

Background

$\log(x+1) \approx x$ if $|x| \ll 1$
Variance and covariance:
 $\text{Var}(X) = \text{Cov}(X, X) = \mathbb{E}[(X - \mathbb{E}X)^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$
 $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)^T] = \mathbb{E}[XY^T] - \mathbb{E}[X]\mathbb{E}[Y]^T$
 $\text{Cov}(aX + bY, cW + dV) = ac\text{Cov}(X, W) + ad\text{Cov}(X, V) + bc\text{Cov}(Y, W) + bd\text{Cov}(Y, V)$

Integration by parts: $\int f \cdot g' = f \cdot g - \int f' \cdot g$
Geometric sums & series: $\sum_{k=0}^n ar^k = a(\frac{1-r^{n+1}}{1-r}) \xrightarrow{n \rightarrow \infty} \frac{a}{1-r}, |r| < 1$
Hoeffding's inequality:

For $Z_i \in [0, 1]$ iid, $\mathbb{P}(\frac{1}{n} \sum_{i=1}^n Z_i - \mathbb{E}[Z] \geq t) \leq \exp(-2nt^2)$.

For $Z_i \in [a_i, b_i]$ the upper bound is $\exp(-\frac{2nt^2}{\sum_{i=1}^n (b_i - a_i)^2})$

Boole's inequality - union bound $\mathbb{P}\left(\bigcup_i A_i\right) \leq \sum_i \mathbb{P}(A_i)$

Cauchy-Schwarz inequality: $v^T u \leq \|v\|_2 \|u\|_2$

Hölder's inequality: $v^T u \leq \|v^T u\|_1 \leq \|v\|_p \|u\|_{p^*}$,

where $1/p + 1/p^* = 1$, equality if $|v|^p = \gamma |u|^p$ *

Jensen's inequality: f convex, then $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$
 $\Rightarrow \log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$

Markov's inequality: $\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$

KL Divergence Let Q, P be prob. distr. of a continuous RV $x \in \mathbb{R}^d$, with densities q, p . Then $KL(Q||P) = \int_{-\infty}^{\infty} q(x) \log \frac{q(x)}{p(x)} dx$

Entropy For $X \in \{x_1, ..., x_n\}$, $H(X) = -\sum P(x_i) \log P(x_i)$

$$\left(\sum_i a_i\right)^2 = \sum_i a_i^2 + 2 \sum_{i < j} a_i a_j$$

For square matrices A, B : $\det A = \det A^t$, $\det AB = \det A \cdot \det B$,
 $\det A^{-1} = \frac{1}{\det A}$

Moment Generating Function $M_X : \mathbb{R}^n \rightarrow \mathbb{R}$, $M_X(t) = \mathbb{E}_x[\exp(t \cdot x)]$.
 $M_{X+Y} = M_X \cdot M_Y$

For $x \sim \mathcal{N}(\mu, \Sigma)$, we have $M_X(t) = \exp(t \cdot \mu + \frac{1}{2} t^T \Sigma t)$

Losses

Cross entropy loss $H(x) = -\sum_{c=1}^M y_{o,c} \log P(x_{o,c})$

Distributions

Standard normal: $p(x|\mu, \sigma^2) = \frac{e^{-(x-\mu)^2/(2\sigma^2)}}{\sqrt{2\pi\sigma^2}}$

Multivariate normal: $x \sim \mathcal{N}(\mu, \Sigma)$

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{k}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu))$$

Exponential: $p(x|\lambda) = \lambda e^{-\lambda x}$

Bernoulli: $p(x|p) = p^x (1-p)^{1-x}$

Binomial: $p(x|n, p) = \binom{n}{x} p^x (1-p)^{n-x}$

Poisson: $p(x|\lambda) = \frac{\lambda^x \exp[-\lambda]}{x!}$

Connectionism

McCulloch & Pitts neuron: $f(x; \sigma, \theta) = \begin{cases} 1, & \sum_{i=1}^n \sigma_i x_i \geq \theta \\ 0, & \text{else} \end{cases}$

$x \in \{0, 1\}^n$, $\sigma \in \{\pm 1\}^n$, $\theta \in \mathbb{Z}$
Disjunctive Normal Form: $f(x; \sigma, \theta) = \bigvee_{I \in \mathcal{I}} (\bigwedge_{i \in I} x_i \bigwedge_{i \notin I} \bar{x}_i)$,
 $\mathcal{I} = \{I : \sum_{i \in I} \sigma_i \geq \theta\}$ (OR of ANDs)

Turing Type A machine (NAND):
 $y(t+1) = 1 - x_1(t)x_2(t)$, $x_1(t), x_2(t) \in \{0, 1\}$

Turing Type B machine:
 $A(1, \text{NAND}(x_1, x_1)) = x_1$, $A(0, \text{NAND}(x_1, x_1)) = 1$
 $\text{NAND}(x_2, ..., x_n) = \begin{cases} \text{NAND}(x_2, ..., x_n) & \text{if } A \leftarrow 1 \\ \text{NAND}(x_1, ..., x_n) & \text{else} \end{cases}$

Perceptron

$(x, \theta) \rightarrow \text{sgn}(x \cdot \theta)$, update rule: $\Delta \theta = \begin{cases} 0, & y(x \cdot \theta) \geq 0 \\ yx, & \text{otherwise} \end{cases}$

Path of updates always zig-zag since $\Delta \theta \cdot \theta < 0$ for updates
Update rule is SGD for the loss: $l(x; y, \theta) = \max\{0, -yx \cdot \theta\}$

Lemma (Norm Growth): (x^t, y^t) perceptron mistakes inducing

updates $\Delta \theta^t, \theta^s = \sum_{t=1}^s \Delta \theta^t$. Then:
 $\|\theta^s\|^2 \leq \sum_{t=1}^s \|x^t\|^2$ (prove by induction!)
Cor: If $\|x^t\| \leq 1$ then $\|\theta^s\| \leq \sqrt{s}$
Def (Linear Separability): \mathcal{S} linearly separable with margin $\gamma > 0$
if $\exists \theta^*, \|\theta^*\| = 1 : yx \cdot \theta^* \geq \gamma > 0 \quad \forall (x, y) \in \mathcal{S}$

Novikov's convergence theorem: Converges in at most γ^{-2} steps (if assume $\|x^t\| \leq 1$). Show $\gamma s \leq s$ and start with fact that $\Delta \theta^t \theta^* = y^t x^t \theta^* \geq \gamma + \theta^s = \sum \Delta \theta^t + \text{Cauchy-Schwarz}$
Cover's theorem: Dichotomies possible with linear separators: $C(s, n) = 2 \sum_{i=1}^{n-1} \binom{s-1}{i}$
Prove by showing recurrence relation $C(s+1, n) = C(s, n) + C(s, n-1)$ using Pascal's rule: $\binom{n-1}{k} + \binom{n-1}{k-1} = \binom{n}{k}$

VC dimension: largest set S that can be shattered. Corollary of Cover's thm: n points can be shattered by linear functions in n dimensions (meaning can realize all 2^n points). $m > n$ can not.
Willshaw Memory

Hebb rule: $\Delta \theta_{ij}^t \propto x_i^t x_j^t$ (neurons that fire together, wire together)
r-sparse Boolean vectors $\mathbb{B}_r^n = \{x \in \{0, 1\}^n | \sum x_i \leq r\}$

Upper bound on number of patterns $s : s \leq \frac{n^2}{r \log n}$,
 $\log \binom{n}{r} \approx r \log n$ is pattern information and n^2 total number of bits

Binary memory matrix: $\Theta_{ji} = \min\{1, \sum_{t=1}^s y_j^t x_i^t\}$, $\Theta_{ji} \in \{0, 1\}^{n \times n}$
Alternatively: $\Theta = \min(1, \sum_{t=1}^s y^t (x^t)^T)$

Retrieve y from given x : $z = \Theta x$, $y_j = \begin{cases} 0 & z_j < r \\ 1 & \text{else} \end{cases}$

Monotonicity: $x^t \mapsto y \geq y^t, \forall (x^t, y^t) \in \mathcal{S}$ (i.e. get at least as much as ask for). Proof: $\Theta = \min(1, \sum_{\tau} \Theta^{\tau}) \geq \min(1, \Theta^t) = \Theta^t$

$z^t = \Theta x^t \geq \Theta^t x^t = y^t (x^t)^T x^t = r y^t$

Maximal capacity: $\max_q I(\text{patterns}) = (\log 2)n^2 \approx 0.693n^2$ in the limit, for $r = \log n$. Half of θ_{ji} will be 0 and half will be 1 as $n \rightarrow \infty$

Hopfield Networks

Idea: Store patterns in Θ , for recovery from corrupted patterns
 $\Theta = \sum_{t=1}^s [x_t x_t^T - I_n] \in \mathbb{Z}^{n \times n}$
Asynchronous update: $\hat{x}_i \leftarrow \text{sign}(\sum_{j \neq i} \theta_{ij} \hat{x}_j + \theta_{i0})$ (Hopfield dynamics)
Synchronous updates can lead to limit cycles
Lyapunov function: $= -\sum_{i,j=1}^n \theta_{ij} x_i x_j - \sum_{i=1}^n \theta_{i0} x_i$

Hebb -Hopfield $\Theta \approx \sum_{t=1}^s [x^t x^t - I] \in \mathbb{R}^{n \times n}$. Use fact that for \hat{x} with

hamming-dist k to x , $x \cdot \hat{x} - 1 = (n - 2k - 1)$ to show $-x, x$ both attractors. Capacity for num of patterns is $\frac{n}{n} \leq \alpha^* \approx 0.138$. For Hopfield (without Hebb's rule, but with optimal weights) $\alpha^* = 2$

Linear Networks

Linear unit: $u(x; \theta) = x \cdot \theta$, Affine unit: $u(x; \theta, b) = x \cdot \theta + b$
Level sets: $u(x + \Delta x, \theta) = (x, \theta)$, $\Delta x \perp \theta$
Linear unit defines: direction of change via $\frac{\theta}{\|\theta\|}$, rate of change

via $\|\theta\|$
Learning algo for linear units: $\Delta \theta = \eta(y - x \cdot \theta)x$
Homogeneity: $f(ax) = \alpha f(x)$
Additivity: $f(x + y) = f(x) + f(y)$
 $\Rightarrow f$ is linear with these properties
 f, g linear $\Rightarrow f \circ g$ linear
Affine functions: $f(x + \beta y) = \alpha f(x) + \beta f(y)$, $\forall x, y : \forall \alpha, \beta : \alpha + \beta = 1$

