

# Connections

**1.1 Background**  
Connections aims to build intelligent machines inspired by systems of interconnected neurons.

- McCulloch & Pitts Neuron - 1940:**
  - A neuron with  $n$  inputs can be characterized as a threshold element with Boolean signature  $\sigma \in \{-1, 1\}^n$  and threshold  $\theta$ .
  - A neuron computes its action or output as a Boolean function of its input  $x \in \{0, 1\}^n$  according to  $output = \begin{cases} 1 & \text{if } \sum_{i=1}^n \sigma_i x_i \geq \theta \\ 0 & \text{otherwise} \end{cases}$

**Issue:** The Boolean logic model of neural networks fails to capture biological systems' adaptive and learning capabilities.

**Solution:** Subsequent contributions (the Thesaurus maze-navigating machine and Minsky's SNARC), demonstrated early concepts of learning and adaptability but achieved limited practical impact.

**1.2 Perceptron, 1958 and 1969**  
The perceptron marked a milestone in machine learning by introducing supervised learning for binary classification.

**Setting:** Feature vector  $x \in \mathbb{R}^n$ , class membership  $y \in \{-1, +1\}$ , sample of  $s$  training examples  $S = \{(x_i, y_i) : 1 \leq i \leq s\} \subset \mathbb{R}^n \times \{\pm 1\}$

**Goal:** Learn discriminant function  $f$  that classifies examples in  $S$   
**Model:** Linear threshold unit with synaptic weights  $w \in \mathbb{R}^n$  and threshold  $\theta \in \mathbb{R}$ . Formally,  $f[w, b](x) = \text{sign}(x \cdot w + b)$ , where  $x \cdot w = \sum_{i=1}^n x_i w_i$ . Here,  $\text{sign}(0) = 0$

**Decision Boundary:** The decision boundary is the set of all  $x$  for which  $\sum_{i=1}^n w_i x_i + b = 0$ . Note that this is the Hesse normal form of a hyperplane in  $\mathbb{R}^n$  with normal vector  $n = \frac{w}{\|w\|}$ .

**Distances:**

- Signed distance from any point  $x_0$  to hyperplane:  $\frac{x_0 \cdot w + b}{\|w\|}$

**Geometric Margin:**

- Definition:**  $\gamma[w, b](x, y) = \frac{y(x \cdot w + b)}{\|w\|}$  (the sign of  $\gamma$  encodes the correctness of the classification + correct, - wrong)
- $\gamma$ -Separation:** A classifier  $f[w, b]$  is said to  $\gamma$ -separate the sample  $S$ , if  $\forall (x, y) \in S : \gamma[w, b](x, y) \geq \gamma$
- Maximum Margin (MM) Classifier:**  $f[w^*, b^*]$  is a classifier with maximal separation margin:  $(w^*, b^*) \in \text{argmax}_{(w, b)} \gamma[w, b](S)$ , where we define  $\gamma[w, b](S) = \min_{(x, y) \in S} \gamma[w, b](x, y)$
- Uniqueness of MM Classifier:** The MM classifier is unique, but its parameters are only unique up to positive scaling

**Version Space:** The version space is the set of all parameters of a classifier  $f[\theta]$  that are compatible with the training sample (predict all labels correctly):  $V(S) = \{\theta \mid \forall (x, y) \in S : f[\theta](x) = y\}$ .

**Linear Separability:** The version space of a linear threshold classifier is  $V(S) = \{(w, b) : \gamma[w, b](S) > 0\} \subset \mathbb{R}^{n+1}$ . We call  $S$  linearly separable if and only if  $V(S) \neq \emptyset$  ( $\Rightarrow$  perceptron learning aims to find a perfect linear separator, and does not explicitly aim (and is not guaranteed) to find classification with low training error).

- Perceptron Learning:**
- Present training examples repeatedly in some arbitrary order, making sure every example is repeated after finite number of steps.
  - Follow the update rule:  $f[w, b](x) \neq y \Rightarrow w \leftarrow w + yx, b \leftarrow b + y$ . Note how the perceptron only learns from mistakes!

**Perceptron weights:**  $w^0 \in \text{span}(x_1, \dots, x_n) \Rightarrow w^0 \in \text{span}(x_1, \dots, x_n), \forall t$

**Convergence:** General setting: Set  $b = 0$ , initialize  $w^0 > 0$ .

- Lemma** Let  $w, \|w\| = 1$ , such that  $\gamma[w](S) = \gamma > 0$ . Then  $w^t \cdot w \geq t\gamma$ .
- Lemma** Let  $R = \max_{x \in S} \|x\|$ , then  $\|w^t\| \leq R/\gamma t$
- Theorem** Let  $S$  be  $\gamma$ -separable,  $R = \max_{x \in S} \|x\|$ . The perceptron converges in at most  $\lceil R^2/\gamma^2 \rceil$  steps.
- Exercise** Let  $(x^i, y^i)$  be sequence of mistakes, inducing updates  $\Delta \theta^i, \delta^i = \sum_{i=1}^n \Delta \theta^i$ . Then  $\|\theta^0\|^2 \leq \sum_{i=1}^n \|\delta^i\|^2$ .

**Combinatorics and Geometry:** General setting: Set  $b = 0$ .

- Number of ways to classify  $s$  points in  $n$ -space with linear classifier:**  $C(s, n) = |\{y \in \{-1, 1\}^s \mid \exists w \in \mathbb{R}^n : \forall (x, y) \in S : (x \cdot w) > 0\}|$
- General position:** Let  $\{x_1, \dots, x_k\} \subset \mathbb{R}^n$  are in general position if and only if  $\forall S \subseteq \mathcal{X}$  with  $|\mathcal{S}| \leq n$  is linearly independent.
- Cover's Theorem:** Let  $S \subset \mathbb{R}^n$  be a set of  $s$  points in general position. Then  $C(s+1, n) = 2 \sum_{i=0}^n \binom{s}{i}$
- Corollary:**  $C(s, n) = 2^n$  for  $s \leq n$ .

$$\text{Asymptotic Shattering: } \frac{C(n, n)}{2^n} = \begin{cases} 1 & \text{if } n \leq n, \\ 1 - O(e^{-n}) & \text{if } n < 2n, \\ e^{-2n} & \text{if } e = 2n, \\ \tilde{O}(e^{-n}) & \text{otherwise.} \end{cases}$$

**Interpretation:** There two growth regimes. One where the number of realizable labelings grows exponentially, and one where the growth function of realizable labelings starts go get constrained by geometry.

**1.3 Parallel Distributed Processing**  
**Layers and Activation Functions:** The layer transfer function is  $f(W, B)(x) = \phi(Wx + B)$ .  $W = (w_1, \dots, w_m)^T$ ,  $\Phi(z) = (\phi(z_1), \dots, \phi(z_m))^T$ . Note how the linear threshold unit is a special case with  $\phi = \text{sign}$ .

- Delta Rule:**
- The original perceptron update rule has disadvantages:
    - 0/1 classification error not differentiable at decision boundary and has zero derivatives elsewhere.
    - We do not learn anything from correctly classified points.
  - The delta learning rule can be formulated as follows:  
 $\Delta w_{ij} = \eta \delta_j \phi'(w_i \cdot x_j)$ ,  $\delta_i = (y_i - f_i)$ ,  $\eta > 0$

**Remarks:**

- Update of a weight vec. for a unit is always in dir. of a training point
- The update is scaled by the residual ( $\delta$  or delta)  $\delta_i$  of the unit
- The sensitivity of the unit factors in (derivative of activation)
- $\eta$  is a design choice
- The delta rule is nothing but a gradient step with step size  $\eta$ , i.e.,  $\Delta w_{ij} = -\eta \frac{\partial \text{loss}}{\partial w_{ij}}$  with squared loss  $L = \frac{1}{2} \sum_{i=1}^n (y_i - y_i^t)^2$
- Generalized delta rule for deep architectures:  $\delta_i = \frac{\partial \text{loss}}{\partial y_i}$ , derived by propagating error terms from the outputs back towards the inputs

**Multi-Layer Perceptron:** Particular neural network architecture that can be written as (for  $n$  inputs,  $m$  hidden units, and one real output):  
 $f(w, W)(x) = \sum_{i=1}^m w_i \phi_i(x), \phi_i = \phi(w_i \cdot x), \phi(z) = \frac{1}{1 + e^{-z}}$ .

**1.4 Hopfield Networks: Classical model of an associate memory**  
**Idea:** Define a parameterized energy function via second order interaction of  $n$  binary neurons:

$E(x) = -\frac{1}{2} \sum_{i,j} w_{ij} x_i x_j + \sum_i b_i x_i, x_i \in \{-1, +1\}, 1 \leq i \leq n$ .  
Where we constrain the weights as follows:  $w_{ij} = w_{ji}$  ( $\forall i, j$ ),  $w_{ii} = 0$  ( $\forall i$ ). Think of the weights as quantifying the interaction strength between neurons, and of the biases as thresholds (set to zero in the following).  
**Dynamics:** After initialization, the state of each neuron is iteratively (in some fixed order) updated as follows:  
 $x_i \leftarrow \begin{cases} +1 & \text{if } E(\dots, x_{i-1}, -1, x_{i+1}, \dots) \leq E(\dots, x_{i-1}, +1, x_{i+1}, \dots) \\ -1 & \text{else} \end{cases}$ ,

where updates can be performed synchronous or asynchronous. Note that asynchronous updates guarantee that the energy function will never increase. Note after a finite number of steps, a fixed point is reached, which is the pattern of the Hopfield network.

**Simplified update Rule:** Instead of evaluating the  $(n+1)/2$  summands in  $E$ , we can just calculate  $n$  terms:  $H_i = \sum_j w_{ij} x_j$  and follow the update rule  $x_i = \text{sign}(H_i)$ . Here  $\text{sign}(0) = 1$ .

**Hebbian Learning:** Given a set  $S$  with patterns  $x^i \in \{\pm 1\}^n$ ,  $(1 \leq i \leq s)$ . Then choose the weights as  $w_{ij} = \frac{1}{s} \sum_{i=1}^s x_i^j x_i^j$ , which in matrix notation can be written as  $W = \frac{1}{s} \sum_{i=1}^s x^i x^{iT}$  ( $\Rightarrow$  neurons that fire together wire together)

**Memorization:** For a pattern to be included in associative memory, the pattern must define a meta-stable state:  $x^i = \text{sign}(\sum_{j=1}^n w_{ij} x_j^i)$  (the dynamics will not diverge from the pattern).

**Cross-talk term between patterns:** Using the couplings learned with Hebbian learning, we get:  $\sum_{j=1}^n w_{ij} x_j^i x_j^j = x^i \cdot x + \frac{1}{s} \sum_{j=1}^n \sum_{k=1}^s x_{i,j} x_{k,j} x_j^j$ , where  $C_j^i = \frac{1}{s} \sum_{k=1}^s x_{i,j} x_{k,j} x_j^j$  is called the cross-talk between patterns. For  $x^i$  to be stable, we require this term to be strictly less than 1 in magnitude for all neurons.

**Capacity:** Assume patterns to be random signs (probability of  $\frac{1}{2}$  for  $-1$  and  $+1$ ). Assume that  $n \rightarrow \infty$ . Then  $C_1^i \approx^{\text{approx.}} N'(0, \frac{n}{s})$ . Then the error probability of any sign to be incorrectly flipped becomes:  $P_e = \text{Prob}(\sum_j C_j^i C_j^i) \approx \int_0^\infty e^{-\frac{u^2}{2}} dz = \frac{1}{2} [1 - \text{erf}(\sqrt{n/2s})]$  ( $\Rightarrow$  ratio  $s/n$  controls asympt. err. rate).  
 $\Rightarrow$  It can be shown that a phase transition occurs at  $s/n \approx 0.138$ . Beyond that limit many errors occur. If one requires that a pattern is retrieved with high probability, then one gets a sub-linear capacity bound of  $s \leq \frac{n}{2\log_2 n}$ .

**Spurious patterns:** Simplest spurious pattern: sign inverted version of stable pattern is also stable as the network is sign symmetric.

## 2. Linear Networks

### 2.1 Regression Models

- Least Squares** Assume:  $f[w](x) = w^T x$
- Criterion (MSE):  $h[w] = \frac{1}{2s} \sum_{i=1}^s (w^T x_i - y_i)^2$
  - Matrix form  $h[w](S) = \frac{1}{2s} \|Xw - y\|^2$ , Normal eq.:  $\forall h = 2X^T X - 2X^T y$
  - $w^* = X^+ y \in \text{argmin}_{w \in \mathbb{R}^n} h[w]$  where  $X^+$  is the Moore-Penrose inverse of  $X \Rightarrow \lim_{\delta \rightarrow 0} (X^T X + \delta I)^{-1} X^T$
  - For very large datasets, SGD can be used to train:

- $w_{i+1} = w_i + \eta(y_i - w_i^T x_i) x_i$ , where  $x_i \stackrel{iid}{\sim} \text{Uniform}(1, \dots, s)$
- LS model as neg-log-likelihood function corresponds to Gaussian homoscedastic noise model:  $y_i = w^T x_i + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$
- LS assumes iid Gaussian noise on targets
- Regularized:  $h_1[w] = h_{\text{old}}[w] + \frac{\lambda}{2} \|w\|^2$  s.t.  $w^* = (X^T X + \lambda I)^{-1} X^T y$
- $\lambda \rightarrow 0$  LS and  $\lambda = \text{Lasso}$ ,  $\lambda = \text{Ridge}$

**Logistic function** is  $f(z) = \frac{1}{1 + e^{-z}}$ ,  $\sigma(z) = \sigma(-z) = 1$ , as  $\sigma'(-z) = \frac{e^{-z}}{1 + e^{-z}}^2$ .

It means  $\sigma = \frac{1}{2}$  is an odd function. Derivatives are polynomials in  $\sigma$ .  
 $\sigma' = \sigma(1 - \sigma) = \sigma(1 - \sigma)(1 - 2\sigma)$   
 $\sigma'$  is monotonically increasing, convex in  $\mathbb{R}$  and concave in  $\mathbb{R}$ .  
Units  $1$  and  $\sigma$  are often referred to as probabilities of Bernoulli event  
(Negative) logistic log likelihood with respect to binary outcome is equivalent to cross entropy with respect to degenerate target distribution:  
 $\ell(y, z) = -y \log(\sigma(z)) - (1 - y) \log(1 - \sigma(z))$   
 $\ln(\sigma(z)) = -\ln(1 + e^{-z})$ ,  $\ln(\sigma(z))' = 1 - \sigma$ ,  $\ln(\sigma(z))'' = -\sigma(1 - \sigma) < 0$   
Implying:  $(1 - \log \sigma)$  strictly convex  $(2)$  upper-bounds the heavy-side function  $(0-1)$  error

**Logistic Regression**

- Minimizes logistic loss over  $S$ :  $h[w] = \frac{1}{2s} \sum_{i=1}^s \ell(y_i, w^T x_i)$
- No closed form solution, but can be optimized with SGD
- Iterates remain in the span of  $\{w_0, x_1, \dots, x_n\}$ , and  $\forall i, V_i = \sigma(w^T x_i) - y_i |x_i$ . Note: predicted probability minus 0/1 target appears in the gradient.

