

1. History

1.1 Perceptron

Threshold Unit

$f[w, b](x) = \text{sign}(x \cdot w + b)$ where $x \cdot w := \sum_{i=1}^n x_i w_i$

Decision Boundary

$$x \cdot w + b = 0 \Leftrightarrow \frac{x \cdot w}{\|w\|} + \frac{b}{\|w\|} = 0$$

$$x \cdot w - b = \begin{pmatrix} x \\ -1 \end{pmatrix} \cdot \begin{pmatrix} w \\ b \end{pmatrix} =: \tilde{x} \cdot \tilde{w}, \quad \tilde{x}, \tilde{w} \in \mathbb{R}^{n+1}$$

Geometric Margin

$$\gamma[w, b](x, y) := \frac{y(x \cdot w + b)}{\|w\|}$$

Maximum Margin Classifier

$(w^*, b^*) \in \operatorname{argmax}_{w, b} \gamma[w, b](S)$ with $\gamma[w, b](S) := \min_{(x, y) \in S} \gamma[w, b](x, y)$

Perceptron Learning

if $f[w, b](x) \neq y$: update $w \leftarrow w + yx$, and $b \leftarrow b + y$
 $w_0 \in \operatorname{span}(x_1, \dots, x_s) \Rightarrow w_t \in \operatorname{span}(x_1, \dots, x_s) (\forall t)$

Convergence

$\exists w, \|w\| = 1$, that $\gamma[w](S) = \gamma > 0 \Rightarrow w_t \cdot w \geq t\gamma$.

$$R = \max_{x \in S} \|x\| \Rightarrow \|w_t\| \leq R\sqrt{t}$$

$$\cos \angle(w, w_t) = \frac{w \cdot w_t}{\|w_t\|} \geq \frac{t\gamma}{\sqrt{t}R} = \frac{\sqrt{t}\gamma}{R} \leq 1 \Rightarrow t \leq \frac{R^2}{\gamma^2}$$

Covers Theorem

$C(s+1, n) = 2 \sum_{i=0}^{n-1} \binom{s}{i}$
 $C(S, n)$: Number of ways to separate S with n dimensions.
 $C(s, n) = 2s$ for $s \leq n$
Phase transition at $s = 2n$. For $s > 2n$ empty version space is the exception.

1.2 Hopfield Networks

Hopfield Model

$E(X) = -\frac{1}{2} \sum_i \sum_{j \neq i} w_{ij} X_i X_j + \sum_i b_i X_i$ where $X_i \in \{-1, +1\}$
 $w_{ij} = w_{ji} \ (\forall i, j)$, $w_{ii} = 0 \ (\forall i)$: Interaction strengths

Hebbian Learning

$$x^t \in \{\pm 1\}^n \ (1 \leq t \leq s), \ w_{ij} = \frac{1}{n} \sum_{t=1}^s x_i^t x_j^t, \ W = \frac{1}{n} \sum_{t=1}^s x^t (x^t)^\top$$

2. Feedforward Networks

2.1 Linear Models

Linear regression

$$h[w](S) = \frac{1}{2s} \|Xw - y\|^2, \ \nabla h = 2X^\top Xw - 2X^\top y$$

Moore-Penrose inverse solution

$w^* = X^+ y \in \operatorname{argmin}_w h[w]$ where $X^+ := \lim_{\epsilon \rightarrow 0} (X^\top X + \epsilon I)^{-1} X^\top$

SGD update

$w_{t+1} := w_t + \eta(y_{i_t} - w_t^\top x_{i_t})x_{i_t}$ with $i_t \stackrel{\text{iid}}{\sim} \text{Uniform}(1, \dots, s)$

Gaussian noise model

$y_i = w^\top x_i + \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$. Least squares \equiv neg log likelihood.