Linear Autoencoder

Def: $x \mapsto z \mapsto y$, $z = Cx, y = Dz$, $C, D^T \in \mathbb{R}^{m \times n}, m < n$
Loss: $l(x) = \frac{1}{2} \|x - y\|^2$
Matrix notation: $X, Y \in \mathbb{R}^{n \times s} : \theta = (C, D) \rightarrow \min \frac{1}{2s} \|X - DCX\|_F^2$
 $\frac{\delta l(x)}{\partial C} = D^T (y - x)x^T \in \mathbb{R}^{m \times n}$
 $\frac{\delta l(x)}{\partial D} = (y - x)x^T C^T \in \mathbb{R}^{n \times m}$
 $\text{rank}(DC) \leq \min\{\text{rank}(C), \text{rank}(D)\} \leq m < n \Rightarrow \text{rank}(Y) \leq m$
Eckhart-Young Theorem: $\|X - X_r\|_F = \min_{\text{rank}(Y) \leq r} \|X - Y\|_F$
 $C = U_m^T, D = U_m \Rightarrow DCX = X_m$

Gradients of deep linear networks

Assume X, Y centered, X whitened: $X \mapsto \Lambda^{-\frac{1}{2}} U^T X$, s.t. $\frac{1}{s} XX^T = I$
Then LS problem $\Theta \rightarrow \min \frac{1}{s} \|Y - \Theta X\|_F^2$ can be written:
 $\Theta \rightarrow \min \frac{1}{2} \|\Theta - \Gamma\|_F^2$, $\Gamma := \frac{1}{2s} XY^T$
Let $\Theta = DC$, then: $\frac{\delta l}{\partial C} = D^T (\Theta - \Gamma)$, $\frac{\delta l}{\partial D} = (\Theta - \Gamma)C^T$
Now re-write: $DC - \Gamma = U(\tilde{D}\tilde{C} - \Sigma)V^T$, $\tilde{D} = U^T D, \tilde{C} = CV$
Then: $\frac{\delta l}{\partial \tilde{C}} = (\tilde{D}\tilde{C} - \Sigma)^T \tilde{D}$, $\frac{\delta l}{\partial \tilde{D}} = \tilde{C}(\tilde{D}\tilde{C} - \Sigma)C^T$
Let d_r be rows of \tilde{D} and c_r be columns of \tilde{C} . Then can minimize:
 $\tilde{l}(C, D) = \frac{1}{2} \sum_r (d_r \cdot c_r - \sigma_r)^2 + \frac{1}{2} \sum_{r \neq q} (d_r \cdot c_q)^2$ (cooperative and competitive terms, second orthogonalizes)

Sigmoid Networks

Ridge function if can be written as a composition of an affine function: $f = \phi \circ u$, $f(x; \theta) = \phi(x \cdot \theta)$, preserves level sets and directional sensitivity
But not rate of change: $|\nabla_x f(x; \theta)| = |\phi'(x \cdot \theta)| \cdot \|\theta\|$
Threshold units like Heaviside and sign gives no derivative info
 \Rightarrow **Logistic unit:** $\sigma(z) = \frac{1}{1 + \exp[-z]}$, $\sigma^{-1}(t) = \log \frac{t}{1-t}$ (log-odds)
 $\sigma'(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z)$, smooth since polynomials in σ

Hyperbolic tangent: $\tanh(z) := \frac{e^z - e^{-z}}{e^z + e^{-z}} = 2\sigma(2z) - 1$
 $\tanh'(z) = 1 - \tanh^2(z)$
Sometimes preferred since range $(-1, 1)$ symmetric around 0
Logistic regression
Cross-entropy loss $\{y \in \{-1, +1\}\} : l(x, y; \theta) = -\log \sigma(yx \cdot \theta)$
 $\Rightarrow \nabla_{\theta} l(x, y) = -\sigma(-yx \cdot \theta)yx$
Cross-entropy loss $\{y \in \{0, 1\}\} :$
 $l(x, y; \theta) = -y \log \sigma(x \cdot \theta) - (1 - y) \log(1 - \sigma(x \cdot \theta))$
Softmax
 $\sigma_i^{\max}(x; \Theta) = \frac{\exp[x \cdot \theta_i]}{\sum_{j=1}^k \exp[x \cdot \theta_j]}$, $\Theta = \{\theta_1, ..., \theta_k\}$

Over-parameterized (can add/subtract any constant vector)
Equal to sigmoid unit for $k = 2$, $\theta = \theta_1 - \theta_2$
 $\nabla_{x_j} \sigma_i^{\max} = \begin{cases} \sigma_i^{\max}(1 - \sigma_i^{\max}), & i = j \\ -\sigma_i^{\max} \sigma_j^{\max}, & i \neq j \end{cases}$

Softmax regression: $l(x, y; \Theta) = -y \cdot \log \sigma^{\max}(x; \Theta), y \in \{e_1, ..., e_k\}$
Alternatively: $l(x, y; \Theta) = -\sum_{i=1}^k [y_i x \cdot \theta_i + \log(\sum_{j=1}^k \exp[x \cdot \theta_j])]$
 $\Rightarrow \nabla_{\theta_i} l(x, y; \Theta) = (\sigma_i^{\max} - y_i)x$

Approximation Theory

Definitions and notation
 $f \approx g \iff \text{approx-err}(f, g) = \inf\{g \in \mathcal{G} | \|f - g\|_{\infty} = 0\}$
Uniform convergence: $(g_m) \xrightarrow{\infty} f \iff \forall \epsilon > 0 : \exists m \geq 1 : \|g_m - f\|_{\infty} < \epsilon$
Follows that: $\mathcal{G} \supset g_m \xrightarrow{\infty} f \Rightarrow f \approx g$
Denseness: $\mathcal{G} \subset \mathcal{F}$ dense in $\mathcal{F} \iff \mathcal{F} \approx \mathcal{G} \iff \forall f \in \mathcal{F}, f \approx g$
Closure is all functions that can be approximated by $\mathcal{G} : \text{cl}(\mathcal{G})$
 \mathcal{G} **universal approximator** $\iff C(S) \approx \mathcal{G}(S) \forall$ compact $S \subset \mathbb{R}^n$ (compact = closed and bounded)

Weierstrass Theorem
Polynomials \mathcal{P} are dense in $C([a, b]) \forall a, b \in \mathbb{R}$
Universal approximation theorem (1d)
Let $\sigma \in C^{\infty}(\mathbb{R})$ smooth & not polynomial,
 $\mathcal{G}_\sigma^1 = \{g : g(x) = \sigma(ax + b) \quad a, b \in \mathbb{R}\}, H_\sigma^1 = \text{span}(\mathcal{G}_\sigma^1)$,
then H_σ^1 is a universal approximator (can't pick polynomial because then we form polynomials of limited degree)

Ridge function theorem
 $\mathcal{G}_\sigma^n = \{g : g(x) = \sigma(x \cdot \theta) \theta \in \mathbb{R}^n\}, \mathcal{G}^n = \bigcup_{\sigma \in C(\mathbb{R})} \mathcal{G}_\sigma^n, H_\sigma^n = \text{span}(\mathcal{G}^n)$
Then H^n is a universal function approximator
(Problem: here we can pick any combination of ridge functions)
Dimension lifting theorem
 H_σ^1 univ. approx for $C(\mathbb{R}) \Rightarrow H_\sigma^n$ univ. approx for $C(\mathbb{R}^n) \quad \forall n \geq 1$
Barron's theorem (number of units required)
Gradient regularity condition: $C_g = \int \|\omega\| \cdot |\hat{g}(\omega)| d\omega < \infty$
If g differentiable, then $\nabla \hat{g}(\omega) = \omega \cdot \hat{g}(\omega)$

Theorem: Let σ be bounded, monotonic s.t. $\lim_{t \rightarrow \infty} \sigma(t) = 1$ & $\lim_{t \rightarrow -\infty} \sigma(t) = 0$. Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ with $C_g < \infty$, & $r > 0$. Then $\exists (f_m(x))_{m=1}^{\infty}$ sequence defined as $f_m(x) = \sum_{j=1}^m (\beta_j \sigma(\theta_j + b_j) + b_0)$, s.t.
 $\int_{\mathbb{R}^n} (g(x) - f_m(x))^2 \mu(dx) \leq O(\frac{1}{m})$,
 $r_B = \{x \in \mathbb{R}^n : \|x\| \leq r\}$, μ probability measure. Independent of $n!$ no curse of dimensionality when approx. certain functions
Benefits of depth
Function $g(x) = \psi(\|x\|)$ has exponential advantages in approximation with 2 layers vs 1:
Upper bound on how much space covered by 1 layer: me^{-n}

Backpropagation

Notation
Network $F = F_{k:1} = F_k \circ F_{k-1} \circ ... \circ F_1$
Width of a layer: $\text{width}_l = n_l := \dim(\text{range}(F_l))$
Using activations: $z_l := F_{l:1}(x) = (F_l \circ F_{l-1:1})(x) = F_l(z_{l-1})$
Jacobian map for $F : \mathbb{R}^n \rightarrow \mathbb{R}^m : \partial F = (\partial_{ij} F) : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$
Chain rule for maps: $\underbrace{\partial(G \circ F)}_{\mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}} = \underbrace{(\partial G \circ F)}_{\mathbb{R}^n \rightarrow (\mathbb{R}^{m \times k}, \mathbb{R}^{k \times n})}$

The Jacobian is $\partial F = \prod_{l=k}^1 \partial F_l \circ F_{l-1:1}$, $\partial F(x) = \prod_{l=k}^1 \partial F_l(z_{l-1})$
With loss function: $f = l \circ F$
Gradient information provides direction of steepest descent:

$\lim_{\eta \rightarrow 0} \arg \min_{\theta} g_{\eta}[\theta] = -\frac{\nabla_{\theta} f(x; \theta)}{\|\nabla_{\theta} f(x; \theta)\|}$
Derivative wrt parameter in layer: $h' \circ g$
Derivative wrt parameter in next layer: $(\partial h \circ g) \cdot g'$