### 2.2 Layers & Units

**Linear  $F$  is a parametrized map (affine transformation + nonlinearity)**  
Width  $M$ :  $F: \mathbb{R}^{n(n+1)} \times \mathbb{R}^m \rightarrow \mathbb{R}^m$   
 $\text{parameter}$   $\text{input}$   $\text{output}$   
 $F[\theta] = \text{vec}(W, b)(x) \equiv \phi(Wx + b)$

- Composition**  $F = F^1[\theta^1] \circ \dots \circ F^L[\theta^L]$ ,  $F^i[\theta^i] = \phi^i(W^i x + b^i)$
- Computing  $F^i = F^i(x^{i-1})$  is dominated by matrix-vector product
- Units:** component functions of  $F$  parametrized by  $w$  and  $b$  defining family of functions  $f[w, b](x) = \phi(w^T x + b)$
- Invar. of units:**  $f(x) = f(x + \delta x) \forall \delta x \perp w$ . Hence, level sets of a unit are orthogonal to  $w$ . For  $f_L(x) = \{x : \phi(w^T x + b) = z\}$  we have  $w \perp L_f$  (nonlinearity only affects the rate of change in the direction of  $w$ )
- Symmetry:** Units can be arbitrarily numbered: for Permutation matrix  $P$ , we have  $F(x) = P^{-1} \phi(P(Wx + b)) \Rightarrow$  map invariant under simultaneous permutation of in- and out-weights of a layer.
- Also implies parametrization are never unique in feedforward network
- Nesting:** Layer can be nested within layer of width  $m+1$  by adding a unit with weight  $(v = 0)$  zero (barren), by having the same unit twice where the weights sum to old weight, or  $w = b = 0$

**Activation functions**

- Linear:  $\phi = id$ , Sigmoid:  $\phi = \sigma = \frac{1}{1 + e^{-x}} \in (0, 1)$ , ReLU:  $\phi = \max(0, z)$

- Loss functions:** Squared  $\frac{1}{2} \|y - \hat{y}\|^2$ , binary  $\mathcal{L}$  loss (see below)
- Combination of PMFs  $p$  and  $q$  is  $-\sum_i p_i \log q_i \Rightarrow$  Define CE for  $p = y \in \{1, \dots, m\}$ ,  $Q(y, q) = -\log(q_y)$
- $q_y$  outputs of Softmax  $= \frac{e^{z_y}}{\sum_i e^{z_i}} \approx$  probability of class  $i$

CE loss can be written in terms of logits  $z$ :  $\ell(y; z) = [-z_y + \log \sum_i e^{z_i}] \frac{1}{m}$   
Train risk:  $\mathcal{H}(S)$ . Population risk  $E_P[h(\theta)]$  where  $(x, y) \sim P$ : can be estimated with held-out data or CV.

### 2.3 Theorems and Results:

- Lin. networks don't gain represent. power from depth ( $W^2 W^1 \Rightarrow W$ )
- Lin. layers can reduce rep. power: Project to a low dim rep ( $m_2 < m_1$ ):  
 $\rightarrow \text{rank}(F) = \text{rank}(W^L \dots W^1) \leq \min_{i=1}^L m_i \leq n$  - reduction in representational power is limited by the layer of smallest width
- Linear. AE** attempt to learn identity to find an embedding by learning linear map  $G(x) = V W x \rightarrow h(x) = \frac{1}{2} \|V W x - x\|^2$ ,  $W, V^T \in \mathbb{R}^{m \times n}$   
 $\text{rank}(V W) \leq m$
- Any optimal  $W$  will map data  $n$ -space to subspace spanned by  $m$ -PCs of sample patterns

**Residual Layers:** Often, however non-trivial to learn identity. So, initialize layers around identity so re-learning not necessary  $\rightarrow$  residual layer

- Add  $x$  to  $F$  output directly to ease learning of identity (Skip connection), often hard to find  $F[W, b](x) = x + \phi(Wx + b) - \phi(0)$ , s.t.  $F[0, 0] = id$
- Compositions with residual layers will change number of paths to grow exponentially:  $(2 \text{ layers}) F[\theta](x) = x + \phi(W^1 x) + \phi(W^2 x) + \phi(W^2 x) + \phi(W^1 x)$

- doubles  $2$  summands for  $3$  layers and so on
- Finally note that as  $x$  does not depend on  $[W, b]$ : gradient map remains unchanged  $\nabla F = \nabla \phi(Wx + b)$ . Activations change though, naturally.
- Linear projection can more easily preserve existing features. However, Residual layers do not reduce dimensionality as  $W$  is square matrix, so one often projects to allow flexibility of residual layers wrt to dimensions:  $F[V, W, b] = \phi(Wx + b) + Vx$  where  $V, W \in \mathbb{R}^{m \times n}$
- Normal res layer recovered for  $n = m$  and  $V = I$
- DenseNet  $x^{i+1} = F^{i+1}(x^i, \dots, x^0 = x)$ . Issue: dimensionality  $\uparrow$
- Sigmoid Networks:**  $\phi(z) = \sigma(z)$  ( $\tanh(z) = 2\sigma(2z) - 1 \rightarrow \tanh$  and  $\sigma$  lead to equally expressive networks, tanh preferred for sign symmetry)
- Universal Approx:** Let  $\Phi_n = \text{span}(\sigma(w^T x + b)) w \in \mathbb{R}^n, b \in \mathbb{R}$

$\Rightarrow$  can approximate any continuous functions arbitrarily well  $g$  sufficiently wide, single hidden layer MLPs? YES We want to use uniform guarantees for all points in a compact domain  $K$ . So, we obtain *uniform norm*. Let:

- $\|f\|_{\infty, K} = \sup_{x \in K} |f(x)|$  and  $d_K(f, g) := \|f - g\|_{\infty, K}$

- Assume we have a function class  $\mathcal{G}$  with which we want to approx.  $f$  over  $K$ . Then
- Approximation error of function class  $\mathcal{G}$ :  $d_K(f, g) = \inf_{g \in \mathcal{G}} d_K(f, g)$

- If  $d_K(f, g) = 0 \Rightarrow f$  is approximated by  $g$  on  $K$
- Alternatively,  $f$  approximated by  $\mathcal{G}$  for  $\exists \{g_1, \dots\} \in \mathcal{G}$  then  $\forall \epsilon > 0 : \exists M \geq 1 : \forall m \geq M : \|\mathcal{G}_m - f\|_{\infty, K} < \epsilon$
- Func. class is approx. by  $\mathcal{G}$  over  $K$  if  $\forall f \in \mathcal{F} : d_K(f, g) = 0$

$\Rightarrow$  If this holds for all compact sets  $K_j$ , then  $\mathcal{G}$  is a univ. approx. for  $\mathcal{F}$

- A point  $x$  is in  $\text{cl}(G)$ , if for every neighborhood, there is a point  $G$  within an  $\epsilon$ -distance of  $x \rightarrow$  It must hold that  $F \subseteq \text{cl}(G)$  i.e. the largest function class that can be approximated by  $G$  is its closure
- $\rightarrow$  For every  $f$  and  $m$  there is a simple MLP of width  $m$  which will achieve error  $\epsilon = \epsilon_m$  which can be made arbitrarily small by  $\uparrow m$
- Def: Let  $C(\mathbb{R}) = \{f : \mathbb{R}^n \rightarrow \mathbb{R}\}$  is continuous on  $\mathbb{R}^n$ .
- Weierstrass Thrm:** Polynomials are univ. approx. of  $C(\mathbb{R})$

**Smooth func. Approx:** Let  $\phi \in C^\infty(\mathbb{R})$  with all  $\infty$  differentiable functions, but not a polynomial  $\rightarrow \text{span}(\{\phi(ax + b)a, b \in \mathbb{R}\})$  univ. approx.  $C(\mathbb{R})$   
- For  $n=1$ , an MLP with smooth activ. **which isn't a polynomial** is a univ. approx.

**Lifting Lemma:** Let  $\phi$  be s.t.  $\text{span}(\{\phi(ax + b)a, b \in \mathbb{R}\})$  univ. approx.  $C(\mathbb{R})$  then  $\text{span}(\{\phi(w^T x + b)w \in \mathbb{R}^n, b \in \mathbb{R}\})$  univ. approx.  $C(\mathbb{R}^n)$   
**For any smooth act. (upper bounds on necessary width  $m$  though)**

- Single hidden layer needed: addit. layers could reduce err quicker
- Why not polynomials? span of  $\phi$ : degree  $k$  polynomial is also degree  $k$  polynomials. In other words: Polynomial Activation  $\Rightarrow$  Polynomial Network. With  $L$  layers, entire network is polynomial in  $x$  of finite degree  $k^L$  and univ.-variate finite degree polynomials are not dense in  $C(\mathbb{R}^n)$  in compact!

**How many units needed for approx. acc?** ? Answer: **Barron's Thrm** for any  $f$  s.t. :  $\int_{\mathbb{R}^n} |f(x)| dx < \infty$  let us define Fourier transform:  $\hat{f}(w) = \int_{\mathbb{R}^n} e^{-2\pi i w \cdot x} f(x) dx, f : \mathbb{R}^n \rightarrow \mathbb{C}$ . Barroon Thrm applies for  $f$  fulfilling regularity condition  $C_f := \|\omega\| |\hat{f}(w)| d\omega < \infty$  and  $r > 0$

- Then  $\sum_{k=1}^n \mu(k)$  sum of one hidden layer MLPs  $\{g_m\}_{m \in \mathbb{N}}$  s.t.  $\int_{\mathbb{R}^n} (f(x) - g_m(x))^2 dx \leq O(\frac{1}{m})$ ;  $rB := [x \in \mathbb{R}^n : \|x\| \leq r]$  and  $\mu$  is a meas. on  $rB$
- So, a single hidden layer MLP can approximate functions with finite  $C_f$  at order  $\frac{1}{m}$  and independently of input dimensions  $n$  and data distribution!
- Approximation error of MLPs:  $\alpha \frac{1}{m}$
- MLPs avoid the curse of dimensionality (if  $C_f < \infty$ ) and can increase convergence by increasing width!
  - However, methods that use fixed set of  $m$  basis fct.  $\rightarrow$  best approx. order is  $(\frac{1}{2})^m$

### 2.4 ReLU Networks $\phi(z) = \max(0, z)$

- Leaky relu to obtain gradient info in zero-branch:  $e' w^T x$  for  $x \leq 0$  for  $e \in \mathbb{R}$
- Let  $\Theta(W + b) \in \{0, 1\}^m$ , be the binary activation pattern of ReLU units for  $\Theta \Rightarrow$  Heavyside function
- Space can be partitioned into  $X_0 = \{x : \Theta(Wx + b) = 0\}$  and  $X_1 = \dots = 1$   $\rightarrow$  restricted to any  $X_{0/1}$ , ReLU layer map is affine
- How expressive are ReLU Layers? Answer: **Thrm (Zaslavsky):** Let  $\mathcal{H}$  be a set of  $m$  hyperplanes in  $\mathbb{R}^n$  and  $\mathcal{R}(\mathcal{H})$ : number of connected regions of  $\mathbb{R}^n - \mathcal{H}$ . Then  $\mathcal{R}(\mathcal{H}) \leq \sum_{i=0}^n \binom{n}{i} \Rightarrow R(m) \approx \frac{1}{2} m^n$ . The inequality is an

- equality if hyperplanes in  $\mathcal{H}$  are in general position
- Number of linear cells is a measure of expressive. of ReLU layers. **Complexity grows subexponent.** only once  $m > n \Rightarrow$  inefficiency of shallow ReLU Networks.
- Rep Power of deep ReLU Network: **Thrm (Motufar):** Consider ReLU network with  $L$  layers of width  $n \rightarrow R(m, L) \geq R(m) \left(\frac{m}{n}\right)^{\lfloor n(L-1) \rfloor} \rightarrow$  exponential growth for  $m \geq 2n$ . So represent. adv. for deep MLNs!
- One of very few theoretical results showing advantage of depth in DNNs!!
- Universality of ReLU: **Thrm (Shekhtman):** Piecewise linear functions are universal approximators of  $C(\mathbb{R})$ , then sufficient to show ReLU can re-present arbitrary piecewise functions which is a classic result of Thrm (Lebesgue)

- One then obtains: **Thrm (ReLU Universality):** Networks with one hidden layer of ReLU units are universal function approximators
- What about minimal non-linearity required for univ. approximators? We need functions that generalize ReLUs

- $m$ -Hinge fct.  $g[W, b](x) = \max_{i=1}^m w_i^T x + b_i$  (Maxout units)
- Thrm:** Every continuous piecewise linear function  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  can be written as a signed sum of  $m$ -Hinges with  $m \leq n+1$
- Thrm:** Every continuous piecewise linear function  $f$  can be written as the difference of two polyhedral functions as follows  $f(x) = \max_{i \in \mathcal{I}} \{w_i^T x + b_i\} - \max_{j \in \mathcal{J}} \{w_j^T x + b_j\}$  where  $A \subset \mathcal{I}$  and  $A^c \subset \mathcal{J}$
- Minimality result:** Two max operations are a sufficient element of non-linearity to obtain universal approximators

### 2.5 Exercises

**Sigmoid vs Tanh:**  $F(z) = W_2 \tanh(W_1 x + b_1) + b_2 = W_2 [2\sigma(2W_1 x + 2b_1) - 1] + b_2 = W_2' \sigma(W_1' x + b_1') + b_2'$  with  $W_1' = 2W_1, W_2' = 2W_2, b_1' = 2b_1, b_2' = b_2 - W_2$

**Linear Auto-encoders:** Let  $X = [x_1 \dots x_n] \in \mathbb{R}^{d \times n}$  be the data matrix and  $A = DE$  be a linear autoencoder with  $E \in \mathbb{R}^{k \times d}, D \in \mathbb{R}^{d \times k}$  and  $\hat{X} = AX$  be the reconstructed matrix. Thrm:  $\text{rank}(A) = \text{rank}(DE) \leq \min(\text{rank}(D), \text{rank}(E)) \leq \min(d, k)$  with  $\text{rank}(E) \leq \min(d, k)$  and  $\text{rank}(D) \leq \min(d, k)$

- By Eckart-Yong-Mirsky Thm.:  $\hat{X}^* = \text{argmin}_{\hat{X}} \|\hat{X} - X\|_F^2 = U_k X_k V_k^T$ , where  $X = U_k X_k V_k^T$  is the truncated SVD of  $X$
- The optimal weight matrices are given by the first  $k$  columns (ranked according to dominant singular values) of  $U$  combined with any invertible  $A \in \mathbb{R}^{k \times k}, D^* = U_k A, E^* = A^{-1} U_k^T$

- $\|X - DE\|_F^2$  is convex in  $E$  and  $D$  but not jointly
- Thm. (Weierstrass):** If  $f(x)$  is a continuous function for  $a \leq x \leq b$  and  $f$  is an arbitrary positive quantity, it is possible to construct an approximating polynomial  $P(x)$  s.t.  $|f(x) - P(x)| \leq \epsilon$ ,  $a \leq x \leq b$  where w.l.o.g. one can assume  $0 < a < b < 1$  and  $f(x) > 0$  outside of  $(a, b)$

## 3 Gradient-based learning



**Receptive Fields:** Receptive field of  $x_i^t$ :  $\mathcal{R}_i^t := \{j : w_{ij} \neq 0\}$  where  $W^t$  is the toplitz matrix of the convolution (includes all the elements in the previous layer that effect the calculation of  $x_i^t$ ). Then  $\frac{\partial x_i^t}{\partial x_j^0} = 0$  for  $j \notin \mathcal{R}_i^t$ .

