}' in the right direction, since

$$y_i(w_{t+1}' \cdot x_i') = y_i(w_t' \cdot x_i') + y_i^2(x_i' \cdot x_i') > y_i(w_t' \cdot x_i')$$

Algorithm stops after $T \leq R^2 B^2$ steps where $R := \max_{i=1, \dots, m} \|x_i'\|$

$$B := \min \{ \|w'\| : y_i(w' \cdot x_i') \geq 1 \quad \forall i=1, \dots, m \}$$

$$\text{Since } (w_s, b_s) = \arg\min_{w,b} \sum_{i=1}^m \ln(1 + e^{-y_i(wx+b)})$$

$$l(h_{w,b}, (x, y)) := \ln(1 + e^{-y(wx+b)})$$

$$\begin{aligned} L_{w,b}(s) &:= IP_{w,b}(Y_1=y_1, \dots, Y_m=y_m \mid X_1=x_1, \dots, X_m=x_m) \\ &= \prod_{i=1}^m IP_{w,b}(Y=y_i \mid X=x_i) = \prod_{i=1}^m \frac{1}{1 + e^{-y_i(wx+b)}} \end{aligned}$$

Maximum Likelihood Estimate

$$(w_s, b_s) = \arg\max_{w,b} L_{w,b}(s)$$

It should be clear that

$$\arg\max_{w,b} L_{w,b}(s) = \arg\min_{w,b} -\ln(L_{w,b}(s))$$

Multinomial Logistic Regression

Cross Entropy Loss

Consistency

PAC Condition
(Probably Approximately Correct)

PAC Learnability

PAC Learnability and
Consistency

No Free Lunch Theorem

Learnability via Uniform Convergence

$$l(h_{w,b}, (x,y)) = - \sum_{j=1}^n 1_{\{y_j\}}(y) \ln \left(\frac{e^{w_j \cdot x + b_j}}{\sum_{k=1}^n e^{w_k \cdot x + b_k}} \right)$$

$$h(x) \approx \begin{bmatrix} P(Y=y_1 | X=x) \\ \vdots \\ P(Y=y_n | X=x) \end{bmatrix} \in \mathbb{R}^n$$

output of the hypothesis h is a vector $p \in [0,1]^n$ such that $p_1 + \dots + p_n = 1$

Activation function is softmax $\sigma(z) = \left[\frac{e^{z_j}}{\sum_{k=1}^n e^{z_k}} \right]$

For a tolerance bound $\epsilon > 0$ and failure probability $\delta \in (0,1)$ it holds

$$\mathbb{P}_{\mu^m}(\mathcal{R}_{\mu}(h_S) \leq \inf_{h \in H} \mathcal{R}_{\mu}(h) + \epsilon) \geq 1 - \delta$$

or in short $\mathbb{P}_{\mu^m}(\mathcal{E}_{\text{est}}(H, S) \leq \epsilon) \geq 1 - \delta$

Given a hypothesis class $H \subseteq Y^X$, a loss l we call a learning algorithm $A: U_{\mu \in \mathcal{M}} D^m \rightarrow H$ (universally) consistent if for any distribution μ on D we have with $S \sim \mu^m$

$$\mathcal{E}_{\text{est}}(H, S) \xrightarrow[m \rightarrow \infty]{\mathbb{P}} 0 \iff \lim_{m \rightarrow \infty} \mathbb{P}_{\mu^m}(\mathcal{R}_{\mu}(A(S)) - \inf_{h \in H} \mathcal{R}_{\mu}(h) > \epsilon) = 0 \quad \forall \epsilon > 0$$

If a hypothesis class H is PAC learnable w.r.t. a given loss by a learning algorithm A , then this learning algorithm is consistent for H and l . Converse is not true, because the sample complexity m_H has to apply to any distribution μ

A hypothesis class $H \subseteq Y^X$ is called (agnostic) PAC learnable w.r.t. a given loss l , if there exists a mapping $m_H: (0,1)^2 \rightarrow \mathbb{N}$, a learning algorithm A , such that for any data μ on $X \times Y$ and any $\epsilon \in (0,1)$ and $\delta \in (0,1)$ we satisfy the PAC condition for $h_S = A(S)$ and $m \geq m_H(\epsilon, \delta)$. Smallest such mapping m_H is called sample complexity of H