Ridge regression

$$h_\lambda[w] := h[w] + \frac{\lambda}{2} \|w\|^2, \ w^* = (X^\top X + \lambda I)^{-1} X^\top y$$

Logistic function

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

$$\sigma' = \sigma(1 - \sigma), \ \sigma'' = \sigma(1 - \sigma)(1 - 2\sigma)$$

Cross entropy loss

$$\ell(y, z) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z)) = -\log \sigma((2y - 1)z)$$

Logistic regression gradient

$$\nabla \ell_i = [\sigma(w^\top x_i) - y_i] x_i$$

2.2 Feedforward Networks

Generic feedforward layer

$$F : \mathbb{R}^{m(n+1)} \times \mathbb{R}^n \rightarrow \mathbb{R}^m, \ F[\theta](x) := \varphi(Wx + b), \ \theta := (W, b)$$

Composition of layers

$$G = F_L[\theta_L] \circ \dots \circ F_1[\theta_1] \text{ where } F_l[W^l, b^l](x) := \varphi_l(W^l x + b^l)$$

Layer activations

$$x^l := (F_l \circ \dots \circ F_1)(x) = F_l(x^{l-1}), \ x^0 = x, \ x^L = F(x)$$

Softmax function

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}, \ \text{softmax}(A)_{ij} = \frac{e^{A_{ij}}}{\sum_k e^{A_{ik}}}$$

$$\ell(y; z) = \left[-zy + \log \sum_j e^{z_j} \right] \frac{1}{\ln 2}$$

Residual layer

$$F[W, b](x) = x + [\varphi(Wx + b) - \varphi(0)], \text{ therefore } F[0, 0] = \text{id}$$

Skip connection

Concatenate previous layer back in

2.3 Sigmoid Networks

Sigmoid/Tanh activations

$$\sigma(z) = \frac{1}{1+e^{-z}}, \ \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = 2\sigma(2z) - 1, \ \tanh'(z) = 1 - \tanh^2(z)$$

Barron's Theorem

For f with finite $C_f := \int \|\omega\| |\hat{f}(\omega)| d\omega$, \exists MLP g with width m : $\int_B (f - g_m)^2 \mu(dx) \leq O(1/m)$

2.4 ReLU Networks

ReLU activation

$\varphi(z) := (z)_+ := \max\{0, z\}$. ReLU networks are universal function approximators.

Zaslavsky: Connected regions

$$R(H) \leq \sum_{i=0}^{\min\{n, m\}} \binom{m}{i} := R(m)$$

Montufar: Connected regions

$$R(m, L) \geq R(m) \lfloor \frac{m}{n} \rfloor^{n(L-1)} \ (L: \text{layers}, m: \text{width})$$

3. Gradient-Based Learning

3.1 Backpropagation

Parameter derivatives

$$\frac{\partial x_i^l}{\partial w_{ij}^l} = \varphi_i^l x_j^{l-1}, \ \frac{\partial x_i^l}{\partial b_i^l} = \varphi_i^l \text{ where } \varphi_i^l := \varphi'^l((w_i^l)^\top x^{l-1} + b_i^l)$$

Loss derivatives

$$\frac{\partial \eta}{\partial w_{ij}^l} = \delta_i^l \varphi_i^l x_j^{l-1}, \ \frac{\partial \eta}{\partial b_i^l} = \delta_i^l \varphi_i^l \text{ with } \delta_i^l = \frac{\partial \eta}{\partial x_i^l} \varphi_i^l$$

3.2 Gradient Descent

GD update & flow

$$\theta_{t+1} = \theta_t - \eta \nabla h(\theta_t), \text{ ODE: } \frac{d\theta}{dt} = -\nabla h(\theta)$$

L-smoothness

$$\begin{aligned} \|\nabla h(\theta_1) - \nabla h(\theta_2)\| &\leq L \|\theta_1 - \theta_2\| \ (\forall \theta_1, \theta_2) \\ \lambda_{\max}(\nabla^2 h) &\leq L, \ \ell''(x) \leq L \\ \ell(w) - \ell(w') &\leq \nabla \ell(w')^\top (w - w') + \frac{L}{2} \|w - w'\|_2^2 \end{aligned}$$