Backpropagation
 $\nabla_{\theta_l} F = (\partial F_{k:l+1} \circ F_{l:1}) \cdot (F'_l \circ F_{l-1:1})$
Or with activity vectors: $\nabla_{\theta_l} F(x) = \partial F_{k:l+1}(z_l) \cdot F'_l(z_{l-1})$
(downstream Jacobian \times local Jacobian). $F = G \circ F_{\theta} \circ H$, apply chain rule twice $(\partial G \circ F_{\theta} \circ H) \cdot (F_{\theta} \circ H)' = (\partial G \circ F_{\theta} \circ H) \cdot (F'_{\theta} \circ H)$.
 ∂l is a Jacobi vector. Get computational savings by multiplying it with Jacobi matrices in reverse order (this is exactly backprop!)
 $\nabla_{\theta_l} f(x) = \partial(l \circ F_{k:l+1}(z_l) \cdot F'_l(z_{l-1})) =: \xi_{l+1} \cdot F'_l(z_{l-1})$
Row vectors ξ_i can be calculated through backward iteration:
 $\xi_{k+1} = \partial l(z_k)$, $\xi_l = \xi_{l+1} \partial F_l(z_{l-1})$
Backpropagation algorithm can be written in three steps: (1) forward pass computing z_i^f, s , (2) backward pass computing ξ_i^f, s and (3) local computations computing $\nabla_{\theta_l} f(x)'s$

Automatic differentiation
Reverse mode (backpropagation) and **forward mode**
Forward mode: For each parameter, compute the derivatives for each of the (intermediate) outputs starting from the parameter (no forward pass here!)
For $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$:
Reverse mode more efficient if $M < N$ (scales in $M, \mathcal{O}((V+E)M)$)
Forward mode more efficient if $M > N$ (scales in $N, \mathcal{O}(N(V+E))$)
For DL: $N = \# \text{parameters}$ and $M = 1$ (scalar loss)
Reverse mode requires in general more memory due to storing results from forward pass

Derivatives
Element-wise function (e.g. activation fun): $\frac{\partial}{\partial x} f(x) = \text{diag}(f'(x))$
ReLU wrt pre-activations: $\frac{\partial}{\partial x} \max(0, x) = \text{diag}(H(x))$
ReLU wrt params: $\frac{\partial}{\partial W_{ij}} \max(0, Wx)_k = \max(0, \text{sign}(W_i^T x)) x_j \delta_{ik}$
ReLU wrt activations: $\frac{\partial}{\partial x_j} \max(0, Wx)_i = \max(0, \text{sign}(W_i^T x)) \cdot W_{ij}$

Rectified Networks

Sigmoids suffer from vanishing gradients / unreliable gradient information in deep networks: $\sigma'(z) = \sigma(z)\sigma(-z) \xrightarrow{z \rightarrow \pm \infty} 0$
Rectified units: continuous piecewise linear
Additional benefit of rectified units: computationally faster
Rectified Linear Unit (ReLU)
 $(x, \theta) \mapsto (x \cdot \theta)_+ = \max(0, x \cdot \theta)$
Splits the input space in two half-spaces separated by hyper-

plane $\mathcal{H}_{\boldsymbol{\theta}}^0: \quad \mathcal{H}_{\boldsymbol{\theta}}^+ = \{\mathbf{x} : \mathbf{x} \cdot \boldsymbol{\theta} > 0\}, \quad \mathcal{H}_{\boldsymbol{\theta}}^- = \{\mathbf{x} : \mathbf{x} \cdot \boldsymbol{\theta} < 0\}$

Subderivative: $\partial(z)_+ = \begin{cases} 1 & z > 0 \\ [0, 1] & z = 0 \\ \text{else} & \end{cases}$

Definition of subderivative of convex f at z_0 :
 $\partial f(z_0) = \{c : f(z) - f(z_0) \geq c(z - z_0)\}$

Activation patterns: units can be either active $\mathbf{x} \in \mathcal{H}_{\boldsymbol{\theta}}^+$ or inactive $\mathbf{x} \in \mathcal{H}_{\boldsymbol{\theta}}^-$ with patterns $H(\boldsymbol{\Theta}\mathbf{x}) \in \{0, 1\}^m$

$\|H(\boldsymbol{\Theta}\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\| \leq 2^m$, but is less since not every boolean vector is a valid activation pattern (because not every activation pattern might actually have an input that produces it)

Vanishing gradient: For sigmoids the sensitivity (gradient wrt activation) can go to 0, which means that slight changes in parameters do not lead to any signal. With ReLUs: $\partial(z)_+ = 1 \quad \forall z > 0$

Backprop: If unit z_{lj} inactive then $\nabla_{\boldsymbol{\theta}_{lj}} z_{lj} = \nabla_{z_{l-1}} z_{lj} = 0$

Parameter gradients vanish and Jacobi matrix is sparse:

$$\partial F_l := \hat{\boldsymbol{\Theta}}_l := \begin{cases} \begin{bmatrix} \hat{\boldsymbol{\theta}}_{l1}^T \\ \vdots \\ \hat{\boldsymbol{\theta}}_{lm}^T \end{bmatrix} & \boldsymbol{\theta}_{lj}^T = \begin{cases} 0 & z_{lj} = 0 \\ \boldsymbol{\theta}_{l1}^T & \text{otherwise} \end{cases} \end{cases}$$

Dying ReLUs: If a specific unit is inactive for all inputs (can happen on initialization or during training) \Rightarrow parameters will not be updated. Could prune or re-initialize

Absolute Value Unit (AbsU): $|z|, \quad \partial|z| = \begin{cases} 1 & z > 0 \\ [-1, 1] & z = 0 \\ -1 & z < 0 \end{cases}$

Relation to ReLU: $(z)_+ = \frac{z+|z|}{2}$ and $|z| = 2(z)_+ - z = (z)_+ + (-z)_+$

No sparseness property as in ReLUs, but symmetric

Smooth ReLU approximations: Combine rectification and smoothness

Softplus: $(\mathbf{x}, \boldsymbol{\theta}) \mapsto \log(1 + \exp[\mathbf{x} \cdot \boldsymbol{\theta}]) \in (0, \infty)$

Exponential linear unit: $(\mathbf{x}, \boldsymbol{\theta}) \mapsto \begin{cases} \mathbf{x} \cdot \boldsymbol{\theta} & \mathbf{x} \cdot \boldsymbol{\theta} \geq 0 \\ \exp[\mathbf{x} \cdot \boldsymbol{\theta}] - 1 & \text{else} \end{cases} \in (-1, \infty)$

Leaky ReLU: Gives some gradient information even in (low sensitivity instead of no sensitivity, typical $\epsilon = 0.01$)

$(\mathbf{x}, \boldsymbol{\theta}) \mapsto \begin{cases} \mathbf{x} \cdot \boldsymbol{\theta} & \mathbf{x} \cdot \boldsymbol{\theta} \geq 0 \\ \epsilon \mathbf{x} \cdot \boldsymbol{\theta} & \text{else} \end{cases} \in \mathbb{R}$

Universal function approximators

Theorem: Piecewise linear functions are dense in $C([0, 1])$

Theorem: A piecewise linear function with m pieces can be written as $g(x) = ax + b + \sum_{i=1}^{m-1} c_i(x - x_i)_+$ (sum of ReLU units)

(alternative representation instead with absolute value function)

Corollary: Networks with one hidden layer of ReLU or AbsU are universal function approximators

Minimal non-linearity

k-Hinge functions: $g(\mathbf{x}) = \max_{j=1}^k \{\boldsymbol{\theta}_j \cdot \mathbf{x} + b_j\}$ (aka maxout units)

Representational power: $2 \max\{f, g\} = f + g + |f - g|$

Theorem: Every continuous piecewise linear function can be written as a signed sum of k-Hinges with $k \leq n + 1$

Polyhedral functions: f is polyhedral $\Leftrightarrow \text{epi}(f)$ is polyhedral set

S polyhedral $\Leftrightarrow S$ is finite intersection of closed half-spaces, i.e. $S = \{\mathbf{x} \in \mathbb{R}^n : \boldsymbol{\theta}_j \cdot \mathbf{x} + b_j \geq 0, j = 1, \dots, r\}$

Theorem: If f polyhedral then $\exists A \subset \mathbb{R}^{n+1}, |A| = k$ s.t. $f(\mathbf{x}) = \max_{(\boldsymbol{\theta}, b) \in A} \{\boldsymbol{\theta} \cdot \mathbf{x} + b\}$

Theorem: Every continuous piecewise linear function f can be written as the difference of two polyhedral functions

\Rightarrow **Theorem:** Maxout networks with two maxout units (difference of two k-Hinges) are universal function approximators

Optimization

Losses

Squared loss: $l_y(\mathbf{v}) = \frac{1}{2} \|\mathbf{y} - \mathbf{v}\|^2$

Zero-one loss: $l_y(v) = \begin{cases} 0, & v = y \\ 1, & \text{else} \end{cases}$

Log-loss (multiclass): $l_y(v) = -\log v_y$

Soft target cross-entropy $(\mathbf{y} \in [0, 1]^m)$:
 $l_y(\mathbf{v}) = -\sum_{j=1}^m y_j \log v_j \geq -\sum_{j=1}^m y_j \log v_j =: H(\mathbf{y})$

Probabilistic loss: $l_y(\mathbf{v}) = -\log p(\mathbf{y}; \mathbf{v})$ (e.g. replace with isotropic Gaussian to get squared loss)

Exponential family: $p(\mathbf{y}; \mathbf{v}) = h(\mathbf{y}) \exp[\mathbf{y} \cdot \mathbf{v} - \psi(\mathbf{v})]$, log partition/normalizing function ψ

Can construct losses by taking distributions in the exponential family in the log prob loss and replace \mathbf{v} with $h(\boldsymbol{\theta} \cdot \mathbf{z})$, where \mathbf{z} is produced by a Neural Network

Gradient Descent and Optimization Theory

$\mathbf{x}^{k+1} = \mathbf{x}^k - \eta \nabla f(\mathbf{x}^k)$

Solution to $\dot{\mathbf{x}} = -\nabla f(\mathbf{x})$ is called gradient flow - approximated by GD

Notes on η : Not typical to optimize by line search for neural networks since expensive to evaluate. Typically keep constant or use learning rate schedule. Selecting η is a trade-off between computational speed (larger better) and convergence & approximation of gradient flow (smaller better)

Quadratic Model: 2nd order Taylor approximation yields $f(\mathbf{x} + \Delta \mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^T \nabla^2 f(\mathbf{x}) \Delta \mathbf{x}$

\Rightarrow minimizer is $\Delta \mathbf{x} = -[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$ (Newton's method)

Setting $\nabla^2 f(\mathbf{x}) = \frac{1}{\eta} \mathbf{I}$ yields gradient descent (η curvature of quadratic and step size, smaller η gives more curvature and hence smaller step)

Convex: $f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}), \quad \forall \lambda \in [0, 1]$

FOC of convexity: $f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y})$

Lipschitz smoothness: $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

$\Leftrightarrow f(\mathbf{x}) \leq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$

Strong convexity: $f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2$

Hessian of smooth and strongly convex function: $\mu \mathbf{I} \preceq \nabla^2 f \preceq L \mathbf{I}$

ϵ -Stationarity: $\|\nabla f(\mathbf{x})\| \leq \epsilon$ (to measure local convergence)

PL condition (generalization of strong convexity without convexity):
 $\frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu(f(\mathbf{x}) - f^*), \quad \forall \mathbf{x}, f^* = \min f(\mathbf{x})$

SGD

Assume additive structure $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \eta_k \nabla f_{I(k)}(\mathbf{x}^k), \quad I(k) \sim \text{Unif}(1, \dots, n)$$

Unbiased, $\mathbb{E}[\nabla f_{I(k)}(\mathbf{x})] = \nabla f(\mathbf{x})$.