Given $H \subseteq Y^X$ we have for $A = \text{ERM}_H$ almost surely

$$\mathcal{E}_{\text{est}}(H, S) = \mathcal{R}_{\mu}(A(S)) - \inf_{h \in H} \mathcal{R}_{\mu}(h) \leq 2 \sup_{h \in H} |\mathcal{R}_S(h) - \mathcal{R}_{\mu}(h)|$$

Let X be finite and $|Y|=2$. Further, let l be the 0-1 loss and A be an arbitrary learning algorithm. Then for any sample size $m < |X|/2$, there exists a distribution μ on $X \times Y$ and a hypothesis $h^*: X \rightarrow Y$ such that

$$\mathcal{R}_{\mu}(h^*) = 0 \quad \text{but} \quad \mathbb{P}_{\mu^m}(\mathcal{R}_{\mu}(A(S)) \geq 1/8) \geq 1/2$$

- There is no learning algorithm which succeeds on all tasks.
- We need to restrict to suitable hypotheses classes $H \subseteq Y^X$ for learnability

Uniform Convergence

Restriction

Growth Function

Shattering

VC Dimension

Fundamental Theorem of Learning

Margin

Hard SVM Rule

Given a hypothesis class $H \subseteq \{0,1\}^X$ and a finite set $M = \{x_1, \dots, x_m\} \subseteq X$ we define the **restriction of H to M** by

$$H_M := \{[h(x_1), \dots, h(x_m)] : h \in H\}$$

i.e., the set of all m -bits $b \in \{0,1\}^m$ generated by an $h \in H$ on M

A class $H \subseteq Y^X$ satisfies the **uniform convergence condition** (w.r.t. a loss ℓ) if there exists a **mapping** $m_H^{uc}: (0,1)^2 \rightarrow \mathbb{N}$ such that for any data distribution μ on $X \times Y$, any $\epsilon \in (0,1)$, $\delta \in (0,1)$ we have

$$\mathbb{P}_{\mu^m} \left(\sup_{h \in H} |R_{\mu}(h) - R_S(h)| \leq \epsilon \right) \geq 1 - \delta \quad \forall m \geq m_H^{uc}(\epsilon, \delta)$$

• H satisfies **UC** $\Rightarrow H$ is PAC-learnable with $A = \text{ERM}_H$

Let $M \subseteq X$ be finite and $H \subseteq Y^X$, $|Y|=2$, be a class of **binary hypotheses**. Then we say **H shatters the set M** if its restriction H_M to M satisfies

$$H_M = Y^{|M|} \iff |H_M| = 2^{|M|}$$

For a binary hypothesis class $H \subseteq \{0,1\}^X$ its **growth function** $T_H: \mathbb{N} \rightarrow \mathbb{N}$ is given by

$$T_H(m) := \sup_{M \subseteq X: |M|=m} |H_M|$$

For a class $H \subseteq Y^X$, $|Y|=2$, the following statements are equivalent given the **0-1 loss**:

- 1) H satisfies uniform convergence (UC).
- 2) H is (agnostic) PAC-learnable by $A = \text{ERM}_H$.
- 3) H is (agnostic) PAC-learnable
- 4) H has finite VC dimension.

The **VC (Vapnik-Chervonenkis) dimension** of a hypothesis class $H \subseteq \{0,1\}^X$ is

$$\text{VCD}(H) := \sup \{ |M| : H \text{ shatters } M \subseteq X \}$$

Given **linearly separable sample** s with m pairs of data $(x_i, y_i) \in \mathbb{R}^d \times \{-1, +1\}$, compute $h_{w,b}(s) \in \mathbb{L}_d$ given by

$$(w_s, b_s) = \arg \min_{(w,b) \in \mathbb{R}^{d+1}} \|w\|^2 \quad \text{subject to} \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

The **margin** of a sample $s = ((x_1, y_1), \dots, (x_m, y_m))$ to a **hyperplane** $H_{w,b} := \{x \in \mathbb{R}^d : w \cdot x + b = 0\}$ is the smallest distance of a point $x_i \in \mathbb{R}^d$ to $H_{w,b}$.

$$\gamma_{w,b}(s) := \frac{1}{\|w\|} \min_{i=1, \dots, m} |w \cdot x_i + b|$$

Why is it called
"Support Vector" Machine?