Polyak-Łojasiewicz

$$\frac{1}{2} \|\nabla h(\theta)\|^2 \geq \mu(h(\theta) - \min h) \ (\forall \theta)$$

Convergence rate

$\eta = 1/L \Rightarrow t = \frac{2L}{\epsilon^2} (h(\theta^0) - \min h)$ for ϵ -critical. With PL: $h(\theta^t) - \min h \leq (1 - \frac{\mu}{L})^t (h(\theta^0) - \min h)$

3.3 Acceleration and Adaptivity

Heavy ball momentum

$$\theta_{t+1} = \theta_t - \eta \nabla h(\theta_t) + \beta(\theta_t - \theta_{t-1})$$

Nesterov acceleration

$$\tilde{\theta}_{t+1} = \theta_t + \beta(\theta_t - \theta_{t-1}), \ \theta_{t+1} = \tilde{\theta}_{t+1} - \eta \nabla h(\tilde{\theta}_{t+1})$$

AdaGrad

$$\theta_{t+1}^i = \theta_t^i - \eta_t^i \partial_i h(\theta^t), \ \nu_t^i = \nu_{t-1}^i + [\partial_i h]^2, \ \eta_t^i = \frac{\eta}{\sqrt{\nu_t^i + \epsilon}}$$

Adam

$$g_t^i = \beta g_{t-1}^i + (1 - \beta) \partial_i h, \ \nu_t^i = \alpha \nu_{t-1}^i + (1 - \alpha) [\partial_i h]^2, \ \theta_{t+1}^i = \theta_t^i - \frac{\eta}{\sqrt{\nu_t^i + \epsilon}} g_t^i$$

RMSprop

Adam without momentum term (set $\beta = 0$)

3.4 SGD

SGD update & variance

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla h(\theta^t)(x_{i_t}, y_{i_t}), \quad V[\theta] = \frac{1}{s} \sum_{i=1}^s \|\nabla h(S) - \nabla h(x_i, y_i)\|^2 \\ \text{General: } &O(1/\sqrt{t}), \text{ Strongly convex: } O(\log t/t), \text{ Additionally smooth: } O(1/t) \end{aligned}$$

SGD convergence

General: $O(1/\sqrt{t})$, Strongly convex: $O(\log t/t)$, Additionally smooth: $O(1/t)$

3.5 Function properties

Convexity

$$\begin{aligned} \ell(\lambda w + (1 - \lambda)w') &\leq \lambda \ell(w) + (1 - \lambda) \ell(w'), \ \ell''(x) \geq 0 \ \forall x \\ \ell(w) &\geq \ell(w') + \nabla \ell(w')^\top (w - w') \ (\text{differentiable case}) \end{aligned}$$

Strong convexity

$$\ell(w) \geq \ell(w') + \nabla \ell(w')^\top (w - w') + \frac{\mu}{2} \|w - w'\|_2^2, \ \ell''(x) \geq \mu$$

4. Convolutional Networks

4.1 Convolutions

Convolution definition

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

Fourier property

$$\mathcal{F}(f * g) = \mathcal{F}(f) \cdot \mathcal{F}(g)$$

Discrete convolution

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

Cross-correlation

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

Toeplitz matrices

$$(f * g) = \text{Toeplitz-Matrix}(g) f$$

4.2 Convolutional Networks

Conventions

Padding: Add zeros around input. Stride: Step size of convolution.

Max-Pooling

Take maximum value in windows (size r)

ConvNets for Images

$$y[r][s, t] = \sum_u \sum_{s', t'} w[r, u][s', t'] \cdot x[u][s + s', t + t'] \ (r: \text{output channel}, u: \text{input channel})$$

Parameters count

$$D = \#r \cdot \#u \cdot \#s' \cdot \#t' \ (\text{channels} \times \text{window size})$$