$\mathbf{V}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2$

Polyak Averages: To combat variance around \mathbf{x}^*

$\bar{\mathbf{x}}^{k+1} = \frac{k}{k+1} \bar{\mathbf{x}}^k + \frac{1}{k+1} \mathbf{x}^{k+1}$

Minibatch SGD Sample r functions f_j . Unbiased and reduces variance ar . Smaller r implies more noise, yet better results in non-convex setting. Noise might help avoid getting stuck in bad regions. Batch size depending on concurrency model of GPU, has to fit in GPU memory.

Learning rate For theoretical results typically $\eta_k \propto \frac{1}{k}$, since $\sum \infty \frac{1}{k} = \infty$ (min requirement, to ensure we can approach end point) and $\sum \infty \frac{1}{k^2} \leq \infty$. However in practice: keep step size constant or reduce step size at a small number of points.

Theorem Assume $f = \sum_i f_i$, each f_i L_i -smooth, $\sup_i L_i \leq L$ and f μ -strongly convex. \mathbf{x}^* minimizer of f and $\sigma^2 := V(\mathbf{x}^*) = \frac{1}{n} \sum \|\nabla f_i(\mathbf{x})\|^2$. \mathbf{x}^k SGD iterate generated with $\eta \leq 1/\mu$. Then

$$\mathbb{E} \|\mathbf{x}^k - \mathbf{x}^*\|^2 \leq A^k \|\mathbf{x}^0 - \mathbf{x}^*\|^2 + B,$$

where $A = 1 - 2\eta\mu(1 - \eta L), \quad B = \frac{\eta\sigma^2}{\mu - \eta L}$, the bigger the step size, the more stochasticity. Advantage over SGD: can return a set of visited parameters.

Variance reduction: $\mathbf{x}^{k+1} = \mathbf{x}^k - \eta[\nabla f_j(\mathbf{x}^k) - \nabla f_j(\bar{\mathbf{x}}) + \nabla f(\bar{\mathbf{x}})]$

Gradient clipping: To avoid exploding gradients, can bound the norm of the gradient in the update step by a threshold (e.g. when reaching a "cliff")

Momentum

If gradients start going into a particular direction, keep going in that same direction (inertia in direction we had)

Nesterov's Acceleration Method

$$\mathbf{y}_{k+1} = \mathbf{x}_k + \beta(\mathbf{x}_k - \mathbf{x}_{k-1})$$
$$\mathbf{x}_{k+1} = \mathbf{y}_{k+1} - \eta \nabla f(\mathbf{y}_{k+1})$$

Theorem Let f L -smooth and μ -strongly convex, $\kappa = \frac{L}{\mu}, \quad \beta = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$. Then $f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq L \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa}} \right)^k \|\mathbf{x}^0 - \mathbf{x}^*\|^2$.

Polyak's Heavy ball method

$\mathbf{x}^{k+1} = \mathbf{x}^k - \eta \nabla f(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) \quad \beta \in (0, 1)$.

Adaptivity

Adapt learning rate per parameter or dimension. Advantage in compositional models: adapt step size for different parameters in different layers (non-uniformity of parameters).

AdaGrad Increasing sequence:

$\gamma^k = \gamma^{k-1} + \nabla f(\mathbf{x}^k) \odot \nabla f(\mathbf{x}^k)$.

γ^k will be large for parameters that have received large updates. Use these as pre-conditioner matrix

$\mathbf{x}^{k+1} = \mathbf{x}^k - \eta \Lambda^k \nabla f(\mathbf{x}^*)$,

where $\Lambda^k = \text{diag}(\lambda_i^k)$, with $\lambda_i^k = \frac{1}{\sqrt{\gamma_i^k} + \delta}, \quad \delta > 0$

If something has received very significant updates in the past, then we decay the learning rate faster.

Adam Momentum + adaptivity

$\mathbf{m}^k = \beta \mathbf{m}^{k-1} + (1 - \beta) \nabla f(\mathbf{x}^k), \quad \beta \in [0, 1], \quad \mathbf{m}^0 = 0$

$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{\eta}{1 - \beta^k} \mathbf{m}^k$.

Since $\mathbf{m}^0 = 0$ estimate is biased, $\frac{1}{(1 - \beta^k)}$ corrects for this.

$\gamma^k = \alpha \gamma^{k-1} + (1 - \alpha) [\nabla f(\mathbf{x}^k) \odot \nabla f(\mathbf{x}^k)]$

$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{\eta}{1 - \beta^k} \Lambda^k \mathbf{m}^k$

$\Lambda^k = \text{diag}(\lambda_i^k)$, with $\frac{1}{\lambda_i^k} = \sqrt{\frac{\gamma_i^k}{1 - \alpha^k}} + \delta$.

As AdaGrad, not parametrization invariant (matrix depends on basis of chosen parametrization). Typical values $\beta = 0.9, \alpha = 0.99$.

AMSGrad: enforce monotonic increase of γ to ensure convergence, $\gamma^{k+1} = \max\{\gamma^k, \gamma^{k+1}\}$

Compressed stochastic gradients

SignSGD: $\mathbf{x}^{k+1} = \mathbf{x}^k - \eta_k \text{sign}(\nabla f_{I(k)}(\mathbf{x}^k))$. Reduces computational and communication complexity.

Convolutional Networks

Convolution Operator

Integral operators Kernel $H : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad -\infty \leq t_1 \leq t_2 \leq \infty$

$(Tf)(u) = \int_{t_1}^{t_2} H(u, t) f(t) dt$

Fourier Transform

$t_1 = -\infty, \quad t_2 = \infty, \quad H(u, t) = \exp\{-2\pi i t u\}$

$(\mathcal{F}f)(u) := \int_{-\infty}^{\infty} \exp\{-2\pi i t u\} f(t) dt$

Convolution

$(f * h)(u) := \int_{-\infty}^{\infty} h(u - t) f(t) dt = \int_{-\infty}^{\infty} f(u - t) h(t) dt$

Convolutions are **commutative** (exchange function and kernel)

And **shift-equivariant:** $f_{\Delta}(t) := f(t + \Delta) \Rightarrow f_{\Delta} * h = (f * h)_{\Delta}$

Linear shift-equivariant transforms

$T(\alpha f + \beta g) = \alpha T f + \beta T g, \quad \forall f, g; \quad \forall \alpha, \beta \in \mathbb{R}$

$(T f)_{\Delta}(t) = (T f)(t + \Delta)$

Theorem: Any linear translation-equivariant (shift-equivariant) transformation T can be written as a convolution with some h

Discrete convolutions

$f, h: \mathbb{Z} \rightarrow \mathbb{R}$. Define discrete convolution via $(f * h)[u] := \sum_{t=-\infty}^{\infty} f[t] h[u - t]$

Typical choice of h : support over finite window, e.g. $h(t) = 0$ for

$t \notin [t_{\min}, t_{\min}]$. In higher dimensions, replace vectors by matrices or fields $(F * G)[i, j] = \sum_k \sum_l F[i - k, j - l] G[k, l]$

Cross-correlation

$(f * h)[u] := \sum_{t=-\infty}^{\infty} f[t] h[u + t]$

$(f * h) = (\bar{f} * h), \quad \bar{f}[t] := f[-t]$

Toeplitz matrices: Constant on the diagonals (from exercise: parameters of a 1D convolutional layer can be written as a dense layer through a Toeplitz matrix)

Convolutional Neural Networks

Exploiting translation equivariance

Exploiting locality and scale (temporal, spatial...)

Increased efficiency through parameter sharing

Receptive fields and sparse connectivity Activity of one unit doesn't depend on all units from previous layer. Convolved signal **inherits** topology of original signal. Can create longer range dependencies (larger receptive field) by nesting of convolutions

Border handling: through *same* padding (with zeros) to retain dimension or *valid* padding to only retain values from windows fully contained in support of signal

Half (same) padding (Assuming unit strides) For any i and k odd $(k = 2n + 1, \quad n \in \mathbb{N}), s = 1$ and $p = \lfloor \frac{k}{2} \rfloor = n$

$o = i + 2\lfloor \frac{k}{s} \rfloor - (k - 1) = i$.

Full padding In this setting every possible partial or complete superimposition of the kernel on the input feature map is taken into account.

For any $i, \quad k$ and for $s = 1$ and $p = k - 1$

$o = i + k - 1$

Backpropagation

Exploit structural sparseness in computing $\frac{\partial x_i^l}{\partial x_{j'}^{l-1}}$

Receptive field of $x_i^l : I_i^l := \{j : w_{ij}^l \neq 0\}, \mathbf{W}^l$ Toeplitz matrix of the convolution. $\frac{\partial x_i^l}{\partial x_{j'}^{l-1}} = 0$, for $j \notin I_i^l$.

Weight sharing $\frac{\partial \mathcal{R}}{\partial h_j^l} = \sum_i \frac{\partial \mathcal{R}}{\partial x_i^l} \frac{\partial x_i^l}{\partial h_j^l} h_j^l$ kernel weight. Weights are reused for every unit within target layer.