Con: local receptive field make it hard to connect distant features.

**Weight sharing:** When computing the loss gradients  $\partial \mathcal{L} / \partial \theta_{ij}^k$  ( $\theta_{ij}^k$  is a kernel weight) we have to respect the weight sharing:  $\frac{\partial \mathcal{L}}{\partial \theta_{ij}^k} = \sum_l \frac{\partial \mathcal{L}}{\partial x_l^t} \frac{\partial x_l^t}{\partial \theta_{ij}^k}$ . The sum over units appears as the same weight is re-used for every unit within the target layer.

- $V_{ij}(v^T \text{vec}(\sigma(x * w))) = \text{vec}(\sigma(x * w))$
- $\nabla_A v^T \text{vec}(A) = \text{mat}(v)$
- $\nabla_A (x * w) \in \mathbb{R}^{q \times p \times r \times s}$  if  $(x * w) \in \mathbb{R}^{r \times s}$ ,  $w \in \mathbb{R}^{p \times q}$
- $\nabla_u (v^T \text{vec}(\sigma(x * w))) = \text{rot}_{\pi}(x) * B$  where  $B = \text{mat}(v) \odot \sigma'(x * w)$

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(1, 1, 3)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(1, 1, 3)
        self.fc1 = nn.Linear(1 * 1 * 1, 5 * 128)
        self.fc2 = nn.Linear(128 * 5, 128)
        self.fc3 = nn.Linear(128 * 5, 10)