4.3 NLP with ConvNets

Word embedding

$$w \mapsto x_w \in \mathbb{R}^n$$

Conditional log-bilinear model

$$P(\nu | w) = \frac{\exp[x_w^\top y_\nu]}{\sum_\mu \exp[x_w^\top y_\mu]}$$

$$h(\{x_w\}, \{y_\nu\}) = \sum_{(w, \nu)} \ell_{w, \nu}, \ \ell_{w, \nu} = -x_w^\top y_\nu + \ln \sum_\mu \exp[x_w^\top y_\mu]$$

Negative sampling

$$\tilde{\ell}_{w, \nu} = -\ln \sigma(x_w^\top y_\nu) - \beta \mathbb{E}_{\mu \sim D} \ln(1 - \sigma(x_w^\top y_\mu))$$

5. Recurrent Networks

5.1 Simple RNNs

Time evolution

$$z_t := F[\theta](z_{t-1}, x_t), \ z_0 := 0 \ (\forall t)$$

Output map

$$\hat{y}_t := G[\psi](z_t)$$

RNN parameterization

$$F[U, V](z, x) := \varphi(Uz + Vx), \ G[W](z) := \Phi(Wz), \ W \in \mathbb{R}^{q \times m}$$

BPTT

$$\begin{aligned} \frac{\partial h}{\partial z_t^s} &= \sum_{s=t}^T \sum_{k=1}^m \sum_j \frac{\partial \hat{y}_k^s}{\partial z_t^s} \frac{\partial z_j^s}{\partial z_t^s}, \ \frac{\partial \hat{y}_k^s}{\partial z_t^s} = \Phi_k^s w_{kj} \\ \frac{\partial h}{\partial v_{ij}^j} &= \sum_{t=1}^T \frac{\partial h}{\partial z_t^i} \varphi_i^t x_j^t, \ \frac{\partial h}{\partial u_{ij}^j} = \sum_{t=1}^T \frac{\partial h}{\partial z_t^i} \varphi_i^t z_j^{t-1} \end{aligned}$$

Spectral norm

$$\|A\|_2 = \max_{x: \|x\|=1} \|Ax\|_2 = \sigma_1(A)$$

Gradient norms

$$\frac{\partial z_T}{\partial z_0} = \Phi^T U \dots \Phi^1 U. \text{ Vanishes if } \sigma_1(U) < 1/\kappa, \text{ explodes if } \sigma_1(U) \text{ too large.}$$

Bidirectional RNNs

$$\hat{y}_t = \Phi(Wz_t + \tilde{W}\tilde{z}_t)$$

5.2 Gated Memory

LSTM

$$z_t := \sigma(F\tilde{x}_t) \odot z_{t-1} + \sigma(G\tilde{x}_t) \odot \tanh(V\tilde{x}_t), \ \tilde{x}_t := [x_t, \ell_t], \ \ell_{t+1} = \sigma(H\tilde{x}_t) \odot \tanh(Uz_t)$$

GRU

$$z_t = (1 - \sigma) \odot z_{t-1} + \sigma \odot \tilde{z}_t, \ \sigma := \sigma(G[x_t, z_{t-1}]), \ \tilde{z}_t := \tanh(V[\ell_t \odot z_{t-1}, x_t]), \ \ell_t := \sigma(H[z_{t-1}, x_t])$$

5.3 Linear Recurrent Models

Linear state evolution

$$z_{t+1} = Az_t + Bx_t$$

Diagonal form

$$A = P\Lambda P^{-1}, \ \Lambda := \text{diag}(\lambda_1, \dots, \lambda_m), \ \lambda_i \in \mathbb{C}$$

Stability

$$\max_j |\lambda_j| \leq 1$$

Initialization
$\lambda_i = \exp(-\exp(\nu_i) + i\theta_i)$, $e^{\nu_i} = -\ln r_i$ $\theta_i \sim \text{Uni}[0; 2\pi]$, $r_i \sim \text{Uni}[I]$, $I \subseteq [0; 1]$
Advantages
(i) clear long/short range dependencies (ii) no channel mixing required (iii) parallelizable training