Backprop example: Let $(\mathbf{x} * \mathbf{w})_{ij} = \sum_{k=1}^q \sum_{l=1}^q x_{i+q-k, j+q-l} w_{k, l}, \quad 1 \leq i, j \leq d - q + 1$. Let $r = d - q + 1$, then $(\mathbf{x} * \mathbf{w}) \in \mathbb{R}^{r \times r}$. Define: $\mathbf{f}(\mathbf{x}) = \mathbf{v}^T \text{vec}(\sigma(\mathbf{x} * \mathbf{w}))$. Then: $\frac{\partial f(\mathbf{x})}{\partial \mathbf{w}_{ab}} = \sum_{k, l=1}^r \text{mat}(\mathbf{v})_{kl} \sigma'((\mathbf{x} * \mathbf{w})_{kl}) x_{k+q-a, l+q-b} = \dots = (\text{rot}_{\pi}(\mathbf{x}) * (\text{mat}(\mathbf{v}) \odot \sigma'(\mathbf{x} * \mathbf{w})))_{ab}$, where $\text{rot}_{\pi}(\mathbf{A})$ flips rows and columns

Stages

- Non-linearities
- Pooling
- Sub-sampling (strides)

Pooling 1D: $x_i^{\max} = \max\{x_{i+k} : 0 \leq k < r\}$

Pooling 2D: $x_{ij}^{\max} = \max\{x_{i+k, j+l} : 0 \leq k < r, 0 \leq l < r\}$

Stacking convolutional layers (X, Y 3rd order tensors):
 $y[r][s, t] = \sum_u \sum_{\Delta s, \Delta t} w[r, u][\Delta s, \Delta t] x[u][s + \Delta s, t + \Delta t]$

Convolution arithmetic - output dimensions

Dimension-wise: For any input size i , kernel size k , padding p and strides s ,

$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1$

Convolution arithmetic - receptive field

Dimension-wise recursion formula: For receptive field r_l , stride s_l and kernel size k_l

$r_{l-1} = s_l \cdot r_l + (k_l - s_l), \quad r_L = 1$

Deep Gradients

Tricks and approaches to improve learning in deep networks

Short Connectivity

Motivation: Avoid vanishing gradients by using activations from previous layers until adjacent layer learns its weights
Use if get worse performance on train set by going deeper
Residual layers: $g(\mathbf{x}) = \mathbf{x} + f(\mathbf{x}; \boldsymbol{\theta})$
If weights are initialized around zero, the gradient might not be able to propagate through properly. Now Jacobian at initialization ($f \approx 0$) will be: $\mathbf{J}_g = \mathbf{I} + \mathbf{J}_f \approx \mathbf{I}$, which is a better start since the gradient will flow through with identity map
If dimensions do not match: $g(\mathbf{x}) = \mathbf{W}\mathbf{x} + f(\mathbf{x}; \boldsymbol{\theta})$, $\mathbf{J}_g = \mathbf{W} + \mathbf{J}_f \approx \mathbf{W}$ at initialization $f \approx 0$
Can skip an arbitrary number of layers
Empirically see that ResNets allows the model to better take advantage of depth. Common in computer vision.
Note: They improve training, but do not increase representational power!

Dense Connectivity

Connect layer output to all downstream layers (\Rightarrow get more channels in CNNs). Concatenate activations instead of adding as in ResNets
Residual connections: shortcut layers and *add* back in
Skip connections: shortcut layers and *concatenate* back in

Normalization

No agreement on what they do and how they help, just *that* they help
Observation: Poorly calibrated dynamic range of activities (different variances) \Rightarrow poor backpropagation of errors (vanishing gradients)
Batch Normalization

Motivation: have strong dependencies between weights in layers, want to find a suitable learning rate to get similar scale
Broadly used in vision
Dependence on batch size not suitable for all architectures (e.g. RNNs)

Normalize layer l activations for a batch $\mathbf{I} \subseteq [1 : s]$: $\boldsymbol{\mu}^l := \frac{1}{|\mathbf{I}|} \sum_{i \in \mathbf{I}} \mathbf{z}^l[i]$
 $\sigma_i^l := \sqrt{\delta + \frac{1}{|\mathbf{I}|} \sum_{i \in \mathbf{I}} (\mathbf{z}_i^l[i] - \mu_i^l)^2}$
 $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ are functions of the weights: can be differentiated

Normalize (z-score): $\tilde{\mathbf{z}}_i^l := \frac{\mathbf{z}_i^l - \boldsymbol{\mu}_i^l}{\sigma_i^l}$

Regain representational power / expressivity: $\hat{\mathbf{z}}_i^l = \alpha_i^l \tilde{\mathbf{z}}_i^l + \beta_i^l$
Ideally would do BN over whole dataset, but cannot for computational reasons
Weight Normalization
Normalize by the inner product of the weights instead of sd of activations
Easier to use in inference / test mode than BN
Permits distributed optimization as the normalization doesn't include dependencies on batch samples

Layer Normalization
Used in NLP
Fix data point but use population average in layer as reference:

$\boldsymbol{\mu}^l[t] := \frac{1}{m^l} \sum_{i=1}^{m^l} \mathbf{z}_i^l[t]$
 $\sigma^l[t] := \sqrt{\delta + \frac{1}{m^l} \sum_{i=1}^{m^l} (\mathbf{z}_i^l[t] - \boldsymbol{\mu}^l[t])^2}$
Rest same as BN
BN, LN and WN are scale-invariant

Regularization

Intended to lower generalization error but not the training error
Norm-based
Regularized objective $\mathcal{R}_\Omega(\boldsymbol{\theta}, \mathcal{S}) = \mathcal{R}(\boldsymbol{\theta}, \mathcal{S}) + \Omega(\boldsymbol{\theta})$
 L_2 /Frobenius norm for deep networks
 $\Omega(\boldsymbol{\theta}) = \frac{1}{2} \sum_{l=1}^L \boldsymbol{\mu}^l \|\boldsymbol{\theta}^l\|_F^2$, $\boldsymbol{\mu}^l \geq 0$
Weight decay (a.k.a L_2 regularization)
 $\nabla \Omega = \boldsymbol{\mu} \boldsymbol{\theta}$ or $\frac{\partial \Omega}{\partial \theta_{ij}} = \mu \theta_{ij}^l$
Gradient descent update: $\boldsymbol{\theta}_{k+1} = (1 - \eta \mu) \boldsymbol{\theta}_k - \eta \nabla \mathcal{R}(\boldsymbol{\theta}_k)$

favors weights of small magnitude (shrinkage). Network behavior won't change much, avoiding learning local noise from data.
 $\boldsymbol{\theta}$ still moves in the direction of the gradients, but also shrinks (if η is small enough s.t. $1 - \eta > 0$). Forces network to learn only features which are seen often across the training set.
 \Rightarrow shrink weights where gradient vanishes!
Quadratic Taylor approximation of \mathcal{R} around optimal $\boldsymbol{\theta}^*$ + first order optimality condition of \mathcal{R}_Ω : $\boldsymbol{\theta} = \mathbf{Q}(\boldsymbol{\Lambda} + \boldsymbol{\mu} \mathbf{I})^{-1} \boldsymbol{\Lambda} \mathbf{Q}^T \boldsymbol{\theta}^*$, $\mathbf{H} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T$,
 $\lambda_i \gg \mu$: vanishing effect \Rightarrow very small regularization and $\theta_i \approx \theta_i^*$
 $\lambda_i \ll \mu$: shrinking effect \Rightarrow large regularization and $\theta_i \approx 0$
Take-away: Small regularization in directions of large eigenvalues, which correspond to directions with a lot of curvature (steep), and vice versa
Ridge: Linear regression + L_2 -regularization = ridge regression
 $\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X} + \mu \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$
Constrained optimization view: solve $\min_{\boldsymbol{\theta}: \|\boldsymbol{\theta}\| \leq r} \mathcal{R}(\boldsymbol{\theta})$, with projected GD
 $\boldsymbol{\theta}(k+1) = \Pi_r[\boldsymbol{\theta}(k) - \eta \nabla \mathcal{R}(\boldsymbol{\theta}(k))]$, $\Pi_r(\mathbf{v}) := \min\{1, \frac{r}{\|\mathbf{v}\|}\} \mathbf{v}$
Benefit: Constraints do not affect initial learning where weights are small, only become active once weights are large

Early stopping

Stop learning after small number of iterations. Rely on validation on data. If we can choose k, η s.t.
 $(1 - \eta \boldsymbol{\Lambda})^k \stackrel{1}{\approx} \boldsymbol{\mu}(\boldsymbol{\Lambda} + \boldsymbol{\mu} \mathbf{I})^{-1}$, which for $\eta \lambda_i \ll 1$, $\lambda_i \ll \mu$ can be achieved in $k = \frac{1}{\eta \mu}$ steps. This indicates that a stronger regularization corresponds to stopping training earlier and vice versa.

Ensemble methods

Bagging: Bootstrap samples \mathcal{S}_k^b : sample r times from \mathcal{S} with replacement (on average 2/3 distinct examples). Train model on each sample, then average model output probs $p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}^k)$:
 $p(\mathbf{y} | \mathbf{x}) = \frac{1}{K} \sum_{k=1}^K p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}^k)$
Almost always beneficial but expensive (not often used for NNs)
Knowledge Distillation - best of both worlds
Idea: Transfer knowledge from a complex [ensemble] model (source) into a simpler one (target). Useful if complex model is expensive to evaluate \Rightarrow can deploy on less powerful hardware
Train a simple model to learn the soft outputs of a trained complex model, using a high value of the temperature parameter in softmax

Dropout
Idea: Randomly drop subsets of units for better robustness
Keep probability π_i^l for unit i in layer l . All models share same weights
Realizes an ensemble: $p(\mathbf{y} | \mathbf{x}) = \sum_{\mathbf{Z}} p(\mathbf{Z}) p(\mathbf{y} | \mathbf{x}, \mathbf{Z})$, \mathbf{Z} zeroing mask
To predict, can sample and average. But to avoid sampling blow-up can use heuristic:
Weight Rescaling to avoid sampling 10-20 times, set $\tilde{\theta}_{ij}^l \leftarrow \pi_j^{l-1} \theta_{ij}^l$

Can be shown to be a (sometimes exact) approximation to a geometrically averaged ensemble

Data Augmentation

Powerful trick to generate larger training set & improve robustness
Generate virtual examples by applying transformations τ to each training example $(\mathbf{x}, \mathbf{y}) \mapsto (\tau(\mathbf{x}), \mathbf{y})$. PCA, cropping, resizing, rotations, reflections, etc.
Can inject noise to inputs, weights (regularizing effect) and targets (create soft targets as labels might be noisy)