net = Net()
optimizer = optim.Adam(net.parameters())
```

- #### 4.4 Inception Network
- Deep stacking, uses many channels for accuracy. Dimension reduction in between convolutions (compression of m channels) using a  $1 \times 1 \times m$  convolution with  $k \leq m$ :  $x_{ij}^t = \sigma(W_{kij})$ ;  $W \in \mathbb{R}^{k \times m}$
- $1 \times 1$  part means there is no spatial conv. performed
  - Not committing to a certain window size, but solve trade-off problem a part of the learning by multiple processing paths
  - Softmax layers to improve learning dynamics (shortcut error backprop)

- #### 4.5 Embeddings
- Enable processing of non-numeric sequences via DNN's by embedding symbols in vector space.  $\Omega$  is the alphabet of symbols;  $\Omega \ni \omega \mapsto x_\omega \in \mathbb{R}^n$ . Embedding  $\{x_\omega\}$  are learnable parameters.
- Word2vec:** Learn input and output embeddings:  $\Omega_{in} \ni \omega \mapsto x_\omega \in \mathbb{R}^n$  and  $\Omega_{out} \ni v \mapsto y_v \in \mathbb{R}^n$ .

Combine:  $P(v|\omega) = \frac{\exp(x_\omega^T v)}{\sum_{y \in \Omega_{out}} \exp(x_\omega^T y)}$  and use to predict an output word  $v$  in neighbourhood of an input word  $\omega$

## 5 Recurrent Neural Networks

Idea: Provide more flexibility than CNNs to model temporal/sequential data (must specify filter width in CNNs).

- #### 5.1 Setup
- Observation sequence:**  $x^1, x^2, x^3, \dots$ , with  $x^t \in \mathbb{R}^n$  and **Hidden states:**  $z^1, z^2, z^3, \dots$ , with  $z^t \in \mathbb{R}^n$ . There are two operators:
- Connecting  $x$  to hidden states:**  $F: U, V(z, x) \mapsto \psi(Uz + Vx)$ , with  $U \in \mathbb{R}^{m \times m}$ ,  $V \in \mathbb{R}^{m \times n}$ .
- Output map (optional):**  $G[W](z) \equiv \psi(Wz)$ , where  $W \in \mathbb{R}^{q \times m}$ .
- Rem. 1:** Both  $F$  and  $G$  are *time-invariant*, (i.e., shared parameters).
- Rem. 2:** With a fixed time horizon  $1 \leq T$ , the unrolled RNN is equivalent to a feedforward neural network with  $T$  hidden layers. The main difference is the sharing of parameters between layers.

- 3.3. (Bi-directional RNN) Hidden state evolution** does not have to follow direction of time. One can define a reverse order sequence  $z^t = \phi(\bar{U}z^{t+1} + \bar{V}x^t)$  and modify the output map to be  $\hat{y}^t = \bar{\psi}(z^{t+1} + \bar{W}\bar{z}^t)$ .
- Rem. 4:** (Stacked RNN) Possible to divide up by connecting layers horizontally:  $z^t = \phi(U^t z^{t-1} + V^t z^{t-1})$ , where  $z^{t,0} \equiv x^t$ . Outputs are generated based on the last hidden activations  $z^{t,L}$

- #### 5.2 Backpropagation through time
- Error signals at the outputs:  $\delta_k^t = \frac{\partial \mathcal{L}}{\partial y_k^t}$
  - Train the **output weights:**  $\frac{\partial \mathcal{L}}{\partial w_{k,j}} = \sum_{i=1}^n \delta_k^i \psi_k' x_{ij}^t$ , with  $\psi_k' \equiv \psi(w_k^t)$
  - Train the **input weights:**  $\frac{\partial \mathcal{L}}{\partial w_{ij}} = \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial y_k^i} \phi_k' x_{ij}^t$ , with  $\phi_k' \equiv \phi(v_k^t)$
  - with  $\phi_k' \equiv \phi(v_k^t)$
  - which needs...  $\frac{\partial \mathcal{L}}{\partial z_i^t} = \sum_{k=1}^n \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial y_k^i} \sum_{j=1}^n \frac{\partial y_k^i}{\partial z_j^t} \frac{\partial z_j^t}{\partial z_i^t}$ , with  $\frac{\partial y_k^i}{\partial z_j^t} = \psi_k' w_{kj}$
  - Train the state-to-state weights:  $\frac{\partial \mathcal{L}}{\partial w_{ij}} = \sum_{i=1}^n \frac{\partial \mathcal{L}}{\partial z_i^t} \frac{\partial z_i^t}{\partial z_j^{t-1}}$

- #### 5.3. (Instability) of RNNs
- In backprop. through time, the norm of the gradients is difficult to control. We have:  $\frac{\partial z^t}{\partial z^0} = \Phi^T U \dots \Phi^1 U$ , with  $\Phi^t = \text{diag}(\phi(w_1^t z^{t-1} + v_1^t x^t))$ .
- For typical activations, diagonal sensitivity matrices are bounded  $\Phi^t \leq \alpha \mathbb{I}$  ( $\alpha = 1$  for ReLUs,  $\alpha = 1/4$  for logistic activations). We have:
- $$\left\| \frac{\partial z^t}{\partial z^0} \right\|_2 \leq (\alpha \alpha_1) (U)^T, \text{ which ...}$$

- $\dots \rightarrow 0$  as  $T \rightarrow \infty$  if we assume  $\alpha_1(U) < \frac{1}{\alpha}$  (**Vanishing gradients**)
  - $\dots \rightarrow \infty$  as  $T \rightarrow \infty$  if we assume  $\alpha_1(U) > \frac{1}{\alpha}$  (**Exploding gradients**)
  - $\rightarrow$  **Makes it difficult to model long-term dependencies**
- #### 5.4 Gated Memory
- $\rightarrow$  *Avoid short-term fluct. by controlling when memory is kept or overwritten*
- A gating unit takes the following form:  $z^t = \alpha(GC) * z$ , where  $z^t \in \mathbb{R}^m$ ,  $\alpha(GC) \in (0,1)^m$  and  $*$  denotes pointwise multiplication.  $\sigma$  is the vector-valued logistic function, and  $\zeta$  is a control input determining the relative preservation of information contained in  $z$ . In the limit  $\sigma_i \rightarrow 0$ ,  $z_i$  is forgotten, while the limit  $\sigma_i \rightarrow 1$  preserves  $z_i$

- Long Short-Term Memory (LSTM):**
  - Uses two gating units, a **forget gate** and a **storage gate**:  $z^t = \sigma(Fx^t) * z^{t-1} + \sigma(Gx^t) \cdot \tanh(Vx^t)$ . Here,  $x^t = [x^t, \zeta^t]$  (concatenated), and  $\zeta^{t+1} = \sigma(Hx^t) \cdot \tanh(Vx^t)$  is a control signal sequence.
  - $\rightarrow$  An LSTM has **5 weight matrices in total**: 3 gating matrices ( $F, G, H$ ), 2 propagation matrices ( $V, U$ ) (6 incl. output weight matrix:  $W$ ).
- Gated Recurrent Unit (GRU):**

The GRU combines the forget and storage gates of an LSTM into a convex combination  $z^t = (1 - \sigma) * z^{t-1} + \sigma * z^t$ , with  $\sigma = \sigma(G[x^t, z^{t-1}])$ , and  $z^t = \tanh(W[C * z^{t-1}, x^t])$ , where  $C' = \sigma(H[z^t, x^t])$ .

- Adv.  $\zeta^t$ :** can be computed recursively without recursion.
- $\rightarrow$  Only **3 weight matrices in total** (4 incl. output weight matrix).
- #### 3. Minimal Gated Units:
- Idea: Remove dependencies on previous hidden states from gatin mechanism:  $\sigma = \sigma(Gx^t)$ , and  $z^t = Vx^t$ , which also removes the squeezing of the dynamical range via tanh

- #### 5.5 Linear Recurrent Models
- Background:** Classical RNNs lack sufficient parallelization during learning. Transformers, on the other hand, scale quadratically and not just linearly in context length.  $\rightarrow$  Linear recurrent models like S4 or Mamba avoid non-linear dependencies along the chain of hidden states.
- Linear Recurrent Unit (LRU):**
 $z^{t+1} = Az^t + Bx^t$ . The idea is that we can parameterize and run this recurrence over the complex numbers in a way that the architecture & parameterization guarantee stability:
    - Diagonalize the evolution matrix  $A$**  (assuming it is non-defective) over  $\mathbb{C}$  as follows:  $A = PAP^{-1}$ , where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ , and  $\lambda_i \in \mathbb{C}$ .
    - Perform a change of basis:**  $C^{t+1} = \Lambda C^t + Cx^t$ , with  $C^t \equiv P^{-1}z^t$ , and  $C \in \mathbb{C}^{m \times n}$ . We make this equation our recurrence.
    - Stability of this (linear) system** requires the modulus of all eigenvalues to be bounded:  $\max_i |\lambda_i| \leq 1$ , where  $|a + ib| = \sqrt{a^2 + b^2}$  (modulus)  $\rightarrow$  spectral radius of  $A$  (or  $\Lambda$ )
    - Control the spectral radius:** Use the polar form to represent complex number  $z = r(\cos(\phi) + i \sin(\phi))$ , with modulus  $r = |z| \geq 0$ , and phase  $\phi \in [0, 2\pi)$ . Stable parameterization with double exp. for the  $\lambda$ s:  $\lambda_i = \exp(-\exp(v_i) + i\phi_i) \in (0,1)$ . This can be done by sampling at initialization from the ring defined by  $I$ :
    - Move the non-linearities into the output map** to compensate for the simpler linear recurrence:  $y^t = \text{MLP}(\text{RE}(Gz^t))$ ,  $G \in \mathbb{C}^{k \times m}$ .  $\rightarrow$  The resulting model is universal as a seq-2-seq model.

- #### 5.6 RNN Models for Sequences
- Generic model  $p(y^1 T | x^{1:T}) \approx \prod_{t=1}^T p(y^t | x^{1:t-1})$ . Then we have  $x^{1:t} \xrightarrow{F} z^t$  and to get a distribution and not a deterministic output, the deterministic computation of the RNN is augmented with a final probabilistic step with a parametric family  $z^t \mapsto \mu^t$ ,  $\mu^t \mapsto p(y^t; \mu^t)$ . This could be a softmax:  $p(y; \mu) = \frac{\exp(\mu_y)}{\sum_{\mu'} \exp(\mu_{\mu'})}$ ,  $y \in \{1, \dots, k\}$ ,  $\mu \in \mathbb{R}^k$ .
- Note. We need to feed back output of stochastic generation because otherwise,  $y^t$  will depend on  $y^{1:t-1}$  only through  $z^t$  which is too strong an independence assumption to make.

- #### 5.7 Teacher Forcing
- Fixing the problem that predicted sequences may quickly diverge from the target one due to long-range dependencies, resulting in error signal of little value. Teacher forcing feeds back the target sequence  $y^t$  during training, instead of the generated outputs. **Adv.:** No lasting effect of prevs. during training as they get overwritten by ground-truth. **Disadv.:** Creates divergence between training & testing because it can only be applied during training (**Exposure Bias**).
- #### 5.8 seq-2-seq
- Encoder-Decoder based architecture. The encoder maps input sequence to dimensional activation vector via encoder RNN without outputs:  $z^0 \ni 0$ ,  $(x^t, z^{t-1}) \mapsto z^t \quad \forall 1 \leq t \leq T$ . Then, a decoder RNN without inputs is used:  $\bar{z}^0 \ni z^T$ ,  $(\bar{y}^{t-1}, \bar{z}^{t-1}) \mapsto \bar{z}^t$ ,  $\bar{z}^t \mapsto y^t$ .
- As before, outputs are fed back and next output is generated stochastically based on the hidden state of the RNN.

## 6 Transformers

- #### 6.1 Token Paradigm
- Given is a sequence of tokens in the form of their respective embeddings  $X = [x_1, \dots, x_T] \in \mathbb{R}^{n \times T}$ . The fundamental idea is to map the non-contextualized embeddings  $X$  to contextualized representations  $\Xi = [\xi_1, \dots, \xi_T] \in \mathbb{R}^{m \times T}$ .

- #### 6.2 Attention Mixing
- An attention mechanism produces numbers  $\alpha_{s,t}$  that are then used to convexly combine the input representations into new representations  $\xi_t = \sum_s \alpha_{s,t} Wx_s$ , where  $\alpha_{s,t} \geq 0$  and  $\sum_t \alpha_{s,t} = 1$ .

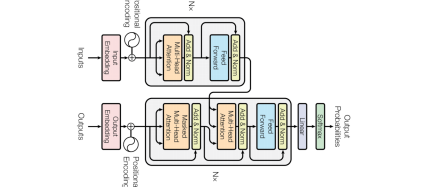
- Remarks:**
- The attention weights  $\alpha_{s,t}$  have a **source** (where attention emerges,  $s$ ) and a **target** (where attention extracts information,  $t$ )
  - In matrix notation, we write  $A = (a_{s,t}) \in \mathbb{R}^{T \times T}$ , such that  $\Xi = WXAT$ .
  - Attention mechanism mixes information across columns of  $X$  via  $AT$ , and mixes dimensions linearly via  $W$ .

- #### 6.3 Query-Key Matching
- Question:** How to learn  $A$ ?
- Project  $X$  to query matrix  $Q = U_Q X$  and key matrix  $K = U_K X$ , where  $U_Q, U_K \in \mathbb{R}^{n \times q}$
  - Produce matching matrix by matching keys and queries via inner products:  $Q^T K = X^T U_Q^T U_K X \in \mathbb{R}^{T \times T}$  (note that  $\text{rank}(U_Q^T U_K) \leq q$ )
  - Soft-max transformation:  $A = \text{softmax}(q^T Q^T K)$  where  $a_{s,t} = \frac{\exp(iQ^T K)_{s,t}}{\sum_{k=1}^T \exp(iQ^T K)_{s,k}}$

**Inverse temperature  $\beta$ :** Controls the entropy of the softmax and is not learned, but usually chosen as  $\beta = 1/\sqrt{q}$  (to standardize the size of  $Q^T K$  dot products to be independent of  $q$ ). **Higher  $\beta$  leads to lower entropy** of the softmax and the **attention becomes more selective/focused** on the highest-scoring queries/keys (and vice-versa for lower  $\beta$ )

- #### 6.4 Multi-Headed Attention
- Idea:** Run multiple attention models in parallel.
- Replicate the matrices  $U_Q, U_K, W$  a total of  $r$  times
  - Perform attention-based propagation  $X \mapsto \Xi_i$  for  $i = 1, \dots, r$
  - Concatenate the matrices  $\Xi_i$  along the feature dimension
- Concrete numbers from *Attention Is All You Need* are: feature space dimension  $n = 512$ , number of heads  $r = 8$ , group/key space dimension  $q = 64$ .
- Remark:** Originally proposed values for  $r, q$  and  $n$  ensure that after concatenating the different heads, the output dimension stays the same ( $8 \cdot 64 = 512$ ). The reasons for this trick:
- Residual connections: Would not be possible with different dimensions.
  - Efficiency costs: Less parameters to estimate.
  - Scalability: The output dimension would increase with every layer by a factor of  $r$ , if we would not use this trick.

- #### 6.5 Feature Transformation
- Issue:** Attention mechanism only performs convex combination of columns of  $X$ , which is insufficient (even if the relative weights are non-linear functions of  $X$ ).
- Solution:** Gain representational power by **applying a per-token** transformation with a feedforward network. Formally,  $X \mapsto \Xi \mapsto F(\Xi)$ , where  $F(\Xi) = (F(\xi_1), \dots, F(\xi_T))$  operates on the columns of  $\Xi$ .
- #### 6.6 Positional Encoding
- Issue:** Attention mechanism operates on (unordered) set of tokens ( $\leadsto$  value of token position is lost)
- Solution:** Assign each non-negative integer  $t$  to a low-dimensional vector and add this vector to the respective token embedding  $x_t$ .
- Sinusoidal Encoding:**  $p_{tk} = \begin{cases} \sin(\omega_k t) & k \text{ even} \\ \cos(\omega_k t) & k \text{ odd} \end{cases}$ , where  $\omega_k = C^{n/k}$  and  $n$  is the feature space dimension ( $= 512$  in the original paper)
- #### 6.7 Masked and Cross-Attention
- Masked Attention:** In cases where transformers operate in coupled encoder-decoder pairs, the decoder typically operates in autoregressive fashion (output at  $t$  depends on encoder, but also on produced token representations for  $s < t$ ). Masked attention restricts the attention to the past. **Cross-Attention:** Attention to representations produced by the encoder. Cross-Attention is implemented across layers of the same depth. Note that in cross-attention the queries emerge from the decoder, whereas the keys and values are derived from the encoder states.
- #### 6.8 Transformer



- #### 6.9 Large Language Models (BERT)
- Acronym:** Bidirectional Encoder Representations from Transformers.
  - Bidirectional:** BERT processes text in both directions simultaneously.
  - Pre-Training Tasks:**
    - Masked Language Modeling** (Cloze test-like training where 15% of words are masked for prediction). The to be predicted words are replaced by MASK token in 80% of cases, by random token in 10% of cases, and left unchanged for the remaining 10% of cases.
    - Sentence pair classification. Specifically, classifying if two sentences are consecutive or not.
    - Tokenization:** BERT operates on tokenized text. Specifically, the WordPiece tokenizer was used. Additionally, a CLS token (and its final embedding) was used for classification problems. Also, a SEP token was used for marking the separation of two sentences.

- #### 6.10 Vision Transformers
- Idea:** Decompose image into non-overlapping patches as the tokens of an image (typically into patches of size  $16 \times 16$  pixels). Note how usually we flatten the tokens into vectors. Formally,  $\mathbb{R}^{p \times p \times q} \mapsto \text{patch}^t \mapsto x^t \equiv V \text{vec}(\text{patch}^t) \in \mathbb{R}^n$ ,  $V \in \mathbb{R}^{n \times p^2 q^2}$
- Issues:** This pre-processing ignores the 2D structure of pixel arrangements within a patch. If enough data, this is not really an issue!
- Remark:** Early ViTs have used two layer GELU (Gaussian error linear unit) MLP. The GELU activation is given by,  $\phi(z) = z \text{Prob}(Z \leq z)$ ,  $Z \sim \mathcal{N}(0,1)$  and then apply standard DNN problems: variable-length inputs if sets have different cardinality & arbitrary means of ordering of elements

- #### 6.11 Complexity
- Let  $n$  be the sequence length,  $d$  the representation dimension,  $k$  the kernel size of convolutions and  $r$  the size of the neighborhood in restricted self-attention