Karush-Kuhn-Tucker Conditions

Soft SVM Rule

Hinge Loss

Kernel SVM

Hard SVM Rule in F

Soft SVM Rule in F

Representer Theorem

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be **differentiable**, $g_i(w) = a_i^T w + c_i$, $a_i \in \mathbb{R}^d$, $c_i \in \mathbb{R}$ for $i=1, \dots, m$ and consider

$$w^* \in \operatorname{argmin}_{w \in \mathbb{R}^d} f(w) \quad \text{subject to: } g_i(w) \leq 0 \quad \forall i=1, \dots, m$$

Then there exists coefficients $\alpha_i \geq 0$, $i=1, \dots, m$, such that
 $\nabla f(w^*) + \sum_{i=1}^m \alpha_i \nabla g_i(w^*) = 0$ and $\alpha_i g_i(w^*) = 0 \quad \forall i=1, \dots, m$

Consider now the SVM rule as a special case of the optimization task:
 $f(w, b) = \|w\|^2$, $\Rightarrow \nabla f(w, b) = (2w, 0)$
 $g_i(w, b) = 1 - y_i(\langle w, x_i \rangle + b) \Rightarrow \nabla g_i(w, b) = -y_i(x_i, 1)^T$ } KKT yields $w_s = \sum_{j \in J} \alpha_j x_j$

The weight vector $w_s \in \mathbb{R}^d$ learned by the hard SVM rule is composed of special data points $x_j \in \mathbb{R}^d$:

$$w_s = \sum_{j \in J} \alpha_j x_j \quad J \subseteq \{1, \dots, m\} := \{i: y_i(w_s \cdot x_i + b_s) = 1\}, \alpha_j \in \mathbb{R}$$

the vectors x_j with $y_j(w_s \cdot x_j + b_s) = 1$ are called **support vectors** of w_s .

The support vectors are exactly those data points x_j which have the smallest distance to the hyperplane H_{w_s, b_s} :

$$y_i(w_s \cdot x_i + b_s) = 1 \iff d(x_i, H_{w_s, b_s}) = \gamma_{w_s, b_s}(s)$$

For $w' = (w, b) \in \mathbb{R}^{d+1}$ and $x \in \mathbb{R}^d$ and $y \in \{-1, +1\}$ let

$$l_{\text{hinge}}(w', (x, y)) := \max\{0, 1 - y(w' \cdot x')\}, \quad x' = (x, 1)$$

Given a sample S with m data pairs $(x_i, y_i) \in \mathbb{R}^d \times \{-1, +1\}$, parameter $\lambda > 0$

Compute $h_{w_s, b_s} = \text{SYM}_{\text{soft}}(S; \lambda) \in \mathcal{L}_d$ given by

$$(w_s, b_s, \xi_s) \in \operatorname{argmin}_{(w, b, \xi) \in \mathbb{R}^{d+1+m}} \lambda \|w\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i$$

subject to: $y_i(w \cdot x_i + b) \geq 1 - \xi_i$ and $\xi_i \geq 0 \quad \forall i=1, \dots, m$

Compute $h_s(x) := \operatorname{sgn}(w_s \cdot \psi(x) + b_s)$ by

$$(w_s, b_s) = \operatorname{argmin}_{w \in F, b \in \mathbb{R}} \|w\|^2$$

subject to

$$y_i(w \cdot \psi(x_i) + b) \geq 1 \quad \forall i$$

The idea is to map the original inputs $x \in X$ into a (much) larger **feature space** F , and apply linear hypotheses $h_F: F \rightarrow \{-1, +1\}$ in F for classification

The corresponding embedding $\psi: X \rightarrow F$ is called **feature map** and the resulting **nonlinear hypotheses** on X are then $h_F \circ \psi$.

The outcome $w_s \in F$ of the hard and soft SVM rule in F as well as any

$$(w_s, b_s) \in \operatorname{argmin}_{w \in F, b \in \mathbb{R}} f(w \cdot \psi(x_1) + b, \dots, w \cdot \psi(x_m) + b) + R(\|w\|)$$

for arbitrary $f: \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ and **strictly increasing** $R: [0, \infty) \rightarrow \mathbb{R}$ can be represented as

$$w_s = \sum_{i=1}^m \alpha_{s,i} \psi(x_i), \quad \alpha_{s,i} \in \mathbb{R}$$

Compute $h_s(x) := \operatorname{sgn}(w_s \cdot \psi(x) + b_s)$ by

$$(w_s, b_s) = \operatorname{argmin}_{w \in F, b \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i[w \cdot \psi(x_i) + b]\} + \lambda \|w\|^2$$