Task Augmentation

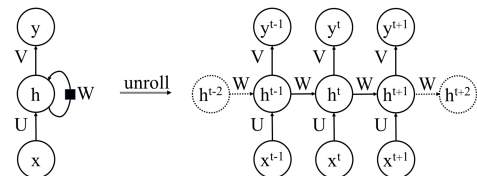
Pre-training: Pre-train parts of model on more generic task, where training data is cheaper. E.g. pre-trained word embeddings or image feature maps. After this: fine-tune
Semi-supervised learning: when not all data is labeled, create a prediction problem for unlabeled data. Define a generative model with corresponding log-likelihood. Generally more effective and more expensive than pre-training. Ex: predict relative position of image patch, predict coloring of gray-scale image

Multi-task learning: share (lower-level) representations across tasks and learn these jointly

Recurrent Neural Networks

Simple RNNs

Given a sequence of observations $\mathbf{x}^1, \dots, \mathbf{x}^s$, not iid
Discrete time evolution of hidden state space sequence:
 $\mathbf{h}^t = F(\mathbf{h}^{t-1}, \mathbf{x}^t; \boldsymbol{\theta})$ \Rightarrow Markov property and time-invariant
 $\mathbf{h}^t = F(\mathbf{h}^{t-1}, \mathbf{x}^t; \boldsymbol{\theta}) = \sigma \circ \bar{F}(\mathbf{h}^{t-1}, \mathbf{x}^t; \boldsymbol{\theta})$, $\bar{F}(\mathbf{h}, \mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\mathbf{h} + \mathbf{U}\mathbf{x} + \mathbf{b}$
Optionally produce outputs: $\mathbf{y}^t = H(\mathbf{h}^t; \boldsymbol{\theta})$, $H(\mathbf{h}, \boldsymbol{\theta}) := \sigma(\mathbf{V}\mathbf{h} + \mathbf{c})$



Hidden state can be thought of as a noisy memory, compresses relevant aspects of sequence: $(\mathbf{x}^1, \dots, \mathbf{x}^{t-1}) \mapsto \mathbf{h}^t$ (conceptually)
For any fixed length s , the unrolled recurrent net corresponds to a feedforward net with s hidden layers. Difference to MLP: sharing of parameters and inputs processed sequentially.

Backpropagation (through time :P): sum over time steps
With $\hat{\sigma}_i^t := \sigma'(\bar{F}_i(\mathbf{h}^{t-1}, \mathbf{x}^t))$:

$$\frac{\partial \mathcal{R}}{\partial w_{ij}} = \sum_{t=1}^s \frac{\partial \mathcal{R}}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial w_{ij}} = \sum_{t=1}^s \frac{\partial \mathcal{R}}{\partial h_i^t} \cdot \hat{\sigma}_i^t \cdot h_j^{t-1}$$
$$\frac{\partial \mathcal{R}}{\partial u_{ik}} = \sum_{t=1}^s \frac{\partial \mathcal{R}}{\partial h_i^t} \cdot \frac{\partial h_i^t}{\partial u_{ik}} = \sum_{t=1}^s \frac{\partial \mathcal{R}}{\partial h_i^t} \cdot \hat{\sigma}_i^t \cdot x_k^t$$

Example with setting: $y_{1:T}$ ground-truth outputs and $\hat{y}_{1:T}$ predictions, $a_t = F(x_t, h_{t-1}, y_{t-1}; \boldsymbol{\theta})$, $h_t = \sigma(a_t)$, $\hat{y}_t = G(h_t, \phi)$, $L_t = H(y_t, \hat{y}_t)$, $L = \sum_{t=1}^T L_t$. Derivative of L wrt $\boldsymbol{\theta}$:

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = \sum_{t=1}^T \frac{\partial L_t}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial \boldsymbol{\theta}} = \sum_{t=0}^T (\prod_{j=t}^{i+1} \frac{\partial h_j}{\partial a_j} \frac{\partial a_j}{\partial h_{j-1}}) \frac{\partial h_i}{\partial a_i} \frac{\partial a_i}{\partial \boldsymbol{\theta}}$$

Exploding/Vanishing gradients: let output in last step $\mathbf{y} = \mathbf{y}^s$
Then $\nabla_{\mathbf{x}^t} \mathcal{R} = [\prod_{j=t+1}^s \mathbf{W}^T \mathbf{S}(\mathbf{h}^j)] \cdot \mathbf{J}_H \cdot \nabla_{\mathbf{y}} \mathcal{R}$, $\mathbf{S}(\mathbf{h}^t) = \text{diag}(\hat{\sigma}_1^t, \dots, \hat{\sigma}_{n_h}^t)$
Take spectral norm, and use $\|\mathbf{A}\mathbf{B}\|_2 \leq \|\mathbf{A}\|_2 \cdot \|\mathbf{B}\|_2$. If $\sigma_{\max}(\mathbf{W}) < 1$:

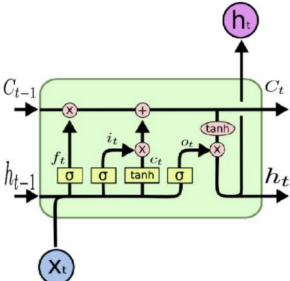
$$\|\nabla_{\mathbf{x}^t} \mathcal{R}\| \leq \sigma_{\max}(\mathbf{W})^{s-t} \cdot \|\mathbf{J}_H \cdot \nabla_{\mathbf{y}} \mathcal{R}\| \xrightarrow{(s-t) \rightarrow \infty} 0$$

Conversely may explode if $\sigma_{\max}(\mathbf{W}) > 1$
Bi-directional network: reverse order $\mathbf{g}^t = G(\mathbf{x}^t, \mathbf{g}^{t+1}; \boldsymbol{\theta})$
Interweave the two hidden state sequences

Deep recurrent network: hierarchical hidden states of l layers
 $\mathbf{h}^{t,l} = F(\mathbf{h}^{t,l-1}, \mathbf{x}^t; \boldsymbol{\theta})$
 $\mathbf{h}^{t,l} = F(\mathbf{h}^{t-1,l}, \mathbf{h}^{t,l-1}; \boldsymbol{\theta})$, $l = 2, \dots, L$

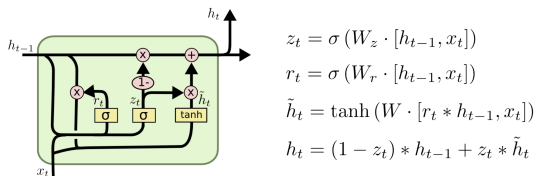
Memory Units

Addressing the problem of vanishing/exploding gradients
Gated units to learn long-term dependencies
Difficult to understand what units learn, resource-hungry and slow in learning
LSTM



Content \mathbf{C} has similar effect as ResNet for vanishing gradients
Forget gate: $f_t = \sigma(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f)$
Input to memory:

$i_t = \sigma(\mathbf{W}_i \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i)$, $\tilde{\mathbf{C}}_t = \tanh(\mathbf{W}_{\tilde{\mathbf{C}}} \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_{\tilde{\mathbf{C}}})$
Updating memory: $\mathbf{C}_t = f_t * \mathbf{C}_{t-1} + i_t * \tilde{\mathbf{C}}_t$
Output gate:
 $o_t = \sigma(\mathbf{W}_o [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o)$, $h_t = o_t * \tanh(\mathbf{C}_t)$
Gates: input gate, output gate and forget gate
Gated Recurrent Unit (GRU)



Gates: update gate and reset gate
Here h is both state and content, computed as a convex combination of old and new information
Learning sequences - Seq2Seq

Goal: learn $p(\mathbf{y}^1 : \mathbf{T} | \mathbf{x}^1 : \mathbf{T}) \approx \prod_{t=1}^T p(\mathbf{y}^t | \mathbf{x}^1 : \mathbf{t}, \mathbf{y}^1 : \mathbf{t}-1)$
Naive implementation: $p(\mathbf{y}^t)$ depends on $\mathbf{y}^1 : \mathbf{t}-1$ only through \mathbf{h}^t \Rightarrow remedy by feeding back previous outputs!

Teacher forcing: When training, compute loss on predicted output but feed back in correct output. During prediction feed back predicted outputs. Improves learning BUT gives exposure bias

Encoder-Decoder model:

Encoder: $(\mathbf{x}^1, \dots, \mathbf{x}^T) \mapsto \mathbf{z}$, $\mathbf{z} = \mathbf{h}^T$ (RNN)
Decoder $\mathbf{z} \mapsto (\mathbf{y}^1, \dots, \mathbf{y}^S)$ (RNN with output feedback)
Can also be used for image captioning with CNN encoder

Attention

Instead of compressing the whole sequence into a vector, determine where to look

Self attention

Gating function: Given query $\boldsymbol{\xi} \in \mathbb{R}^m$ and values $\mathbf{x}^t \in \mathbb{R}^m$,

$$f_\phi(\boldsymbol{\xi}, (\mathbf{x}^1, \dots, \mathbf{x}^s)) = \frac{1}{\sum_j \exp[\phi(\boldsymbol{\xi}, \mathbf{x}^j)]} \left[\frac{\exp[\phi(\boldsymbol{\xi}, \mathbf{x}^1)]}{\exp[\phi(\boldsymbol{\xi}, \mathbf{x}^s)]} \right]$$

Simplest choice for similarity function when $m = n$: $\phi(\boldsymbol{\xi}, \mathbf{x}) = \boldsymbol{\xi} \cdot \mathbf{x}$
Transfer function: $F(\boldsymbol{\xi}, (\mathbf{x}^1, \dots, \mathbf{x}^s)) = [\mathbf{x}^1 \dots \mathbf{x}^s] \cdot f_\phi(\boldsymbol{\xi}, (\mathbf{x}^1, \dots, \mathbf{x}^s))$
i.e. computes convex combination of inputs wrt attention weights

RNN with attention (Seq2Seq): combining ideas of encoding information relevant for the future (RNNs) and selecting what is relevant in retrospective (attention)

Attend to hidden state sequence $(\mathbf{h}_1^e, \dots, \mathbf{h}_n^e)$ of encoder with query $\boldsymbol{\xi}^i$ produced as hidden states by decoder. Use this as input into decoder, which produces hidden states $(\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^t)$ and output sequence $(\mathbf{y}^1, \dots, \mathbf{y}^t)$