- | Layer-Type                  | Complexity per Layer     | Sequential Operations | Maximum Path Length |
|-----------------------------|--------------------------|-----------------------|---------------------|
| Self-Attention              | $O(n^2 \cdot d)$         | $O(1)$                | $O(1)$              |
| Recurrent                   | $O(n \cdot d)$           | $O(n)$                | $O(n)$              |
| Convolutional               | $O(k \cdot n \cdot d^2)$ | $O(1)$                | $O(\log(n))$        |
| Self-Attention (restricted) | $O(r \cdot n \cdot d)$   | $O(1)$                | $O(n/r)$            |

## 7 Geometric Deep Learning

- #### 7.1 Sets & Points
- Goal:** realize some function  $\{x_1, \dots, x_M\} \subset \mathbb{R}$ ,  $f: 2^{\mathbb{R}} \rightarrow \mathcal{Y}$ . Naively, we could turn the set into an ordered  $N$ -tupel  $\{x_1, \dots, x_M\} \mapsto (x_1, \dots, x_M) \in \mathbb{R}^M$  and then apply standard DNN problems: variable-length inputs if sets have different cardinality & arbitrariness of ordering of elements
- Invariance and Equivariance:**
- $f$  is **order-invariant** iff  $f(x_1, \dots, x_M) = f(x_{n_1}, \dots, x_{n_M})$   $x \in \mathbb{R}^M$ ,  $\forall n \in S_M$ ,  $n$ : permutation,  $S_M$ : symmetric group
  - $f$  is **order-equivariant** iff  $f(x_1, \dots, x_M) = (y_1, \dots, y_M) \implies f(x_{n_1}, \dots, x_{n_M}) = (y_{n_1}, \dots, y_{n_M})$   $x \in \mathbb{R}^M$ ,  $\forall n \in S_M$
- More general, for some transformation  $T \in \mathcal{T}$ , where  $\mathcal{T}$  is some group of transformations, and a function  $f$ , we call  $f$
- invariant** iff  $f(T(x_1, \dots, x_M)) = f(x_1, \dots, x_M)$   $\forall x \in \mathbb{R}^M, \forall T \in \mathcal{T}$
  - equivariant** iff  $f(T(x_1, \dots, x_M)) = T(f(x_1, \dots, x_M))$   $\forall x \in \mathbb{R}^M, \forall T \in \mathcal{T}$
- Examples: conv. operator is equivariant to translations; att. mechanism permutation equivariant with regard to token sequence since softmax is equivariant operation; Euclidean norm is invariant under rotations
- Deep Sets:**
- Sum is invariant under permutations (commutativity):  $\sum_{i=1}^M x_{n_i} = \sum_{i=1}^M x_i$ ,  $\forall M \in \mathbb{N}$ ,  $\forall M, \forall n \in S_M \leadsto$  by allowing for other (fixed or DNN) maps  $\phi: \mathbb{R} \rightarrow \mathbb{R}^N$  and  $\varphi: \mathbb{R}^N \rightarrow \mathbb{R}$  we can construct a larger class of invariant function as follows:  $f(x_1, \dots, x_M) = \varphi(\sum_{i=1}^M \phi(x_{n_i}))$ , where  $\phi, \varphi$  are independent of  $M$ , so  $f$  is defined over an arbitrary number of inputs.
  - Can construct invariant mappings using different agg. functions such as (point-wise) max. & min., etc:  $f(x_1, \dots, x_M) = \varphi(\max_{i=1}^M \phi(x_{n_i}))$

- Once we have invariant  $f$ , this can be turned into an equiv. map using  $\rho: \mathbb{R} \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ ,  $(x_{n_i}, \sum_{i=1}^M \phi(x_{n_i})) \mapsto y_{n_i}$  and apply pointwise for each  $x_{n_i}$ . Hence,  $\rightarrow$  down-stream processing  $p$  depends on an invar. argument (here: aggregated as sum) and  $x_{n_i}$  itself (via shortcut connection)
- fixed  $N$  (width of hidden representations) is sufficient if one allows for discontinuous mappings in the limit of  $M \rightarrow \infty$ , but more practically realizable mappings may require  $N \geq M$

**PointNet:** Sparse measurements at specific 3D points i.e. point cloud  $\leadsto$  desirable that a DNN obeys permutation in-/equivariance

- #### 7.2 Graph Convolutional Networks
- Basic Set:** of Vertices  $\mathcal{V} = \{v_1, \dots, v_M\}$ . We associate a feature vector  $x_{v_i}$  with every node and consider undirected edges  $\mathcal{E} = \{e_1, \dots, e_E\} \subseteq \mathcal{C} \in 2^{\mathcal{V}}: |e| = 2$ . Idea: graph may encode similarities between nodes and corresponding dependencies between their outputs.
- Adj. Matrix:**  $A = (a_{nm})$ ,  $a_{nm} = \mathbb{1}_{\{v_n, v_m\} \in \mathcal{E}}$ .  $A$  is symmetric and has 0 diagonal (i.e. no self-loops)

- Permut. Matrix:**  $P \in \{0, 1\}^{M \times M}$  s.t.  $\sum_{i=1}^M p_{nm} = \sum_{n=1}^M p_{nm} = 1$   $\forall m$  Cauchy's two line notation for permutations:  $P = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix} \leftrightarrow \pi = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \end{pmatrix}$ ;  $\pi^{-1} = \begin{pmatrix} 1 & 4 & 2 & 3 \\ 2 & 2 & 3 & 1 & 4 \end{pmatrix}$   $\pi \pi^{-1} = \text{id}$  always inverse elements in the group  $S_M$
- Permutation matrices are orthogonal i.e.  $P^{-1} = P^T$  or  $PP^T = P^T P = I$  as  $(PP^T)_{nm} = \sum_k P_{nk} P_{mk} = \delta_{nm}$ . Let  $\Pi_M$  be the group of permutation matrices, then  $\sum_{\pi \in \Pi_M}$  of a function of a graph can be defined as  $f(X, A) \stackrel{!}{=} f(PX, PAP^T)$ ,  $\forall P \in \Pi_M$  and equiv. as  $Pf(X, A) \stackrel{!}{=} f(PX, PAP^T)$ ,  $\forall P \in \Pi_M$

- Invariant Core:** Multiset of features in neighborhood of node  $i$ :  $X_{n_i} = \{x_{v_n} : \{v_n, v_m\} \in \mathcal{E}\}$ ,  $\{|\}\}$  = multiset i.e. for given node  $v_m$  it's the multi-set of feature vectors from its neighbors. Choose any  $\varphi: (x_{n_1}, x_{n_2}) \mapsto \phi(x_{n_1}, x_{n_2}) \in \mathbb{R}^p$ . We can then apply this function at every node to arrive
- $$x \mapsto \phi(X): \begin{pmatrix} \phi(x_1, x_1) \\ \vdots \\ \phi(x_M, x_M) \end{pmatrix} \in \mathbb{R}^{M \times N}$$

Note that any pair of isomorphic graphs (i.e. graphs that be bijective under edge re-numbering) are pressed in the same manner.  $\phi$  can be designed e.g. by transforming node features and then aggregating by summation:  $\phi(x_{n_1}, x_{n_2}) = \phi(x_{n_1}, \sum_{x \in x_{n_1}} \phi(x))$ . If globally invariant function is desired, one can add a final invariant aggregation step. **Graph-Learning Scenarios:** *node classification:* make a discrete/real-valued prediction at every node; *graph classification:* classify the entire graph; *link prediction:* predict missing edges in graph

- Coupling Matrix:** In convolutional graph NNS, the aggregation over neighborhoods is performed with a fixed sets of weights known as coupling matrix (typically derived from adj. matrix  $A$ ). Often we add self-loops to  $A$  and normalize with the degree matrix as follows:  $\bar{A} = D^{-1/2}(A + I)D^{-1/2}$ ,  $D = \text{diag}(d_1, \dots, d_M)$ ,  $d_m = 1 + \sum_{n=1}^M a_{nm}$ . Rmk: The normalization balances out degree variability between the nodes but the activation of a node propagates to its neighbors. E.g. for a 1-dim (cyclic) chain, a regular graph of degree 2, then  $\bar{A} = 1/3A$

- GCN Layers:** Let  $W$  be weight matrix and  $\sigma$  an activation function. Then, one step of propagation in a convolutional GNN:  $Z = \sigma(\bar{A}XW)$ ,  $W \in \mathbb{R}^{N \times d}$ . Note: Multiplication by  $\bar{A}$  from the left sums over nodes, based on the edge structure of the graph, whereas multiplication by  $W$  from the right sums over features. Iteratively, we have for two-layer GCN:  $Z = \text{softmax}(\bar{A} \text{ReLU}(\bar{A}XW)W_1)$

- ! Control the spectrum of the coupling matrix  $\bar{A} \mapsto \text{diag}(|\lambda(\bar{A})| \leq 1$ ,  $\rightarrow$  guarantees stability in terms of repeated propagation in a deep GCN
- Limitations of GCN:** GCNs require that depth is equal to diameter of graph to propagate between all pairs of nodes  $\leadsto$  deep networks  $\rightarrow$  issues:

- Over-smoothing:** Repr. at nodes may become indistinguishable
- Over-quashing:** Bottleneck effect of how much long-range information can be encoded in fixed-size representation  $\leadsto$  can improve learning from long-range dependencies by modifying the graph structure to allow for smaller diameters

- #### 7.3 Spectral Graph Theory:
- Laplacian Operator:** is  $\Delta f := \sum_{i=1}^n \frac{\partial^2 f}{\partial x_i^2}$ ,  $f: \mathbb{R}^N \rightarrow \mathbb{R}$ . Think of Laplace operator at  $x$  as measuring local variation from  $f$ -average in a vanishingly small neighborhood of  $x$ : for balls of radius  $h \rightarrow 0$   $\bar{f}(x, h) - f(x) = \frac{h^2}{2n} \Delta f(x) + o(h^2)$
- Graph Laplacian:** with adjacency matrix  $A$  & degree matrix  $D$ :  $L := D - A$ ,  $(Lx)_n = \sum_{i=1}^n a_{nn}(x_n - x_{n_i})$ ; but in practice often better to use (symmetric) degree-normalized Laplacian:  $L = I - D^{-1/2}AD^{-1/2} = \bar{D} - \bar{A}$   $D^{-1/2}(D - A)D^{-1/2}$ . Note:  $L \& \bar{L}$  are symmetric, positive semi-definite and weakly diagonally dominant i.e.  $d_n = \sum_m a_{nm} = \sum_m |a_{nm}|$
- Graph Fourier Trafo:**  $L = D - A = UAU^T$ ,  $\Lambda := \text{diag}(\lambda_1, \dots, \lambda_M)$ ,  $\lambda_i \geq \lambda_{i+1}$  with  $U$  orthogonal. The columns of  $U$  can be considered to be the graph Fourier basis. **Graph Convolution:** Pointwise multiplication in the Fourier domain:  $x * y := U(U^T U x) \odot (U^T y)$ . The filtering operation from 1-d signals can be generalized:  $G_\Omega(L)x = U G_\Omega(\Lambda)U^T x \leadsto O(M^3)$  since full eigendecomposition of  $L$ . **TRICK: Polynomial Kernels** i.e.  $U(\sum_{k=0}^n \alpha_k \lambda^k U^T) = \sum_{k=0}^n \alpha_k L^k$ . If  $A$  sparse ( $\implies L$  sparse), then highly efficient scaling of arithmetic operations with  $|\mathcal{E}|$  as only powers of Laplacian have to be computed. spectral conv. can then be expressed as multiple equivariant local layers in the vertex domain. Polynomial order  $K$  defines size of neighborhood (typically 2-5). These filters are isotropic and have fewer parameters. and lower express. than in traditional convs.
- Polynomial Kernel Networks:** Channel indices denoted by  $i, j$ , then pre-activation of graph-convolutional layer can be written as  $x_i^{t-1} = \sum_j p_{ij}(L)x_j^t + b_i$ ,  $p_{ij}(L) = \sum_{k=0}^K \alpha_{ijk} L^k$ . The number of parameters is essentially (ignoring  $b_i$ ) the product of the channel dimensions times  $K + 1$ .

- #### 7.4 Graphs and GNNs
- Coupling with Attention:** Define a coupling matrix



**8.3 Co-adoption:** One unit depends on the presence of other units which leads to overfitting  $\Rightarrow$  dropout unit  $i$  with prob  $(1 - \pi_i)$

**Ensemble** Each instance of Drop-out defines a sub-network.  $p(y|x) = \prod_{i \in [0,1], p(b|p(y|x);b)} R := \# \text{units and } p(b) = \prod_{i=1}^n \pi_i^{b_i} (1 - \pi_i)^{1-b_i}$

$\Rightarrow$  Dropout is like SGD on ensemble, pick a random sub-network (of  $2^R$  possible) and perform a gradient step.

**Test time:** scale  $w_{ij} \leftarrow \pi_j w_{ij}$  by prob of upstream unit being active

#### 8.4 Convexity, Smoothness and GD

**$\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$**  is convex if  $\forall w, w' \in \mathbb{R}^d, \forall \lambda \in [0, 1] : \mathcal{L}(\lambda w + (1 - \lambda)w') \leq \lambda \mathcal{L}(w) + (1 - \lambda) \mathcal{L}(w')$

- For a diff. fct  $\ell : \mathcal{L}(w) \geq \ell(w') + \nabla \ell(w')^T (w - w') \Leftrightarrow \ell$  is convex
- If  $\ell$  is diff. &  $L$ -smooth, then  $\mathcal{L}(w) \leq \mathcal{L}(w') + \nabla \ell(w')^T (w - w') + \frac{L}{2} \|w - w'\|^2$
- Strongly convex:  $\mathcal{L}(w) \geq \mathcal{L}(w') + \nabla \ell(w')^T (w - w') + \frac{\mu}{2} \|w - w'\|^2$
- GD step  $s_t := w_{t+1} - w_t = -\frac{1}{\lambda} \nabla \ell(w_t)$  (with  $\lambda > 0$ ) minimizes a quadr. reg. 1st-order Taylor e.:  $s_t = \underset{\argmin}{\text{argmin}}_{\substack{(w_t) + s^T \nabla \ell(w_t) + \frac{\lambda}{2} \|s\|_2^2}}$

$\Rightarrow$  If  $\ell$  is  $L$ -smooth, and  $\lambda \geq \frac{1}{L}$ , then  $s_t$  is guaranteed to decrease the objective, i.e.,  $\ell(w_{t+1}) \leq \ell(w_t)$   
GD strictly follows the gradient in negative direction, so its trajectory is always perpendicular to the level set

## 9 Regularization & Propagation

#### 9.1 Neural Network with Skip connections

Let  $x \in \mathbb{R}^d$  be the input to the network,  $y_i \in \mathbb{R}^d$  the intermediate representations,  $\Theta_i \in \mathbb{R}^{d \times d}$  the trainable parameters,  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  a smooth activation function, and  $L : \mathbb{R}^d \rightarrow \mathbb{R}^+$  the (smooth) loss function. Also, let  $\alpha, \beta \geq 0$  be a real parameter that controls the strength of the residual branch. Consider the residual network:

$$y_1 = \alpha f(x, \Theta_1) + x, \quad y_2 = \alpha f(y_1, \Theta_2) + y_1 \quad \dots \quad y_n = \alpha f(y_{n-1}, \Theta_n) + y_{n-1}$$

with loss  $L(Y_n)$

**For  $n = 1$**  have  $\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}(y_1)}{\partial y_1} \frac{\partial y_1}{\partial x} = \frac{\partial \mathcal{L}(y_1)}{\partial y_1} \left( \alpha \frac{\partial f(x, \Theta_1)}{\partial x} + \mathbb{I}_d \right) \in \mathbb{R}^{1 \times d}$ .

**For general  $n > 1$**  we have for some  $1 \leq k < n$ :  $\frac{\partial \mathcal{L}}{\partial \Theta_k} = \frac{\partial \mathcal{L}}{\partial y_k} \dots \frac{\partial \mathcal{L}}{\partial y_{k+1}} \dots \frac{\partial \mathcal{L}}{\partial y_n} = \frac{\partial \mathcal{L}(y_n)}{\partial y_n} \left( \alpha \frac{\partial f(y_{n-1}, \Theta_n)}{\partial y_n} + \mathbb{I}_d \right) \dots \left( \alpha \frac{\partial f(y_{k+1}, \Theta_{k+1}}{\partial y_k} + \mathbb{I}_d \right) \cdot \alpha \frac{\partial f(y_k, \Theta_k)}{\partial \Theta_k}$

**Note:** We add  $\mathbb{I}$  as the "diagonal terms", which helps for optimization.

#### 9.2 Ridge Regression

The ridge regression solution for a linear model is given by  $\hat{w} = \underset{\argmin}{\text{argmin}}_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 = (X^T X + \lambda I)^{-1} X^T y$ . Using the SVD

$X = UDV^T$ , we get for the predictions  $\hat{y} = \hat{X} \hat{w} = \sum_{i=1}^d u_i \frac{d_i^2}{d_i^2 + \lambda} u_i^T y$ , where  $u_i$  are the left-singular vectors of  $X$ .

**For OLS**, we get predictions  $\hat{y} = \hat{X} \hat{w} = X(X^T X)^{-1} X^T y = \sum_{i=1}^d u_i u_i^T y$ .  $\rightarrow$  In ridge regression, the contributions of the left-singular vectors of  $X$  associated with small singular values are shrunk. In OLS we just project onto the column space of  $\hat{X}$ .

**Remark** (Connection to MAP inference): Consider the model  $y = Xw + \epsilon$ , with  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$  and  $w \sim \mathcal{N}(0, \tau \mathbb{I})$ . Our objective using MAP is  $\hat{w} = \underset{\argmax}{\text{argmax}}_{w \in \mathbb{R}^d} p(w|X, y) := \underset{\argmin}{\text{argmin}}_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + \frac{\sigma^2}{\tau} \|w\|_2^2$ , i.e., with a normal noise and an isotropic normal prior on the parameters, we recover ridge regression. Note that:  $\tau \downarrow \Rightarrow \lambda \uparrow$  (a strong prior belief that  $w$  is close to 0 implies a stronger regularization)

#### 9.3 Lasso Regression

The lasso regression solution for a linear model is given by  $\hat{w} = \underset{\argmin}{\text{argmin}}_{w \in \mathbb{R}^d} \|Xw - y\|_2^2 + \lambda \|w\|_1$ . Consider the  $L^1$ -regularized second-order approx. of an arbitrary loss fct. around an optimal point  $\theta^*$ :

$$\mathcal{R}_{L^1}(\theta) \approx \mathcal{R}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T H (\theta - \theta^*) + \lambda \|\theta\|_1.$$

**Assuming** that  $H = \text{diag}(h_1, \dots, h_d)$ , with  $h_i > 0$ , this is equivalent to  $\mathcal{R}_{L^1}(\theta) \approx \sum_{i=1}^d \left[ \frac{h_i}{2} (\theta_i - \theta_i^*)^2 + \lambda |\theta_i| \right] + \text{const.}$  We find the following expression for those  $\{\theta_i\}_1^d$  that minimize the approximation:

$$\theta_i = \text{sign}(\theta_i^*) \cdot \max\{0, |\theta_i^*| - \frac{\lambda}{h_i}\}.$$
 We distinguish the following cases:

- $0 < \theta_i^* \leq \lambda/h_i \Rightarrow$  optimal  $\theta_i$  is 0: Regularization term overwhelms loss functions, pushing the value of  $\theta_i$  to zero
- $0 < \lambda/h_i < \theta_i^*$ : Regularization shifts the value of  $\theta_i$  by a distance of  $\lambda/h_i$