Kernel Trick

Hard Kernel SVM Rule

Soft Kernel SVM Rule

Polynomial Kernel

Gaussian Kernel

Train, Validation, Test

Cross Validation

Gradient Descent

Compute $h_s(x) := \text{sgn}(\sum_{i=1}^m \alpha_{s,i} k(x_i, x) + b_s)$ by

$$(\alpha_s, b_s) = \underset{(\alpha, b) \in \mathbb{R}^{n+1}}{\text{argmin}} \alpha^T K \alpha$$

subject to

$$y_i (k_i \alpha + b) \geq 1 \quad \forall i$$

Given a feature map $\psi: X \rightarrow F$ we define a **kernel** $K: X \times X \rightarrow \mathbb{R}$ by

$$k(x, y) := \psi(x) \cdot \psi(y)$$

Thus we can express everything by K , e.g.,

$$w_s \cdot \psi(x_i) = \sum_{j=1}^m \alpha_{s,j} k(x_i, x_j), \quad h_s(x) = \text{sgn}(\sum_{i=1}^m \alpha_{s,i} k(x_i, x) + b_s)$$

Let $X = \mathbb{R}^d$ and for a $q \in \mathbb{N}$ set $k(x, y) := (1 + x \cdot y)^q$

As a feature mapping $\psi: \mathbb{R}^d \rightarrow \mathbb{R}^{(1+d)^q}$ we then define

$$\psi(x) := \left(\prod_{i=1}^q x_{j_i} : j = (j_1, \dots, j_q) \in \{0, \dots, d\}^q \right)$$

It then holds with Euclidean inner product in \mathbb{R}^d of $F = \mathbb{R}^{(1+d)^q}$

$$k(x, y) = (1 + x \cdot y)^q = \sum_{j \in \{0, \dots, d\}^q} \prod_{i=1}^q x_{j_i} y_{j_i} = \psi(x) \cdot \psi(y)$$

Compute $h_s(x) := \text{sgn}(\sum_{i=1}^m \alpha_{s,i} k(x_i, x) + b_s)$ by

$$(\alpha_s, b_s) \in \underset{(\alpha, b) \in \mathbb{R}^{n+1}}{\text{argmin}} \lambda \alpha^T K \alpha + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y_i (k_i \alpha + b)\}$$

Training: Determine hypothesis h_s via numerical computation of a learning rule $h_s \approx A(s)$, e.g. ERM rule $A = \text{ERM}_H$

Validation: Choose the most appropriate class H or hyperparameter

Test: Estimate the generalization error of the learned hypothesis

$$k(x, y) := e^{-\gamma \|x - y\|^2} \text{ with scaling parameter } \gamma > 0$$

Gaussian kernel evaluates the similarity of the features with respect to their distance $\|x - y\|$, where the distance enters nonlinearly

The Gaussian kernel is also called the **RBF kernel**, where RBF stands for **radial basis function**

To small γ may not allow a good fit and too large γ will lead to overfitting

Given a starting vector $w_0 \in \mathbb{R}^p$, calculate for $k = 0, 1, 2, \dots$

$$w_{k+1} = w_k - \eta_k \nabla F(w_k)$$

with corresponding **step sizes** $\eta_k > 0$

1. Divide the dataset into k equal-sized subsets (folds)

2. For each fold:

- Use $k-1$ folds for training
- Use the remaining fold for validation

3. Rotate the validation fold and repeat the process k times

4. Compute the average performance metric

Lipschitz - smooth

λ -strongly convex

Stochastic Gradient Descent

Subgradient
Descent

Artificial Neural Networks

Characteristics of FNN

Structure of CNN

ϵ -table

A differentiable function $F: \mathbb{R}^p \rightarrow \mathbb{R}$ is called λ -strongly convex if for $\lambda > 0$ we have

$$F(v) \geq F(w) + \nabla F(w)^T (v-w) + \frac{\lambda}{2} \|v-w\|^2 \quad \forall w, v \in W$$

Strongly convex functions possess at most one minimum

A differentiable function $F: \mathbb{R}^p \rightarrow \mathbb{R}$ is called L -smooth if for $L > 0$ we have

$$\|\nabla F(w) - \nabla F(v)\| \leq L \|w-v\| \quad \forall v, w \in W$$

L -smoothness guarantees a decrease in the objective function value for gradient descent