Memory Networks

Recurrent attention model over possibly large external memory
Recursive associative recall: Given query \mathbf{q} (e.g. question), find best matching memory cell i , use its content \mathbf{m}_i and \mathbf{x} to generate new query - repeat

Transformers

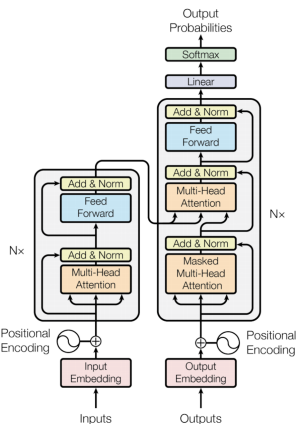
Attention is all you need

Key-Value attention map: given KV pairs $(\mathbf{x}^i, \mathbf{z}^i)$,
 $F(\boldsymbol{\xi}, ((\mathbf{x}^1, \mathbf{z}^1), \dots, (\mathbf{x}^s, \mathbf{z}^s))) = [\mathbf{z}^1 \dots \mathbf{z}^s] \cdot f(\boldsymbol{\xi}, (\mathbf{x}^1, \dots, \mathbf{x}^s))$
 \Rightarrow keys determine where to look, values what features to extract
Similarity function: $\phi(\boldsymbol{\xi}, \mathbf{x}) = \frac{\boldsymbol{\xi} \cdot \mathbf{x}}{\sqrt{m}}$, $\boldsymbol{\xi}, \mathbf{x} \in \mathbb{R}^n$

Assuming entries sampled from standard Gaussian, this gives standard scaling. Good since softmax can be sensitive to large input values (which kills gradient and slows down learning)

Multi-headed attention: $G(\boldsymbol{\xi}, (\mathbf{x}^t, \mathbf{z}^t)_{t=1}^s) = \mathbf{W} \begin{bmatrix} F_1(\boldsymbol{\xi}, (\mathbf{x}^t, \mathbf{z}^t)) \\ \vdots \\ F_h(\boldsymbol{\xi}, (\mathbf{x}^t, \mathbf{z}^t)) \end{bmatrix}$

$F_j(\boldsymbol{\xi}, (\mathbf{x}^t, \mathbf{z}^t)) = F(\mathbf{W}_j^q \boldsymbol{\xi}, (\mathbf{W}_j^k \mathbf{x}^t, \mathbf{W}_j^v \mathbf{z}^t))$
Alt: $\mathbf{Y} = \text{softmax}(\mathbf{X} \mathbf{W}_q \mathbf{W}_k^T \mathbf{X}^T) \mathbf{X} \mathbf{W}_v, \mathbf{W}_{\{q,k,v\}} \in \mathbb{R}^{d \times hd}, \mathbf{X}, \mathbf{Y} \in \mathbb{R}^{L \times d}$



- 1. Input embedding** Accounts for meaning
- 2. Positional Encodings** Accounts for position in the input sentence. Normally use sinusoidal functions. For position t and feature k :

$$p_{tk} = \begin{cases} \sin(t\omega_k) & k \text{ even} \\ \cos(t\omega_k) & k \text{ odd} \end{cases} \quad \omega_k = C^{k/n}, C = 10000$$
- 3. Creating Masks** In decoder, to prevent peaking ahead.
- 4. Multi-Head Attention Layer** Split embedding into N heads and perform attention on each. Three types: Self, masked and mixed.
- 5. Feed-Forward layer**

Applications in NLP

Construct word embeddings that reflect the context

ELMo: Input fixed embeddings on character level, build word representations with CNNs, then stack left-to-right LSTM and right-to-left LSTM. Output (for each word) is convex combination of hidden states over layers. Out-of-vocabulary (not restricted to fixed vocabulary). Train in task-specific manner.

BERT: Input sub-word unit embeddings. Leverage attention - take query and retrieve keys & values from words in context. Insight: do not need to learn a language model. Train by word masking task - predict missing word in text. Also do next sentence prediction in pre-training. Can fine-tune for specific downstream task. **BERTology** - representations hierarchical, struggles with negation, incomplete syntactic knowledge etc

GPT-n: Few, one or zero shot learning - add task description & examples to working memory and let model do the rest.

Theory of DNNs

VC Theory

Shattering coefficient: maximum number of ways in which n points can be classified by function class \mathcal{F}
 $\Rightarrow S_{\mathcal{F}}(n) = \sup_{\{x_1, \dots, x_n\}} |\{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\}| \leq 2^n$
 Say that \mathcal{F} **shatters** a set of n points if $S_{\mathcal{F}}(n) = 2^n$

VC-dim The VC-dim of \mathcal{F} is $n \iff$ there is a set of size n that is shattered by \mathcal{F} , and **no** set of size $n+1$ is shattered by \mathcal{F} .

$VC\text{-dim}(\mathcal{F}) \leq \log_2(|\mathcal{F}|)$

Vapnik-Chervonenkis Theorem: For any $\delta > 0$, with probability at least $1 - \delta$:

$$\forall f \in \mathcal{F}, \quad \mathcal{R}(f) \leq \mathcal{R}_n(f) + 2\sqrt{2 \cdot \frac{\log S_{\mathcal{F}}(2n) + \log \frac{2}{\delta}}{n}}$$

PAC Bayes Bounds

Donsker's Theorem: For any $P \gg Q$, P -measurable function ϕ
 $\mathbb{E}_Q[\phi] \leq \text{KL}(Q||P) + \log \mathbb{E}_P[e^{\phi}]$

McAllister theorem \forall fixed P and any $Q, \epsilon \in (0, 1)$ w.prob $\geq \epsilon$

over sample set S : $\mathbb{E}_Q[e_f] - \mathbb{E}_Q[e_{f^*}] \leq \sqrt{\frac{2}{|S|} \left[K L(Q||P) + \ln \left(\frac{2\sqrt{|S|}}{\epsilon} \right) \right]}$

Interpretation for DNNs: Take $P = \mathcal{N}(\theta_0, \lambda I)$ to be prior over parameter space. $Q = \mathcal{N}(\theta, \text{diag}(\sigma^2))$ (learned from data) is a distribution over \mathcal{F} , where our function f is sampled from. Minimize bound together with (surrogate of) empirical risk.

A **stochastic neural network** is a network whose weights are

drawn from a distribution Q each time data is propagated through the network.

Learning on Graphs

Motivation: Data may not be sampled iid but could be graph-structured

Use-cases: Node classification, link prediction, generative modeling

Focus here: GCN from paper *Semi-Supervised Classification with Graph Convolutional Networks*, Kipf & Welling 2017

Idea: To make prediction on node, gather information from context in graph

Let $\mathbf{X} \in \mathbb{R}^{d \times n}$ be graph of n nodes with d features, then define:
 $\mathbf{X}^{l+1} = \sigma(\mathbf{W}^l \mathbf{X}^l \mathbf{Q})$, where $\mathbf{W}\mathbf{X}$ sums over dimensions and \mathbf{XQ} sums over data points / nodes (to learn from neighborhood)

Define degree-normalized matrix $\mathbf{Q} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$, where $\tilde{\mathbf{D}}$ is the diagonal degree matrix of $\tilde{\mathbf{A}} := \mathbf{A} + \mathbf{I}_n$ and \mathbf{A} is adjacency matrix. It holds that:

$$q_{ij} = \frac{a_{ij} + \delta_{ij}}{\sqrt{d_i d_j}}, \quad \tilde{d}_i = 1 + \sum_j A_{ij}, \quad \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

Activations from neighbors are mixed together with respective Q -weights: and $(\mathbf{X}^l \mathbf{Q})_{ij} = (\{x^l[1], \dots, x^l[n]\} \cdot Q)_{ij} = \sum_{k=1}^n x_{ik} q_{kj}$

$\Rightarrow (\mathbf{X}^l \mathbf{Q})_{\bullet j} = \sum_k q_{kj} \cdot x^l[k]$ (sum over neighborhood of node j)

Cyclic chain: For regular lattice graph, Q is Toeplitz matrix and thus related to convolutional layer with window size same as neighborhood size

Linear Shift Invariant Filter: linear function H over graph with adjacency matrix A such that $H(Ax) = A(Hx)$. H is A -shift invariant iff: $H = \theta_0 I + \theta_1 A + \dots + \theta_n A^n$

If we consider $H(\theta_0, \theta_1) = \theta_0 I + \theta_1 (D^{-\frac{1}{2}} \tilde{\mathbf{A}} D^{-\frac{1}{2}})$ and $\theta_0 = \theta_1$ then largest eigenvalue is 2, which can be a source of instabilities.

Motivates $I + D^{-\frac{1}{2}} \tilde{\mathbf{A}} D^{-\frac{1}{2}} \mapsto \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$, which has largest EV 1

Adversarial Robustness

Adversarial examples are small manipulations of input which change the model to change its prediction

Additive adversarial manipulation: Find $x \in \mathbb{R}^d$ s.t. $k(x+r) \neq k(x)$ and $x+r$ similar to x . As proxy for *similar*, use l_p -norm robustness $\|r\|_p$ should be small

Two approaches to measure robustness to adversarial perturbations: unconstrained & constrained perturbations

Unconstrained Perturbations - DeepFool

Find smallest perturbation in l_p -norm that induces mistake

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}^C$ be the logits, i.e. outputs before softmax, and $\text{wlog } k(x) = 1$

Disadvantage of approach: Typically get close to decision boundary which may be exploited to detect them (can e.g. train a classifier to detect them given some confidence threshold)

$$\hat{\rho}_{\text{adv}} = \frac{1}{n} \sum_{i=1}^n \frac{\|r(x_i)\|_2}{\|x_i\|_2}$$

DeepFool - Binary case: $\arg \min_{\|r\|_p} f_2(x+r) - f_1(x+r) > 0$

First-order Taylor approximation around x , rewrite constraint: $(\nabla_x f_1(x) - \nabla_x f_2(x))^T r + f_1(x) - f_2(x) > 0 \Rightarrow$ linearized classifier
 Solve iteratively until $k(x+r_1 + \dots + r_n) \neq k(x)$
 Optimum in each step, $1/p + 1/p^* = 1$:

$$r^* = \frac{f_1(x) - f_2(x)}{\|\nabla_x f_1(x) - \nabla_x f_2(x)\|_{p^*}^{p^*}} \cdot |\nabla_x f_1(x) - \nabla_x f_2(x)|^{p^*-1} \odot \text{sgn}(\nabla_x f_1(x) - \nabla_x f_2(x))$$