- $-\lambda/h_i < \theta_i^*$  and  $\theta_i^* < -\lambda/h_i < 0$  are analogous

**Remark:**  $L^2$ -regularization with a diagonal Hessian results in scaling the optimal parameters, i.e.,  $\theta_i = \frac{h_i}{h_i + \lambda} \theta_i^*$ , meaning that when  $\theta_i^*$  is nonzero, the optimal  $\theta_i$  will be nonzero as well.

#### 9.4 Connection between early stopping and $L^2$ -regularization

Let  $\lambda$  be the weight decay factor,  $\eta$  the learning rate, and  $\tau$  the (early) stopping time. Then, for a linear model with a quadratic approximation of the loss function, and assuming  $\lambda_i/\lambda \ll 1$  and  $\eta \lambda_i \ll 1$ , where  $\lambda_i$  s are the eigenvalues of the Hessian of the quadratic approximation, it holds that  $\tau \approx \frac{1}{\lambda}$ .

$\rightarrow$  The number of training iterations is **inversely proportional** to the  $L^2$ -regularization parameter  $\lambda$ . Thus, optimizing for many epochs is similar to enforcing a regularization with a relatively low coefficient and vice versa.

#### 9.5 Batch and Weight Normalization

Batch normalization normalizes the data projected on the weight vector before feeding it into the (sufficiently smooth) function  $l : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}^+$  as follows:  $l(w; g, \gamma) = \mathbb{E}_{x,y} [l(\text{BN}(x; w, \gamma), y)]$ , where

$$\text{BN}(x; w, \gamma, g) = \gamma \frac{x^T w - \frac{\mathbb{E}[x^T w]}{\text{Var}(x^T w)}}{\sqrt{\text{Var}(x^T w)}} + \gamma.$$
 We omit  $\gamma$  w.l.o.g. in the following.

With  $x$  being zero-mean, we have  $\text{Var}(x^T w) = w^T S w = \|w\|_2^2$ , where  $S = \text{Var}(x)$ . Then,  $\text{BN}(x; w, \gamma) = g \frac{x^T w}{\|w\|_2} + \gamma$  and thus  $\text{BN}(x^T w; w, \gamma) = \frac{x^T w}{\|w\|_2} + \gamma$ , with  $\tilde{w}/\|w\|_2 := g \frac{w}{\|w\|_2}$ . This means that under certain assumptions, batch normalization corresponds to a **re-parameterization of the weight space**. On the other hand, **weight** normalization parameterizes the weights as  $w = g \frac{\tilde{w}}{\|\tilde{w}\|_2}$ . Hence, WN is the same as BN (under certain assumptions) with the covariance matrix replaced by the identity matrix.

**Remark:** Both BN and WN are linear scale invariant because both perform a scaling of the weights by a norm.

**Remark:** Benefits of WN over BN include: (i) it is easier to use in inference mode; (ii) it permits distributed optimization because the normalization does not include dependencies between batch samples.

## 10 Theory

#### 10.1 Bayesian Linear Regression:

Let  $y = f(x) = e \cdot f(x)$  be  $\Phi(x) = y$ ,  $\Phi = \{\phi(x_1), \dots, \phi(x_n)\}$ ,  $y| \Phi, w \sim \mathcal{N}(\Phi^T w, \Sigma)$ , and  $w \sim \mathcal{N}(0, \Sigma^{-1})$ . Then,

- posterior:  $w_i | x \sim \mathcal{N}\left(Q^{-1} \frac{\partial \mathcal{L}}{\partial w_i}, Q^{-1}\right)$  with  $Q = \left(\frac{\partial^2 \mathcal{L}}{\partial w^2} + \Sigma^{-1}\right)$
- posterior pred.:  $\phi_*^T w \sim \mathcal{N}\left(\phi_*^T Q^{-1} \frac{\partial \mathcal{L}}{\partial w}, \phi_*^T Q^{-1} \phi_*\right)$ , where  $\phi(x^*) = \phi_*$

Note that for  $\Phi \in \mathbb{R}^{N \times n}$ ,  $N \gg n$  the computation of the matrix inverse above is of order  $\mathcal{O}(N^2)$ . Hence, using Woodburry's identity and the pushthrough identity we can reduce the complexity to  $\mathcal{O}(n^2)$ :  $f_*(x_*) = \mathcal{N}(\phi_*^T \Sigma \Phi (\Phi^T \Sigma \Phi + \sigma^2 \mathbb{I}_n)^{-1} y, \phi_*^T \Sigma \Phi (\phi_*^T \Sigma \Phi + \sigma^2 \mathbb{I}_n + \Phi^T \Sigma \Phi)^{-1} \Sigma^T \phi_*)$ , where the feature representation only enters in form of a kernel  $K(x, x_*) = \phi(x)^T \phi(x_*) \in \mathbb{R}$  with  $\phi(x) = \Sigma^{1/2} \Phi(x)$

**Bayesian Ensembling:** We use relative posterior weighting to average outputs of set of DNNs.

$$f(\Theta)(x) = \sum_{i=1}^n \frac{p(\Theta|S)}{\sum_{i=1}^n p(\Theta|S)} f^{(i)}(x) = \sum_{i=1}^n \frac{\exp(-E(\Theta))}{\sum_{i=1}^n \exp(-E(\Theta))} f^{(i)}(x)$$

The above is feasible even if we do not have direct access to the normalized posterior, we can still use unnormalized posterior via energy function  $E$ . **Markov Chain Monte Carlo:** It is a standard approach to sampling from high-dimensional posterior distributions. We define a Markov chain (random sequence) in parameter space  $\theta^0, \theta^1, \theta^2, \dots$ , s.t.  $\theta^{t+1} \theta^t \sim \Pi$ . We choose next parameter based on previous one w.r.t.  $\Pi$ , which is the Markov kernel of the chain. Posterior is the stationary distribution of the Markov chain. Lastly, one usually needs burn-in period for better estimates.

**Metropolis-Hastings:** This is a way of specifying the kernel. We define detailed balance condition and Markov kernel s.t.

$$p(\theta_1 | S) \Pi(\theta_2 | \theta_1) = p(\theta_2 | S) \Pi(\theta_1 | \theta_2).$$

If this condition holds, the resulting Markov chain is time-reversible and has the desired posterior as its unique stationary distribution

**Remark:** Markov chains can have stationary distributions without fulfilling detailed balance.

In MH sampling, we start from initial  $\tilde{\Pi}$  but then modify the transition probability via an acceptance/rejection step to achieve an effective  $\Pi$  meeting detailed balance. Assume that we always want to accept the transition in one of the directions (increasing probability usually), but possibly reject it with some probability in the opposite direction. Thus, we get

$$p(\theta_1 | S) \Pi(\theta_2 | \theta_1) A(\theta_2 | \theta_1) \stackrel{!}{=} p(\theta_2 | S) \Pi(\theta_1 | \theta_2) A(\theta_1 | \theta_2).$$

$$\begin{aligned} &= \Pi(\theta_1 | \theta_2) &= \Pi(\theta_1 | \theta_2) \\ \text{With the desired one-sided structure of acceptance probabilities } A(\theta_2 | \theta_1) &= 1 \vee A(\theta_1 | \theta_2) = 1 \text{ which leads to the unique choice} \end{aligned}$$

$$A(\theta_1 | \theta_2) = \min \left\{ 1, \frac{p(\theta_1 | S) \Pi(\theta_2 | \theta_1)}{p(\theta_2 | S) \Pi(\theta_1 | \theta_2)} \right\}.$$

**Hamiltonian Monte Carlo:** Consider an energy function equal to negative log posterior  $E(\theta) = -\sum_{x_i} \log p(y_i | x_i \theta) - \log p(\theta)$ . This potential function is augmented by momentum  $v$  and kinetic energy term s.t.  $H(\theta, v) = E(\theta) + \frac{1}{2} v^T M^{-1} v$ . Then, the joint probability distribution is given by Gibbs distribution  $p(\theta, v) \propto \exp[-H(\theta, v)]$ . Hamiltonian dynamics reformulate the Euler-Lagrange equations as  $\dot{v} = -\nabla E(\theta)$ ,  $\dot{\theta} = v$ . In Hamiltonian Monte Carlo (HMC), these are discretized using a step size  $\eta$ :  $\theta^{t+1} = \theta^t + \eta v^t$ ,  $v^{t+1} = v^t - \eta \nabla E(\theta^t)$ . HMC relates to gradient descent with momentum and samples the posterior correctly, but it is slow. To scale HMC for stochastic gradients, Langevin Dynamics is used to make the posterior invariant under stochastic updates.

**Langevin Dynamics:** An important contribution was the introduction of Langevin dynamics with **friction** and **noise process**  $\dot{\theta} = v$ ,  $\dot{v} = -\nabla E(\theta) dt - B \sigma^2 dt + \mathcal{N}(0, 2B \sigma^2 dt)$ . Friction reduces the momentum and dissipates kinetic energy. Wiener noise process injects stochasticity. This SDE can be discretized via symplectic Euler, resulting in (**friction**, **stochastic**, **extra noise**):  $\theta^{t+1} = \theta^t + \eta v^t$ ,  $v^{t+1} = (1 - \eta \gamma) v^t - \eta \sigma^2 \nabla E(\theta^t) + \sqrt{2 \gamma \eta} \mathcal{N}(0, \mathbb{I})$ , where  $E$  is the stochastic potential function, which includes an empirical loss over a random minibatch of data. The friction leads to an exponential damping with time.

#### 10.2 Gaussian Processes

A Gaussian Process is a collection of random variables, any finite number of which have a joint Gaussian distribution. It is completely specified by its mean function  $m(x) = \mathbb{E}[f(x)]$  and covariance function  $k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x')))]$  and we write  $f(x) \sim \mathcal{GP}(m(x), k(x, x'))$ .

**Exa:** A Bayesian linear model  $f(x) = \phi(x)^T w$ ,  $w \sim \mathcal{N}(0, \Sigma)$  is a Gaussian Process since for any  $m_1, \dots, \phi(x_1), \dots, \phi(x_m)$  it holds that  $\sum_{i=1}^m \alpha_i f(x_i) = \sum_{i=1}^m \alpha_i \phi(x_i)^T w = (\sum_{i=1}^m \alpha_i \phi(x_i))^T w$  which, as a linear combination of Gaussians, is Gaussian. Hence,  $f(x) \sim \mathcal{GP}(0, \phi(x)^T \Sigma \phi(x'))$ . **Linear Unit & Layer:** Consider a linear unit with  $n$  inputs and a random Gaussian weight vector  $w \sim \mathcal{N}(0, \frac{\sigma^2}{n} \mathbb{I}_{n \times n})$ . The units outputs for a vector of inputs can be written as  $y_i = w_i^T x_i$  and any linear combination of  $y$  writes as  $\sum_{i=1}^n \alpha_i y_i = w^T \tilde{x}$ ,  $\tilde{x} := \sum_{i=1}^n \alpha_i x_i$ . The pre-activation of a single linear unit follows a Gaussian process due to the properties of normal distributions under linear combinations. Its covariance function is the scaled inner product kernel, given as

$$\mathbb{E}[y_i y_j] = \mathbb{E}[(w \cdot x_i)(w \cdot x_j)] = \mathbb{E}[x_i^T w w^T x_j^T] = x_i^T \left( \frac{\sigma^2}{n} \mathbb{I}_{n \times n} \right) x_j = \frac{\sigma^2}{n} x_i^T x_j.$$

This concept applies to single layers of linear units where activations are independent Gaussian, conditioned on the inputs.

**Deep Layers:** In the next layers there is randomness in weights and also in activations. We observe products between random matrices of the type  $W^{l+1} X^l$ ,  $l \geq 1$ , where  $X^l = X^l(W^{l-1}, \dots, W^1)$  is the activation matrix of layer  $l$  and is random. Thus, such products are not normal, although a deep pre-activation process yields normality due to CLT.

**Non-linear Layer Maps:** A non-linear activation function  $\phi$  reshapes pre-activations which are no longer Gaussian. In the GP limit of pre-activations, the mean and covariance of units are sufficient to work with, since activations of layer  $l$  are summed over to form pre-activations of

$l+1$ . CLT reshapes them back to Gaussian such that the mean functions writes as  $\mathbb{E}[W^{l+1} X^l] = \mathbb{E}[\phi(W^l X^{l-1})]$ . DNNs can be interpreted as Gaussian Processes (GPs) in the infinite width limit, using kernel recursion. Then  $K_{f, \mu}^{l+1} = \mathbb{E}[\phi(x^{l+1}) \phi(x^{l+1})^T] = \sigma^2 \mathbb{E}[\phi(\mu_f) \phi(\mu_f)^T]$ ,  $f \sim \mathcal{GP}(0, K^{l+1})$ .

As for NTK, we hereby use kernel regression and Bayesian predictive mean looks like  $f^*(x) = k(x)^T K^{-1} y$ . Adv. of GP view: We can also quantify the conditional variance:  $\mathbb{E}[(y(x') - f^*(x'))^2] = K(x, x') - k(x)^T K^{-1} k(x)$ .

#### 10.3 Empirical Neural Tangent Kernel

$$K(x, x'; \Theta) = \nabla_\theta f(x; \Theta)^T \nabla_\theta f(x'; \Theta) = \sum_{p=1}^P \frac{\partial f(x; \Theta)}{\partial \theta_p} \frac{\partial f(x'; \Theta)}{\partial \theta_p}$$

Using the NTK, we can track how our function changes in the sample during training (functional gradient flow):  $\dot{f} = K(\Theta)(y - f)$ .

**Exa:**  $f(x) = x^T V x$ ,  $w \in \mathbb{R}^m$ ,  $V \in \mathbb{R}^{m \times d}$

$$\nabla_\theta f(x)^T \nabla_\theta f(x') = (x^T V^T w \otimes x^T) \left( \frac{V x'}{w \otimes x'} \right) = (V x, V x') + (w^T w \otimes x^T x') = (V x, V x') + \|w\|_2^2 (x^T x')$$

**Infinite Width Limit** Consider an MLP with wide sequence  $m_l$  and weights initialized as

$$w_{ij}^l = \frac{\sigma_w}{\sqrt{m_l}} w_{ij}^l, \quad b_i^l = \frac{\sigma_b}{\sqrt{m_l}} b_i^l, \quad w_{ij}^{l+1} \sim \mathcal{N}(0, 1).$$

This scaling normalizes weight magnitudes by the square root of layer width (similar to LeCun initialization). In the infinite width limit where  $m_l \rightarrow \infty$  for all hidden layers, under suitable technical conditions (including Lipschitz activation functions), the initial Neural Tangent Kernel (NTK) converges in probability to a deterministic limit  $k(\theta) \xrightarrow{P} k^{\infty}$ . The limiting kernel  $k^{\infty}$  depends only on the initialization distribution, not the specific random parameters. This effectively means there is only one infinite width network at initialization. This result can be intuitively understood through the law of large numbers: as width approaches infinity, the sampled population of units in each layer approaches the population distribution.

**NTK Constancy** Convergence to deterministic NTK  $k^{\infty}$  cannot only be established pointwise, but with a few additional technical assumptions, convergence also holds uniformly over training trajectories s.t.  $k(\theta(t)) \rightarrow k^{\infty}$  uniformly over  $t \in [0, T]$ . **Note:** Activation functions need to have bounded second derivatives and the time-integral of the norms of update directions needs to remain stochastically bounded. The latter assumption can be verified for gradient flow trajectories.

This remarkable property can also be stated (non-rigorously) as follows: the initial NTK stays constant under gradient flow  $\frac{d k(\theta(t))}{dt} = 0$ .

Under NTK constancy, learning in the infinite width limit becomes equivalent to the linearized model, with the solution s.t.  $f^{(\infty)}(x) = k(x) K^{-1} (y - f)$ ,  $k = k^{\infty}$

The emergence of NTK constancy can be explained through vanishing curvature in the infinite width limit, where the Hessian operator norm scales as  $\left| \frac{\nabla^2 f(\theta_0)}{\|\nabla f(\theta_0)\|_2} \right| \ll 1$ . This property relates to the lazy training regime, where DNN training becomes a convex problem in the infinite width limit. For empirical NTK (gram matrices), near-constancy can be established as  $|k(\theta_0) - k(\theta_t)|_F^2 \in \mathcal{O}(1/m)$ ,  $m = m_1 = \dots = m_L$ . This holds under similar assumptions including bounded input domain  $\mathcal{X} = \{x : |x|_2 \leq 1$  and appropriate learning rate conditions.

#### 10.4 Statistical Learning Theory

**VC Learning Theory:** Possible classific. outcomes of a fct. class on sample  $S = \{x_1, \dots, x_n\}$ :  $\mathcal{F}(S) = \{(f(x_1), \dots, f(x_n)) : f \in \mathcal{F}, f \in \mathcal{F}\}$ . VC dimension is the largest possible sample size for which the function class  $\mathcal{F}$  is still able to shatter a sample. Let  $\mathcal{E}(f)$  be the empirical error and  $E(f)$  expected error. The uniform convergence yields VC inequality  $\mathbb{P}(\sup_{f \in \mathcal{F}} |\mathcal{E}(f) - E(f)| > \epsilon) \leq 8 |\mathcal{F}(S)| e^{-\epsilon^2/32}$ , yielding a non-vacuous bound on generalization error  $\pi$  as fixed input (no MCMC needed)  $\rightarrow$  unbiased estimate of gradient and small  $\uparrow$  var

**Optimizing q:** Use stoich. backup (1) reparam. parametrized distr. by fixed unpara. randomness  $\epsilon$ ,  $z \sim \mathcal{N}(\mu, \Sigma) \Leftrightarrow z = \mu + \Sigma^{\frac{1}{2}} \eta$ ,  $\eta \sim \mathcal{N}(0, \mathbb{I})$

**Approach** generiz. for any smooth and integrable  $f : \mathbb{R}^d \mapsto \mathbb{R}$   $\mathbb{E}[f(z)] = \mathbb{E}[V_z f(z)]$ ,  $V_z \mathbb{E}[f(z)] = \frac{1}{2} \mathbb{E}[V_z^2 f(z)]$

#### Randomization Experiments:

- DNNs can perfectly fit training data
- DNNs can possibly memorize training data perfectly (even when true labels are replaced with random labels), and thus learn nothing and cannot generalize.
- Transition from true to random labels increase train time only by small factor. Gradient-based methods easily yields memorization solutions.
- These exist, even when randomly shuffling the pixels of a DNN. DNNs generalize despite infinite capacity, which defies traditional learning theory. They can memorize random data, and ConvNets' inductive bias, like shift invariance, offers a little advantage, showing that these phenomena are not data-dependent.