Given initial state $w_0 \in \mathbb{R}^p$ compute for $k=0,1,2,\dots$

$$w_{k+1} = w_k - \eta_k v_k,$$

$$v_k \in \partial F(w_k)$$

with suitable stepsizes $\eta_k > 0$

Given a starting vector $w_0 \in \mathbb{R}^p$ and an objective function $F(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$, $f_i: \mathbb{R}^p \rightarrow \mathbb{R}$,

calculate for $k=1,2,\dots$ iterates w_{k+1} as follows

- draw realization $i_k \in [m]$ of the uniformly random index variable $I_k \sim U([m])$, $[m] := \{1, \dots, m\}$ where the $I_k, k \in \mathbb{N}$ are stochastically independent,
- for a given deterministic step size $\eta_k > 0$, calculate

$$w_{k+1} = w_k - \eta_k \nabla f_{i_k}(w_k)$$

$$h(x) = p \circ f_{w_L, b_L} \circ \sigma \circ f_{w_{L-1}, b_{L-1}} \circ \sigma \circ \dots \circ \sigma \circ f_{w_1, b_1}(x)$$

Depth: L

Width: $B := \max_{k=0, \dots, L} n_k$

Size: $n := n_0 + n_1 + \dots + n_L$

Architecture: (V, E) with $V = (V_0, \dots, V_L)$ and $E \subseteq \{(v_{k,i}, v_{k+1,j}) : v_{k,i} \in V_k \text{ and } v_{k+1,j} \in V_{k+1}\}$

Number of parameters: $P_{VE} := |E| + |V_1| + \dots + |V_L| \leq L(B^2 + B)$

A feedforward neural network is a hypothesis $h: X \rightarrow Y$ of the form

$$h(x) = p \circ f_{w_L, b_L} \circ \phi \circ f_{w_{L-1}, b_{L-1}} \circ \phi \circ \dots \circ \phi \circ f_{w_1, b_1}(x)$$

where

- $\phi: \mathbb{R} \rightarrow \mathbb{R}$ and $p: \mathbb{R} \rightarrow Y$ are chosen activation functions whose applications are to be understood componentwise,
- given layerwise weight matrices $w_k \in \mathbb{R}^{n_k \times n_{k-1}}$ and bias vectors $b_k \in \mathbb{R}^{n_k}$ $f_{w_k, b_k}(y) := w_k y + b_k$
- with $n_k \in \mathbb{N}$ denotes the size of the k -th layer V_k

Method	ϵ_{app}	ϵ_{est}	ϵ_{opt}
Perceptron			
Logistic Regression			
Hard/Soft SVM rule			
Kernel SVM rule			
FNN			

CNN typically consist of three kinds of layers

- Convolutional layer: extract new features by applying convolutional filters
- Pooling layer: Reduces information and downsamples new features
- Fully connected layer: Learns classification based on new features provided by previous layer

The Universal Approximation theorem or Cybenko theorem states that if we do have a NN specifically a Shallow FNN which has $L=2$, then it learn any kind of pattern or function mostly continuous if we could give it enough neurons and a finite width.

here we do have a continuous input activation function which is a sigmoid activation function. in the Shallow FNN formula we tend to set our learning NN to learn a continuous other activation function with respect to a threshold.

This threshold mean that the Maximum differences of the Shallow FNN and the learned continuous function should not very low and should not be greater than the threshold. this means that the learned continuous function or pattern should be as close as possible to the main Shallow FNN.

We could also tell that bound of differences, that can be the Approximation error too for any $\epsilon > 0$

The Kigler and Lyons theorem says that even if we could have a deep narrow NN, unlike the Cybenko Theorem, we could have a smooth term of the NN too that can approximate the continuous g function. again, the same as above it says that learned continuous pattern should be as close as possible to main NN that have learned. even with having the small fixed with : $d+3$ neurons per Layer!



wide fNN $\rightarrow VC = \infty$



deep narrow

so Going deeper is better than going wider, because it gives: Faster improvement in accuracy - Lower VC dimension not to tend to infinity or even Estimation Error. according to Petersen theorem.

This is a lower bound on how close the neural network's output $h(x)$ can get to the target $g(x)$.

The error shrinks (gets better) as:

Depth L increases \rightarrow very fast (exponentially)

Width B increases \rightarrow much slower (polynomially)

$$\sup_{x \in [0,1]} |g(x) - f(x)| \geq Cg(2B+2)^{-2L-2}$$

$$\textcircled{1} + \textcircled{2} = \sup_{x \in X} |g(x) - h(x)| \leq \epsilon$$

$$\rightarrow \text{NN} : y_{Ki} = \phi \left(\sum_{l=1}^m w_{l,i} y_{K-1,l} + b \right)$$

w.r.t $\phi, \rho, \beta, \alpha$:

$$y_{Ki} = \rho \left(\beta + \sum_{l=1}^m \alpha_l \phi(w_l^T x + b_l) \right)$$