$\nabla_x f_2(x)$

DeepFool - General case: Optimize wrt r_i for each class and update in the direction of smallest r_i
 $x_0 = x$

$$x_k = x_{k-1} + \min_{j \neq k(x)} \left\{ \frac{|h_j(x_{k-1})|}{\|\nabla_x h_j(x_{k-1})\|_{p^*}^{p^*}} \cdot |\nabla_x h_j(x_{k-1})|^{p^*-1} \odot \text{sgn}(\nabla_x h_j(x_{k-1})) \right\}, \text{ where } h_j(x_{k-1}) = f_j(x_{k-1}) - f_{k(x)}(x_{k-1})$$

Constrained Perturbations - PGD

Given $\epsilon > 0$ find perturbation r with $\|r\|_p \leq \epsilon$

Disadvantage of approach: Perturbation found can be larger

than the minimal one within ball

Projected Gradient Descent / Ascent

$$x_0 \sim \mathcal{U}(B_{\epsilon}^p(x))$$

$$x_k = \Pi_{B_{\epsilon}^p(x)}(x_{k-1} + \alpha \arg \max_{v_k: \|v_k\|_p \leq 1} v_k^T \nabla_x l(y, f(x_{k-1})))$$

where $r_k = x_k - x$, $\Pi_{B_{\epsilon}^p(x)}(\bar{x}) = \arg \min_{x^* \in B_{\epsilon}^p(x)} \|\bar{x} - x^*\|_2$ and

$$v^* = \arg \max_{v: \|v\|_p \leq 1} v^T z = \frac{\text{sign}(z) |z|^{p^*-1}}{\|z^*\|_{p^*}^{p^*-1}}$$

Adversarial Training: Most effective method for improving robustness. Can combine with e.g. PGD:

$$\min_{\theta} \max_{\|r\|_p \leq \epsilon} (f_{\theta}(x+r), y)$$

Accuracy on clean data degrades \rightarrow more data

Network overfits to attacks used in training \rightarrow use many iterations for PGD

Takes longer to train \rightarrow explicit regularization

Mode of failure 1: gradient wrt input norm becomes too small

Mode of failure 2: gradient wrt input becomes noisy \rightarrow not find closest perturbation

Benefits beyond robustness and security

Interpretability - robust networks learn more human-interpretable features

Transfer learning - robust networks generalize better across domains

Generative modeling - robust networks perform better in generative tasks

Adversarial training and operator norm regularization

Linearization: $f(x + \Delta x) \approx f(x) + J_f(x) \Delta x$

$$\Rightarrow \frac{\|f(x + \Delta x) - f(x)\|_2}{\|\Delta x\|_2} \approx \frac{\|J_f(x) \Delta x\|_2}{\|\Delta x\|_2} \leq \sigma(J_f(x)) =: \max_{v: \|v\|_2 \leq 1} \|J_f(x) v\|_2$$

So from a robustness perspective, want small operator norm

Data-independent spectral norm reg.: $\sigma(J_f(x)) \leq \prod_{i=1}^L \sigma(W^i)$

Generalizes from train to test set but can be arbitrarily loose

Data-dependent: $\min_{\theta} \mathbb{E}_{(x,y) \sim \tilde{p}} [l(y, f(x)) + \frac{1}{2} \sigma(J_f(x))^2]$,

where $\sigma(J_f(x))$ is computed by the power method

l_p -norm constrained PGA based adversarial training with an l_q -norm loss on the logits of clean and perturbed inputs is equivalent to data-dependent (p, q) operator norm regularization.

VAEs

Linear Factor Analysis

Latent variable prior $z \sim \mathcal{N}(0, I)$, $z \in \mathbb{R}^m$

Linear observation model for $x \in \mathbb{R}^n$: $x = \mu + Wz + \eta$,

$$\eta \sim \mathcal{N}(0, \Sigma), \quad \Sigma := \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$$

Can be shown that: $x \sim \mathcal{N}(\mu, WW^T + \Sigma)$

Non-identifiability, for orthogonal Q : $(WQ)(WQ)^T = WW^T$

Posterior inference:

$$\mu_{z|x} = W^T (WW^T + \Sigma)^{-1} (x - \mu)$$

$$\Sigma_{z|x} = I - W^T (WW^T + \Sigma)^{-1} W$$

MLE: $\theta = (\mu, W) \xleftarrow{\text{max}} \log p(x; \mu, W)$, has no closed form solution

Probabilistic PCA: $\mathcal{W}(x) = \sigma^2 I$

Variational Autoencoders

Generalize factor analysis with depth

Noise variable: $z \sim \mathcal{N}(0, I)$

Density unavailable explicitly, so can't do MLE:

$$\theta \xleftarrow{\text{max}} \sum_{i=1}^n \log p(x[i]; \theta)$$

Max. ELBO: $\log p(x; \theta) \geq \mathbb{E}_{z \sim q(z|x)} [\log p(x|z; \theta)] - \text{KL}(q(z|x)||p(z))$

Can think of $q(z|x)$ as posterior ($= p(z|x)$), restricted to (typically Gaussian) variational family:

$\mathcal{N}(\mu(x), \Sigma(x))$, $\Sigma(x) = \text{diag}(\sigma_1^2(x), \dots, \sigma_n^2(x))$ (output of encoder)

Optimizing over q involves gradients of expectations \Rightarrow stochastic backpropagation / re-parameterization trick:

$$z \sim \mathcal{N}(\mu, \Sigma) \Leftrightarrow z = \mu + \Sigma^{1/2} \eta, \quad \eta \sim \mathcal{N}(0, I)$$

\Rightarrow can now backpropagate wrt μ and $\Sigma^{1/2}$!

Generative Models

Density Estimation

Learn a parametrized model $p_{\theta}(x)$ to be indistinguishable from true generative process $p(x)$

1 - Prescribed model: Density explicitly specified (and accessible). Can do MLE. Challenge: Normalizing model

2 - Implicit model: Transformation of densities, can describe density through integration by substitution

Integration by substitution - change of variables formula:

$x = F(z)$, $p_X(x) = p_Z(F^{-1}(x)) |\det(\partial F)|^{-1}$. p_X pushforward of p_Z

Normalizing Flows

Bijections F which are convenient to compute, invert and calculate $|\det(\partial F)|$

$$F = F_L \circ \dots \circ F_1, \quad F^{-1} = F_1^{-1} \circ \dots \circ F_L^{-1}$$

$$\Rightarrow \det(\partial F) = \prod_{i=1}^L \det(\partial F_i \circ F_{1:i-1}), \quad \det(\partial F^{-1}) = \det(\partial F)^{-1}$$

Log-likelihood: $\log p(x|z) = -\sum_{i=1}^L \log |\det(\partial F_i \circ F_{1:i-1})|$

Linear flow: $F(x) = Ax + b \Rightarrow \mathcal{O}(n^3)$ to compute determinant and inverse. If A diagonal, simple & efficient but not powerful.

If A triangular: $\det(A) = \prod_i A_{ii}$, but A^{-1} expensive...

Invertible linear time flows: $F(z) = z + u\sigma(w \cdot z + b)$

Determinant: $|\det(\partial F)| = |1 + \sigma' u^T w|$

In general, normalizing flows not powerful enough as Jacobian condition is too restrictive \Rightarrow Likelihood-free methods

Generative Adversarial Networks (GANs)

Derive training signal for generator G from classifier D that discriminates data from model-generated samples

For sample x and model label y : $\tilde{p}_{\theta}(x, y) = \frac{1}{2} (yp(x) + (1-y)p_{\theta}(x))$

\Rightarrow balanced binary classification problem

Optimal discriminator - Bayes optimal classifier:

$$q_{\theta}(x) = \frac{p(x)}{p(x) + p_{\theta}(x)} = P(y = 1|x)$$

Train generator by minimizing log-likelihood (=JS for BOC)

$I^*(\theta) = \mathbb{E}_{\tilde{p}_{\theta}} [\log q_{\theta}(x) + (1-y) \log(1-q_{\theta}(x))] = \dots = JS(p, p_{\theta}) - \log 2$

Bayes optimal classifier not accessible, define classification model: $q_{\phi}: x \mapsto [0, 1]$, $\phi \in \Phi$

It holds that: $I^*(\theta) \geq \sup_{\phi \in \Phi} I(\theta, \phi)$ (given a generator, the log-likelihood with the BOC is an upper bound for the log-likelihood with any other classifier)

General objective: $I(\theta, \phi) = \mathbb{E}_{\tilde{p}_{\theta}} [y \log q_{\phi}(x) + (1-y) \log(1-q_{\phi}(x))]$

Saddle-point problem: $\theta^* = \arg \min_{\theta \in \Theta} [\sup_{\phi \in \Phi} I(\theta, \phi)]$

Inner sup is impractical \Rightarrow SGD heuristic:

$$\theta^{t+1} = \theta^t - \eta \nabla_{\theta} I(\theta^t, \phi^t), \quad \phi^{t+1} = \phi^t + \eta \nabla_{\phi} I(\theta^{t+1}, \phi^t)$$

Alternatively, with $D(\cdot)$ the probability discriminator assigns that sample is from real data:

$$\min_G \max_D \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p_Z(z)} [\log(1 - D(G(z)))]$$

Mode collapse (Helvetica scenario): Generator collapses too many values of z to similar values of x and does not capture the diversity in $p(x)$. Avoid by not training G too much without updating D . Can lead to good-looking samples without diversity

Vanishing gradient: If the discriminator becomes too strong the gradient wrt the parameters of the generator vanishes ($\rightarrow 0$), so there will be no learning signal for the generator and it will not improve. Can lead to noisy/blurry samples but good variability \Rightarrow trade-off

Why does minimizing Jensen-Shannon divergence work?

Definition: $JS(Q||P) = \frac{1}{2} \text{KL}(P||\frac{P+Q}{2}) + \frac{1}{2} \text{KL}(Q||\frac{P+Q}{2})$

In MLE we minimize the forward KL-divergence $\text{KL}(P||Q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx$, which leads $q(x)$ to spread its mass over the whole support of $p(x)$, which is why samples from e.g. VAEs tend to be blurry.

If we minimize the reverse KL-divergence $\text{KL}(Q||P)$, $q(x)$ can be set to zero when $p(x) > 0$, which leads to a $q(x)$ that focuses on certain mode(s) of $p(x)$.

Want something in between where whole support of $p(x)$ is covered while $q(x) = 0$ where $p(x) = 0 \Rightarrow$ **JS-divergence**!