- Double Descent:** There is an interpolation threshold at which the DNN "memorizes" training examples with little generaliz. But beyond this point, even larger models start to learn and may eventually reach lower risk. This is called the *over-parametrized regime*, which is not harmed by over-fitting.
- PAC Bayesian Bounds:** We consider a stochastic classifier  $f$  and aim to bound the expected generalization gap. The generalization gap is defined as  $\max(\mathcal{E}_0 - E)$ , where  $\mathcal{E}_0$  is the empirical error w.r.t. the training sample, and  $E$  is expected error w.r.t. the data distribution. We make use of the change of measure inequality.

**Theorem:** For any  $P \geq Q$  and  $P$ -measurable  $f$  (i.e., a random variable),

$$\mathbb{E}_Q[\log P] \leq \text{KL}(Q|P) + \ln \mathbb{E}_P[e^{\phi}]. \text{ Proof: Via Jensen's inequality,}$$

$$\ln \mathbb{E}_P[e^{\phi(f)}] \geq \ln \mathbb{E}_Q[e^{\phi(f)} \frac{p(f)}{q(f)}] \geq \mathbb{E}_Q[\ln e^{\phi(f)} \frac{p(f)}{q(f)}] = \mathbb{E}_Q[\phi(f)] - \mathbb{E}_Q[\ln Q(f) - \ln P(f)], \text{ where the last term is equivalent to } \text{KL}(Q|P).$$

We state famous **PAC-Bayesian Theorem**, ensuring a general rate of  $\mathcal{O}(1/\sqrt{n})$  and there are no hidden constants. This bound is as simple as it gets, but it only bounds stochastic classifier, not a single classifier.

**PAC-Bayesian for DNNs:** Define a Gaussian over parameters  $P = \mathcal{N}(\theta_0, \Lambda^2 I)$ . The setting  $\theta_0 = 0$  and  $\Lambda = 1$  yields a standard Gaussian s.t.  $Q = \mathcal{N}(\theta_0, \text{diag}(\theta_0^2))$ . With  $P$  and  $Q$  being Gaussians with diagonal covariance, one can derive  $\text{KL}(Q|P) : \text{KL}(Q|P) = \sum_i \log \frac{\sigma_i}{\sigma_i} + \frac{\sigma_i^2 - \sigma_i^0}{2 \sigma_i^2} = -\frac{1}{2}$

In the limit  $\Lambda \rightarrow \infty$ ,  $\text{KL}(Q|P) \rightarrow -\sum_i \log \sigma_i$ , favoring large variances for the  $Q$ -ensemble. The PAC-Bayes bound to be minimized is:

$$E_{\text{PAC}}(Q) := \mathbb{E}_Q[\tilde{E}] + \left[ \frac{\sqrt{2}}{2} \left( \text{KL}(Q|P) + \ln \left( \frac{2 \sqrt{e}}{\epsilon} \right) \right) \right]$$

A practical implementation involves generating parameter sequence  $\theta^i$  via SGD, evaluating gradients on perturbed parameters sampled from the  $Q$ -ensemble:  $\theta = \theta - \eta \nabla \mathbb{E}_Q[\tilde{E}]$ ,  $\tilde{\theta} \sim Q(\theta, \sigma)$  The reparameterization trick enables backpropagation:  $\tilde{\theta} = \theta + \text{diag}(\sigma_i) \eta$ ,  $\eta \sim \mathcal{N}(0, \mathbb{I})$  This approach balances between finding parameters  $\theta^i$  that achieve low expected loss and maintaining robustness through the variance in  $Q$ , which affects the KL-divergence term in the bound.

**Wide Local Minima:** Flatness of local minima with better generalization.  $\rightarrow$  Argued that SGD with small mini-batches leads to flatter minima comp. to full batch gradient descent. Also, weight avg. found to be favorably towards finding flatter minima.  $\rightarrow$  entropy SGD uses Langevin dynamics to favor minima with higher entropy.

## 11 Generative Models

#### 11.1 VAEs

- AE:  $x \mapsto z \in \mathbb{C} x$ ,  $\mathbb{C} \in \mathbb{R}^{m \times m}$ , used with  $\mathcal{E}(C, D)$  as  $\frac{1}{2} \|x - DCx\|_2^2$ 
  - Goal: find good embedding

- **GANs vs. Diffusion**
- **Training Obj:** Diffusion: minimize reconstr. loss for noise pred. or data reconstr. GANs: Advers. setup
- **Training Setup:** Diffusion: Stable due to non-advers. setup (requires many training steps, hyperparam tuning), GANs: prone to instability
- **Output Quality:** Diversity may suffer in GANs due to model collapse
- **Gen Process:** Iterative process starting from noise to output, GAN: Single forward pass to create sample
- **Comp cost:** Computationally expensive due to iterative sampling, GANs: Efficient during inference

## 12 Advanced Topics

### 12.1 Theory

**p-Norm Robustness:** For a multi-class classifier  $f$ , the goal of adversarial example perturbation is to find a perturbation  $\eta$  given a pattern  $x$  s.t.  $f(x+\eta) \neq f(x)$  s.t.  $\|\eta\|_p \leq \epsilon$

For  $p = 2$  and a linear model with weight vectors  $w_i$ , we have  $f(x) = \arg\max_i f_i(x)$ ,  $f_i = w_i^T x + b_i$ . For binary classification, the optimal perturbation is then given by  $\eta \propto \text{sign}(f_1(x) - f_2(x))(w_2 - w_1)$ . This generalizes to  $m > 2$  classes ( $f(x) = 1$  w.l.o.g.) by finding most easily confusable class:  $\eta = \arg\min_{i \neq j} \frac{f_i(x) - f_j(x)}{\|w_i - w_j\|_2} (w_j - w_i)$

**Linearization:** Suggested in DeepFool: leverage the linear case to *linearize* a nonlinear model and iteratively find an adversarial perturbation. For  $m = 2$  one iterates solving the problem  $\arg\min_{\Delta y} \|\Delta y\|_2$  s.t.  $(\nabla f_1(x) - \nabla f_2(x))^T \Delta y < f_1(x) - f_2(x)$ , with  $f_1(x) - f_2(x)$  being the margin. This can be generalized to  $p$ -norms.

**Robust Training:** Loss is extended to neighborhoods of training points:  $\ell(f(x), y) \mapsto \max_{\eta: \|\eta\|_p \leq \epsilon} \ell(f(x + \eta), y)$  which yields a minimax problem, where adversary picks worst perturbation  $\eta$  an learner picks best parameter  $\theta$  in response. One can solve inner loop maximization via *projected gradient ascent*. For  $p = 2$ :  $\eta^{t+1} = \epsilon \Pi[\eta^t + \alpha \nabla_{\eta} \ell(f(x + \eta^t), y)]$ ,  $\|\cdot\|_2 \leq \frac{\epsilon}{\sqrt{2}}$ , for  $p = \infty$ :  $\eta^{t+1} = \epsilon \Pi[\eta^t + \alpha \text{sign}(\nabla_{\eta} \ell(f(x + \eta^t), y))]$ ,  $\|\cdot\|_2 \leq \frac{\epsilon}{\sqrt{2}}$ .

A practical method is Fast Gradient Sign Method (FGSM), which performs one iteration of  $\ell_1$ -*l*<sub>∞</sub>-projected gradient descent:  $\eta = \epsilon \text{sign}(\nabla_{\eta} \ell(f(x), y))$ , modifying the pixels by  $\pm \epsilon$  depending on the sign of  $\frac{\partial \ell}{\partial x_i}$  per pixel. FGSM results in  $\approx 2\times$  overhead compared to standard training.

### 12.2 Adversarial Attacks

Aspect	White-Box Attacks	Black-Box Attack
Access to model	Full access (parameters, gradients)	Limited access (input/output only)
Gradient Information	Direct access	Estimated or approximated
Attack methods	FGSM, PGD, CW, etc.	Transferability, query-based, cost-based
Computational cost	Lower (direct gradient usage)	Higher (requires queries or surrogate training)
Strength	Very strong	Relatively weaker, but still effective

**Grey-box:** partial knowledge of target model, which could include some knowledge such as model architecture, training data, setup and outputs, but does not have access to the exact model parameters (weights and biases) or gradients.

**Targeted vs Untargeted:**

- **Untargeted:** aim to fool model such that output is anything different from the correct label
- **Targeted:** aim to fool model to return target label of the attacker's interest

**Types of Attacks:**

- **Gradient-based:** for white-box setting by solving either targeted or untargeted adversarial objective (e.g. FGSM and PGD)
- **Optimization-based:** achieve adversarial objective while minimize perturbation size (e.g. Carlini & Wagner attack)
- **Transferability-Based:** usually for black-box and grey-box; trains surrogate  $\Rightarrow$  generate adversarial examples based on the surrogate using optimization-based methods  $\Rightarrow$  applied to target model
- **Query-based:** Usually for the black-box; if attacker can query the black-box model and observe its outputs (predictions or confidence scores), they can iteratively approximate gradients using optimization-based methods  $\Rightarrow$  applied to target model

Gradients can be estimated using *finite differences*:

Aspect	FGSM	PGD	C&W
Attack type	Single-step	Iterative (multi-step)	Optimization-based
Perturb. constraint	$L_\infty$ -bounded	$L_\infty$ -bounded (often $L_2$ )	Flexible ( $L_\infty, L_1, L_2$ )
Objective	Maximize loss (fast)	Maximize loss iteratively	Minimize perturb.-mistakes
Strength	Weak	Stronger than FGSM	Very strong
Comp. Cost	Low	Medium (multiple iters)	High (complex optim. Problem)
Use case	Quick evaluation	Strong adversarial training	Advanced adversarial attacks

**FGSM:**  $x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y))$ , where  $x$  : original input,  $\epsilon$  : perturbation magnitude (small scalar value),  $\nabla_x J(\theta, x, y)$  gradient of loss; FGSM is a special case of PGD with a single step

**PGD:**  $x_{adv}^{(t+1)} = \text{Proj}_{B_\epsilon(x)}(x_{adv}^{(t)} + \alpha \cdot \text{sign}(\nabla_x J(\theta, x_{adv}^{(t)}, y)))$ , where  $\text{Proj}_{B_\epsilon(x)}$  projects adversarial examples back into the  $\epsilon$ -ball centered at  $x$ .

**Carlini & Wagner:** minimize  $\|x_{adv} - x\|_p + c \cdot f(x_{adv})$ , where  $\|x_{adv} - x\|_p$  is perturbation size e.g.  $L_2$  norm,  $c$ : regularization constant balancing perturbation size and attack success,  $f(x_{adv})$ : function that ensures  $x_{adv}$  is misclassified (often  $f(x_{adv}) \leq 0$  if attack succeeds). Then use optimization algorithm such as SGD or Adam to search for optimal  $x_{adv}$

#### 12.3 Other Adversarial Attacks

- **Backdoor Attacks:** *training-time* attack where an adversary injects hidden "backdoor" patterns into the training dataset. The model behaves normally on clean inputs but produces targeted outputs when presented with inputs containing the backdoor trigger
- **Privacy Inference Attacks:** *test-time* attack where an adversary exploits the model's predictions to infer sensitive information about the training data e.g. *Membership Inference* (determining whether data point was used during model training), *Attribute Inference* (recovering specific attributes of the input data e.g. gender, age, private features)
- **Model Extraction Attacks:** adversary attempts to replicate or approximate the target model by querying it repeatedly and analyzing the outputs to reconstruct the model's behavior

- **Model Inversion Attacks:** *test-time* attack where adversary uses the model's output predictions to reconstruct the input features, potentially revealing private or sensitive data

#### 12.4 Adversarial Defense

- **Denoising based methods:** attempt to use conventional image/signal processing techniques, GANs, autoencoders, or some generative approaches to erase the adversarial perturbations placed on the input samples
- **Randomization based methods:** based on random transformation or regularization with motivation that randomizing adversarial effects introduced in adversarial examples could potentially make DNNs more robust to them
- **Adversarial training based methods:** aim to improve robustness of DNNs by training them with adversarial examples. Motivation: adversarial examples can be considered scarce data residing near the true decision boundary. Including such samples in the training process will enable the models to learn more specific decision boundaries that take these outliers into consideration

Denoising-based easily bypassed by optim.-based approaches; randomization-based have been proven to be exploiting *obfuscated gradient* which gives a false sense of security and can be easily compromised by adv. examples generated through Expectation over Transformation (EoT). Adversarial training has been the most robust defense method

## 13 Miscellaneous

#### 13.1 Activation Functions and Important Derivatives

- $\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(2x) - 1$ ;  $\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{1 + \tanh(x/2)}{2}$
- $\tanh'(x) = 1 - \tanh^2(x)$ ,  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$
- $\text{GELU}(x) = x\Phi(x) \Rightarrow \text{GELU}'(x) = \frac{dx}{dx} \exp(-\frac{(\text{log}x)^2}{2}) + \Phi(ax)$

#### 13.2 Trigonometric Functions

- $\sin(\alpha \pm \beta) = \sin \alpha \cos \beta \pm \cos \alpha \sin \beta$ ;  $\cos(\alpha \pm \beta) = \cos \alpha \cos \beta \mp \sin \alpha \sin \beta$
- $\tan(x) = \frac{\sin(x)}{\cos(x)}$ ;  $\frac{\partial \tan(x)}{\partial x} = \frac{1}{\cos^2(x)} = \sec^2(x)$
- $\sinh(x) = \frac{e^x - e^{-x}}{2}$ ;  $\frac{\partial \sinh(x)}{\partial x} = \cosh(x)$
- $\cosh(x) = \frac{e^x + e^{-x}}{2}$ ;  $\frac{\partial \cosh(x)}{\partial x} = \sinh(x)$
- $\tanh(x) = \frac{\sinh(x)}{\cosh(x)}$ ;  $1 = \sinh^2(x) + \cosh^2(x)$

#### 13.3 Analysis

- $\frac{\partial [x-b]_+}{\partial x} = \frac{x-b}{\|x-b\|_2}$ ;  $\frac{\partial [x-b]}{\partial x} = \frac{x-b}{\|x-b\|_2}$ ;  $\frac{\partial b^T x}{\partial x} = \frac{\partial x^T b}{\partial x} = b$
- $\frac{\partial b^T A x}{\partial x} = A^T b$ ;  $\frac{\partial x^T A}{\partial x} = 2x$ ;  $\frac{\partial x^T A x}{\partial x} = (A^T + A) x$ ,  $A$  symm.  $2Ax$
- $\frac{\partial \log|X|}{\partial X} = X^{-T}$ ;  $\frac{\partial \text{tr}(X^T A)}{\partial X} = A$ ;  $\frac{\partial \text{tr}(X^T X)}{\partial X} = \frac{\partial \text{tr}(XX^T)}{\partial X} = 2X$
- $\frac{\partial x^T A}{\partial A} = xx^T$ ;  $\frac{\partial \text{ReLU}(b^T x)}{\partial x} = \mathbb{I}(b^T x > 0)$ ;  $\frac{\partial \text{softmax}(z_i)}{\partial z_i} = \text{softmax}(z_i)(1 - \text{softmax}(z_i))$
- $\frac{\partial A_1 B A_2}{\partial A} = (A_1 \otimes A_2^T) \frac{\partial B}{\partial A}$ ;  $(A_1 \otimes B_1)(A_2 \otimes B_2) = A_1 A_2 \otimes B_1 B_2$
- $f(x) = w^T V x$ ,  $\frac{\partial f(x)}{\partial V} = w^T \otimes x^T$ ;  $\nabla_V f(x) = w \otimes x$
- $\frac{\partial x \otimes y}{\partial x} = \text{diag}(y)$ ;  $\frac{\partial x \otimes y}{\partial y} = \text{diag}(x)$
- $\sum_{i=0}^n \binom{n}{i} = 2^n$

#### 13.4 Probability Theory

- Markov:  $\phi(\cdot)$  non-decr.& non-neg.  $\leadsto \mathbb{P}(X \geq e) \leq \frac{\mathbb{E}[\phi(X)]}{\phi(e)}$
- Chebyshev:  $\phi(\cdot)$  non-decr.& non-neg.  $\leadsto \mathbb{P}(f(X) \geq e) \leq \frac{\mathbb{E}[\phi(f(X))]}{\phi(e)}$

**Normal Distribution:**

- $p(z) \propto \exp(-\frac{1}{2} z^T Q z + z^T m) \Rightarrow p(z) = \mathcal{N}(z; Q^{-1} m, Q^{-1})$
- $X_1 | X_2 = a \sim \mathcal{N}(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (a - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21})$

#### 13.5 Linear algebra

**Cauchy-Schwarz inequality:**  $\langle x, y \rangle^2 \leq \langle x, x \rangle \langle y, y \rangle$ . With dot prod. as inner prod. and  $\|x\| = \sqrt{\langle x, x \rangle}$ , this means  $|x \cdot y| \leq \|x\|_2 \cdot \|y\|_2 \iff |x \cdot y| \leq \|x\|_2 \|y\|_2$

**Spectral norm definition & properties:**

- $\|A\|_2 = \max_{x: \|x\|=1} \|Ax\| = \sigma_1(A)$ , where  $\sigma_1$  is  $A$ 's largest singular value
- $\|AB\|_2 \leq \|A\|_2 \|B\|_2$ . Also extends to higher order products

**Woodbury identity:**

$(I + UC^T V)^{-1} = I - U(C^{-1} + VU)^{-1} V$ ,  $C$  invertible  
 $(A + UC^T V)^{-1} = A^{-1} - A^{-1} U(C^{-1} + VA^{-1} U)^{-1} VA^{-1}$ ,  $A, C$  invertible

**Pseudoinverse:**

$W_{left}^\dagger = \lim_{\sigma^2 \rightarrow 0} W^T (WW^T + \sigma^2 I)^{-1} W$  and  $W W_{left}^\dagger = I$   
 $W_{right}^\dagger = \lim_{\sigma^2 \rightarrow 0} (W^T W + \sigma^2 I)^{-1} W^T$  and  $W_{right}^\dagger W = I$

**Note** that the identities depending on left or right inverse will have different dimensions!!

**Other Matrix identities:**

$(I + AB)^{-1} A = A(I + BA)^{-1}$  and  $(I + AB)^{-1} = I - A(I + BA)^{-1} B$

#### 13.6 Group Theory

A group is a non-empty set  $G$  with a binary operation on  $G$ , denoted  $\cdot$ , such that the following group axioms are satisfied:

- **(Closure):**  $\forall a, b \in G : a \cdot b \in G$
- **(Associativity):**  $\forall a, b, c \in G : (a \cdot b) \cdot c = a \cdot (b \cdot c)$
- **(Identity Element):**  $\exists e \in G \ \forall a \in G : e \cdot a = a$  and  $a \cdot e = a$ . The element  $e$  is called the *identity element*
- **(Inverse Element):**  $\forall a \in G \ \exists b \in G : a \cdot b = e$  and  $b \cdot a = e$ , where  $e$  is the identity element. Note that for all  $a$ , the *inv. element*  $b$  is unique. The group is called **Abelian** if on top of the group axioms, commutativity holds as well, i.e.,  $\forall a, b \in G : a \cdot b = b \cdot a$

## 14 Coding