6. Attention and Transformers

6.1 Attention

Attention mixing
$\xi_s := \sum_t a_{st} W x_t$, $a_{st} \geq 0$, $\sum_t a_{st} = 1$ $A = (a_{st}) \in \mathbb{R}^{T \times T}$, $\Xi = W X A^\top$

Query-key matching
$Q = U_Q X$, $K = U_K X$ ($U_Q, U_K \in \mathbb{R}^{q \times n}$) $Q^\top K = X^\top U_Q^\top U_K X$ ($Q^\top K \in \mathbb{R}^{T \times T}$, $\text{rank} \leq q$)

Softmax attention
$A = \text{softmax}(\beta Q^\top K)$, $a_{st} = \frac{e^{\beta [Q^\top K]_{st}}}{\sum_r e^{\beta [Q^\top K]_{sr}}}$, usually $\beta = 1/\sqrt{q}$

Feature transformation
$X \mapsto Y \mapsto F(Y)$, $F(\theta)(Y) = (F(y_1), \dots, F(y_T))$

Positional encoding
$p_{tk} = \begin{cases} \sin(t\omega_k) & k \text{ even} \\ \cos(t\omega_k) & k \text{ odd} \end{cases}$, $\omega_k = C^{k/K}$

Transformer architecture
Self-attention: attend to its own values in the past. Cross-attention: decoder attends to encoder output.

Vision transformer patch embedding
$\mathbb{R}^{p \times p \times c} \ni \text{patch}_t \mapsto x_t := V(\text{patch}_t) \in \mathbb{R}^n$ $V \in \mathbb{R}^{n \times (cp^2)}$

GELU activation
$\varphi(z) = z \Pr(z \leq Z)$, $Z \sim \mathcal{N}(0, 1)$

7. Geometric Deep Learning

7.1 Sets and Points

Order-invariance
$f(x_1, \dots, x_M) = f(x_{\pi_1}, \dots, x_{\pi_M}) \, \forall \pi \in S_M$ Permutation invariant sum: $\sum_{m=1}^M x_m = \sum_{m=1}^M x_{\pi m}$, $\forall M$, $\forall \pi \in S_M$

Equivariance
$f(x_{\pi_1}, \dots, x_{\pi_M}) = (y_{\pi_1}, \dots, y_{\pi_M})$

Deep Sets model
$f(x_1, \dots, x_M) = \rho \left(\sum_{m=1}^M \varphi(x_m) \right)$

Max pooling variant
$f(x_1, \dots, x_M) = \rho \left(\max_{m=1}^M \varphi(x_m) \right)$

Equivariant map
$\rho : \mathbb{R} \times \mathbb{R}^N \rightarrow Y$, $(x_m, \sum_{k=1}^M \varphi(x_k)) \mapsto y_m$

7.2 Graph Conv Networks

Feature & adjacency
$X = [x_1^\top; \dots; x_M^\top]$, $A = (a_{nm})$ with $a_{nm} = 1$ if $\{v_n, v_m\} \in E$
Graph invariance
$f(X, A) = f(PX, PAP^\top) \, \forall P$

Graph equivariance
$f(X, A) = P f(PX, PAP^\top) \, \forall P$
Message passing
$\varphi(x_m, X_m) = \varphi(x_m, \bigoplus_{X_m} \Phi(x))$, \bigoplus permutation-invariant
Normalized adjacency
$\bar{A} = D^{-1/2} (A + I) D^{-1/2}$, $D = \text{diag}(d_m)$, $d_m = 1 + \sum_n a_{nm}$
GCN layer
$X^+ = \sigma(\bar{A} X W)$, $W \in \mathbb{R}^{M \times N}$
Two-layer GCN: $Y = \text{softmax} \left(\bar{A} \sigma \left(\bar{A} X W^{(0)} \right) W^{(1)} \right)$

7.3 Spectral Graph Theory

Laplacian operator
$\Delta f := \sum_{n=1}^N \frac{\partial^2 f}{\partial x_n^2}$, $f : \mathbb{R}^N \rightarrow \mathbb{R}$

Graph Laplacian
$L = D - A$, $(Lx)_n = \sum_m a_{nm} (x_n - x_m)$
Normalized Laplacian
$\tilde{L} = I - D^{-1/2} A D^{-1/2}$

Graph Fourier transform
$L = U \Lambda U^\top$, $\Lambda := \text{diag}(\lambda_1, \dots, \lambda_M)$, $\lambda_i \geq \lambda_{i+1}$ Convolution: $x * y = U((U^\top x) \odot (U^\top y))$ Filtering: $G_\theta(L)x = U G_\theta(\Lambda) U^\top x$

Polynomial kernels
$U \left(\sum_{k=0}^K \alpha_k \Lambda^k \right) U^\top = \sum_{k=0}^K \alpha_k L^k$
Polynomial kernel network layer: $x_i^{l+1} = \sum_j p_{ij}(L) x_j^l + b_i$, $p_{ij}(L) = \sum_{k=0}^K \alpha_{ijk} L^k$

7.4 Attention GNNs

Attention coupling
$q_{ij} = \text{softmax}(f^l(u^\top (V x_i; V x_j; x_{ij})))$ s.t. $\sum_j A_{ij} q_{ij} = 1$

Attention propagation
$X^+ = \sigma(Q X W)$

Weisfeiler-Lehman test
Graph isomorphism test: iteratively aggregate node labels from neighborhoods.

8. Tricks of the Trade

8.1 Initialization

Random initialization
$\theta_i^0 \sim \mathcal{N}(0, \sigma_i^2)$ or $\theta_i^0 \sim \text{Uniform}[-\sqrt{3}\sigma_i; \sqrt{3}\sigma_i]$
LeCun initialization
$w_{ij} \sim \text{Uniform}[-a; a]$, $a := 1/\sqrt{n}$, $b_i = 0$. Stabilizes variance.
Glorot initialization
$w_{ij} \sim \text{Uniform}[-\sqrt{3}\epsilon; \sqrt{3}\epsilon]$, $\epsilon := 2/(n + m)$. Stabilizes gradient variance.

He initialization
$w_{ij} \sim \mathcal{N}(0, \epsilon)$ or $\text{Uniform}[-\sqrt{3}\epsilon; \sqrt{3}\epsilon]$, $\epsilon := 2/n$. For ReLU (half units active).

Orthogonal initialization
$\frac{1}{\sqrt{m}} W \sim \text{Uniform}(O(m))$ s.t. $W^\top W = W W^\top = m I$

8.2 Weight Decay

L_2 regularization
$\mu(\theta) = \frac{\mu}{2} \ \theta\ ^2$
GD with weight decay
$\dot{\theta} = -\eta \nabla E(\theta) - \eta \mu \theta$
Local loss landscape
$\theta_\mu^* = (H + \mu I)^{-1} H \theta^*$. Minimum shrunk along small eigenvalue directions.
Optimal weight decay
$\mu = \sigma^2/u^2$. Inverse proportional to signal-to-noise ratio.

8.3 Dropout

Dropout as Ensembling
$p(y x) = \sum_{b \in \{0,1\}^R} p(b) p(y x; b)$, $p(b) = \prod_{i=1}^R \pi_i^{b_i} (1 - \pi_i)^{1-b_i}$
Weight scaling for inference
$\tilde{w}_{ij} \leftarrow \pi_j w_{ij}$

8.4 Normalization

Batch normalization
$\tilde{f} = \frac{f - \mathbb{E}[f]}{\sqrt{\mathbb{V}[f]}}$, $\mathbb{E}[\tilde{f}] = 0$, $\mathbb{V}[\tilde{f}] = 1$
$\tilde{f}[\gamma, \beta] = \gamma + \beta \tilde{f}$ (learnable)
Weight normalization
$f(v, \epsilon)(x) = \varphi(w^\top x)$, $w := \frac{\epsilon}{\ v\ _2} v$
$\partial_\epsilon E = \nabla_w E \cdot \frac{v}{\ v\ _2}$, $\nabla_v E = \frac{\epsilon}{\ v\ _2} \left(I - \frac{w w^\top}{\ w\ _2^2} \right) \nabla_w E$