```

Series 1
# np.sign maps 0 to 0 and is therefore incorrect
y_hat = np.int32(np.inner(X_train[idx, :], theta) > 0) - 1
# perform Perceptron update only if wrongly classified
if y_hat * y_train[idx] < 0:
    theta += self.lr * y_train[idx] * X_train[idx, :]
# compute accuracy
np.mean(np.sign(np.inner(X_train, theta)) == y_train)
```

```

Series 10
class PositionalEncoding(nn.Module):
    def __init__(self, d_model, dropout, max_len=5000):
        super().__init__()
        self.dropout = nn.Dropout(p=dropout)

        position = torch.arange(max_len).unsqueeze(1)
        div_term = torch.exp(torch.arange(0, d_model, 2) *
                              (-math.log(10000.0) / d_model))
        pe = torch.zeros(max_len, 1, d_model)
        pe[:, 0, 0::2] = torch.sin(position * div_term)
        pe[:, 0, 1::2] = torch.cos(position * div_term)
        self.register_buffer('pe', pe)

    def forward(self, x: Tensor) -> Tensor:
        """x: [seq_len, batch_size, embedding_dim]"""
        x = x + self.pe[:, x.size(0)]
        return self.dropout(x)
```

```

class SelfAttention(nn.Module):
    def __init__(self, d: int, heads: int=8):
        super().__init__()
        self.k, self.h = d, heads

        self.Wq = nn.Linear(d, d * heads, bias=False)
        self.Wk = nn.Linear(d, d * heads, bias=False)
        self.Wv = nn.Linear(d, d * heads, bias=False)

        self.unifyheads = nn.Linear(heads * d, d)
```

```

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        b, l, d = x.size()
        h, l_half = h, l // 2

        q = self.Wq(x).view(b, l, h, d).transpose(1, 2)
        .contiguous().view(b * h, l, d)
        k = self.Wk(x).view(b, l, h, d).transpose(1, 2)
        .contiguous().view(b * h, l, d)
        v = self.Wv(x).view(b, l, h, d).transpose(1, 2)
        .contiguous().view(b * h, l, d)

        # (b * h, 1, l)
        w_prime = torch.bmm(q, k.transpose(1, 2)) \
            / np.sqrt(d)
        w = F.softmax(w_prime, dim=-1)
        out = torch.bmm(w, v).view(b, h, l, d)
        out = out.transpose(1, 2).contiguous()
        return self.unifyheads(out.view(b, l, h * d))
```

```

class TransformerBlock(nn.Module):
    def __init__(self, d, heads=8, n_mlp=4):
        super().__init__()
        self.attention = SelfAttention(d, heads)

        self.norm1 = nn.LayerNorm(d)
        self.norm2 = nn.LayerNorm(d)

        self.ff = nn.Sequential(
            nn.Linear(d, n_mlp * d), nn.ReLU(),
            nn.Linear(n_mlp * d, d))
```

```

### Transformer forward: ###
self.token_emb = nn.Embedding(num_tokens, embed_dim)
self.pos_emb = nn.Embedding(max_seq_len, embed_dim)
token_embs = self.token_emb(x)
pos_embs = self.pos_emb(torch.arange(1).to(DEVICE))
.expand(b, 1, d)

final_embs = token_embs + pos_embs
out = self.transformer_blocks(final_embs)
# Compute the mean latent vector for each sequence.
# The mean is applied over dim=1 (time). Shape: [b, d].
out = out.mean(dim=1)
# Classify. Shape: [b, num_classes]
out = self.classification(out)
```

```

### RNN forward: ###
token_embs = self.token_emb(x)
# Feed the embedding into the GRU. Shape: [b, 1, d]
# Use output of the last token as the encoding.
out, final_state = self.rnn(token_embs)
out = out[:, -1]
out = self.classification(out)
```

```

Series 11
class GAN(nn.Module):
    def __init__(self):
        super().__init__()
        # Generator
        self.generator = nn.Sequential(
            nn.ConvTranspose2d(in_channels=latent_dim,
                               out_channels=num_maps_gen*8,
                               kernel_size=4, stride=1, padding=0,
                               bias=False),
            nn.BatchNorm2d(num_maps_gen * 8),
            nn.ReLU(True),
            nn.ConvTranspose2d(
                in_channels = num_maps_gen,
                out_channels = num_maps_gen*4,
                kernel_size=4, stride = 2,
                padding = 1, bias = False),
            nn.BatchNorm2d(num_maps_gen*4),
            nn.ReLU(True),
            [...]
            nn.ConvTranspose2d(
                in_channels = num_maps_gen,
                out_channels = 1,
                kernel_size=4, stride = 2,
                padding = 1, bias = False),
            nn.Tanh()
        )
        # Discriminator
        self.discriminator = nn.Sequential(
            nn.Conv2d(in_channels = 1,
                      out_channels = num_maps_gen,
                      kernels_size = 4, stride = 2,
                      padding = 1, bias = False),
            nn.BatchNorm2d(num_maps_gen*4),
            nn.ReLU(True),
            [...]
            nn.ConvTranspose2d(
                in_channels = num_maps_gen,
                out_channels = 1,
                kernel_size=4, stride = 2,
                padding = 1, bias = False),
            nn.Tanh()
        )
```

```

nn.LeakyReLU(0.2, inplace=True),
[...])
nn.Conv2d(
    in_channels = num_maps_gen*8,
    out_channels = 1,
    kernels_size = 3, stride = 1,
    padding=0, bias = False),
    nn.Sigmoid())

#GAN forward
def gen_forward(self, z):
    img = self.generator(z)
    return self.generator(z)
def disc_forward(self, img):
    pred = self.discriminator(img)
    return pred.view(-1)

gen_optim = torch.optim.Adam( #Adam Gen
    params=model.generator.parameters(),
    lr=generator_lr, betas=(0.5, 0.999))
disc_optim = torch.optim.Adam( #Adam Gen
    params=model.discriminator.parameters(),
    lr=discriminator_lr, betas=(0.5, 0.999))
disc_loss_func = nn.BCELoss()
gen_loss_func = nn.BCELoss()

#GAN Training
for epoch in range(num_epochs):
    model = model.train()
    for batch_idx, (real_data, targets)
        in enumerate(train_loader):
            batch_size = targets.size(0)
            real_data = real_data.to(device)
            targets = targets.to(device)
            # CREATING GROUND-TRUTH
            y_real=torch.ones(batch_size).float().to(d)
            y_fake=torch.zeros(batch_size).float().to(d)
            # TRAIN DISCRIMINATOR
            disc_optim.zero_grad()
            # train the disc.to classify real images
            disc_pred_real=model.disc_forward(real_data)
            .view(-1)
            real_loss = disc_loss_func(disc_pred_real,
                                         y_real)
            # gen imgs frm smpls of the ltnr prior
            z = torch.randn(batch_size, latent_dim, 1, 1,
                             device=device)
            fake_data = model.gen_forward(z)
            # train the disc. to classify fake images
            # detach the comp. graph of the generator
            # s.t. grads are not backprop into the gen)
            disc_pred_fake = model.disc_forward(
                fake_data.detach()).view(-1)
            fake_loss = disc_loss_func(disc_pred_fake,
                                         y_fake)
            disc_loss = 0.5 * (real_loss + fake_loss) # avg
            disc_loss.backward()
            disc_optim.step()
            # TRAIN GENERATOR
            gen_optim.zero_grad()
            # train gen s.t. image is class, as real
            disc_pred_fake = model.disc_forward(
                fake_data).view(-1)
            gen_loss = gen_loss_func(disc_pred_fake, y_real)
            gen_loss.backward()
            gen_optim.step()
            [...]

Series 12
class ForwardProcess: # FIND X_t and noise
    def __init__(self, betas: torch.Tensor):
        self.betas = betas
        self.alphas = 1. - betas
        self.alpha_bar=torch.cumprod(self.alphas,
                                     dim=1)

    def get_x_t(self, x_0: torch.Tensor, t:
               torch.LongTensor) -> Tuple(torch.Tensor, torch.Tensc):
        eps_0 = torch.randn_like(x_0).to(x_0)
        alpha_bar = self.alpha_bar[t-1, None]
        mean = (alpha_bar ** 0.5) * x_0
        std = ((1. - alpha_bar) ** 0.5)
        return (eps_0, mean + std * eps_0)
```

```

class NoiseNetwork1(nn.Module):
    def __init__(self, T):
        super().__init__()
        self.T = T
        self.t_encoder = nn.Linear(T, 1)
        self.model = nn.Sequential(
            nn.Linear(2 + 1, 100),
            nn.LeakyReLU(inplace=True),
            nn.Linear(100, 2),
        )

    def forward(self, x_t, t):
        t_embedding = self.t_encoder(
            F.one_hot(t - 1, num_classes=self.T)
        ).to(torch.float)

        input = torch.cat([x_t, t_embedding], dim=1)
        return self.model(input)

T = 100
beta_start = 0.004
beta_end = 0.02
betas = torch.linspace(beta_start, beta_end, T)

#Training
fp = ForwardProcess(betas)
model = NoiseNetwork(T=T)
optimizer = torch.optim.AdamW(params=model.parameters(),
                                lr=1e-2, betas=(0.9, 0.999), weight_decay=1e-4)
```

```

N = X.shape[0]
for epoch in range(5000):
    with torch.no_grad():
        t = torch.randint(low=1, high=T, size=(N,))
        eps_0, x_t = fp.get_x_t(X, t=t)
        pred_eps = model(x_t, t)
        loss = torch.nn.functional.mse_loss(pred_eps, eps_0)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

#Define reverse process
class ReverseProcess(ForwardProcess):
    def __init__(self, betas: torch.Tensor,
```

```

model: nn.Module):
    super().__init__(betas)
    self.model = model
    self.T = len(betas)

    self.sigma = (
        betas * (1 - torch.roll(self.alpha_bar, 1))
        / (1 - self.alpha_bar)) ** 0.5
    self.sigma[0] = 0.

    def get_x_t_minus_one(self, x_t: torch.Tensor,
                          t: int) -> torch.Tensor:
        with torch.no_grad():
            t_vector = torch.full(size=(x_t.shape[0],),
                                   fill_value=t, dtype=torch.long)
            eps = self.model(x_t, t=t_vector)
            eps *= (1 - self.alphas[t-1])
            / ((1 - self.alpha_bar[t-1]) ** 0.5)
            mean = 1 / (self.alphas[t-1] ** 0.5)
            (x_t - eps)
            return mean + self.sigma[t-1]
                * torch.randn_like(x_t)

    def sample(self, n_samples=1, full_trajectory=False,
               x_t=torch.randn(n_samples, 2)
               trajectory = [x_t.clone()]
               for t in range(self.T, 0, -1):
                   x_t = self.get_x_t_minus_one(x_t, t=t)
                   if full_trajectory:
                       trajectory.append(x_t.clone())
               return torch.stack(trajectory, dim=0) if
               full_trajectory else x_t

Series 13
```

```

def fgsm(model, x, target, eps, targeted=True,
         clip_min=None, clip_max=None):
    # copy of the input Tensor
    input = x.clone().detach_()
    # We have to be able to differentiate
    input.requires_grad_()

    logits = model(input_)
    target = torch.LongTensor([target])
    model.zero_grad_()
    loss = nn.CrossEntropyLoss()(logits, target)
    loss.backward()
    #perform either targeted or untargeted attack
    if targeted:
        out = input_ - eps * input_.grad.sign()
    else:
        out = input_ + eps * input_.grad.sign()
    if (clip_min is not None) or (clip_max is not None):
        out.clamp_(min=clip_min, max=clip_max)
    return out
```

```

def fgsm_targeted(model, x, target, eps,**kwargs):
    return fgsm(model, x, target, eps, targeted=True,
                **kwargs)
def fgsm_untargeted(model, x, label, eps, **kwargs):
    return fgsm(model, x, label, eps, targeted=False,
                **kwargs)
```

```

# x: input image
# label: current label of x
# k: number of FGSM iterations
# eps: size of 1-infinity ball
# eps_step: step size of FGSM iterations
# Each FGSM iteration is projected back to the
# eps-sized 1-inf ball around x
def pgd(model, x, label, k, eps, eps_step, targeted,
        clip_min, clip_max):
    x_min = x - eps
    x_max = x + eps
```

```

# Randomize the starting point x.
x_adv = x + eps * (2 * torch.rand_like(x) - 1)
# Clamp back
if (clip_min is not None) or (clip_max is not None):
    x_adv.clamp_(min=clip_min, max=clip_max)
for i in range(k):
    # FGSM step
    # We don't clamp here (clip_min=clip_max=None)
    # as we want to apply the attack as defined
    x_adv = fgsm(model, x_adv, label, eps_step,
                  targeted)
    # Projection Step
    x_adv = torch.min(x_max, torch.max(x_min, x_adv))
    #if desired: clip the output back to image domain
    if (clip_min is not None) or (clip_max is not None):
        x_adv.clamp_(min=clip_min, max=clip_max)
    return x_adv
```