Layer normalization
$\tilde{f}_i = \frac{f_i - \mathbb{E}[f]}{\sqrt{\mathbb{V}[f]}}$, $\mathbb{E}[f] = \frac{1}{m} \sum_i f_i$, $\mathbb{V}[f] = \frac{1}{m} \sum_i (f_i - \mathbb{E}[f])^2$

Reparameterization trick
$z = \mu + \Sigma^{1/2} \eta$, $\eta \sim \mathcal{N}(0, I)$ $\nabla_\mu \mathbb{E}[f(z)] = \mathbb{E}[\nabla_z f(z)]$, $\nabla_\Sigma \mathbb{E}[f(z)] = \frac{1}{2} \mathbb{E}[\nabla_z^2 f(z)]$

8.5 Model Distillation

Tempered cross entropy
$\ell(x) = \sum_y \frac{q \exp[F_y/T]}{\sum_\nu \exp[F_\nu/T]} \left[\frac{1}{T} G_y - \ln \sum_\nu \exp[G_\nu/T] \right]$

Distillation gradient
$\frac{\partial \ell}{\partial G_y} = \frac{1}{T} \left[\frac{e^{q F_y/T}}{\sum_\nu e^{F_\nu/T}} - \frac{q e^{G_y/T}}{\sum_\nu e^{G_\nu/T}} \right]$

9. Theory

9.1 Infinite Width (NTK)

Neural tangent kernel
$k(x, \xi) = \nabla f(x) \cdot \nabla f(\xi)$, $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ Linearized: $h(\beta)(x) = f(x) + \beta \cdot \nabla f(x)$ with $\beta \approx \theta - \theta_0$
Gradient flow
ODE: $\dot{\theta} = \sum_{i=1}^s (y_i - f_i(\theta)) \nabla f_i(\theta)$ Functional: $\dot{f}_j = \sum_{i=1}^s (y_i - f_i) k^{(\theta)}(x_i, x_j)$, $\dot{f} = K^{(\theta)}(y - f)$

Dual representation
$h(\alpha)(x) = f(x) + \sum_{i=1}^s \alpha_i \nabla f(x_i) \cdot \nabla f(x)$ Optimal: $\alpha^* = K^+(y - f)$, $h^*(x) = k(x) K^+(y - f)$
Infinite width limit

$w_{ij}^l = \frac{\sigma_w}{\sqrt{m}} \omega_{ij}^l$, $b_i^l = \frac{\sigma_b}{\sqrt{m}} \beta_i^l$, $\omega^l, \beta^l \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$
$k^{(\theta)} \rightarrow k^\infty$ for $m_l \rightarrow \infty$

NTK constancy
$\frac{dk^{(\theta(t))}}{dt} = 0$, $f^\infty(x) = k(x) K^+(y - f)$ Near-constancy: $\ k(\theta_0) - k(\theta_t)\ _F^2 \in O(1/m)$
Kernel regression solution
$\hat{F}_\infty = K_0 (K_0 + \lambda I)^{-1} Y$
Function space
$\bar{F} \in \mathcal{H}_K$ (RKHS), $\ \bar{F}\ _{\mathcal{H}_K}^2 = \theta^\top \theta$

9.2 Bayesian DNNs

Parameter prior
$p(\theta) = \prod_{i=1}^d p(\theta_i)$, $\theta_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ $-\log p(\theta) = \frac{1}{2\sigma^2} \ \theta\ ^2 + \text{const}$ (weight decay)

Likelihood
$y_i = f^*(x_i) + \eta_i$, $\eta_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \delta^2)$ $-\log p(S \theta) = \frac{1}{2\delta^2} \ y - f(\theta)\ ^2 + \text{const}$

Posterior
$p(\theta S) = \frac{p(\theta) p(S \theta)}{p(S)}$, $-\log p(\theta S) = E(\theta) + \text{const}$ $E(\theta) = \frac{1}{2\delta^2} \ y - f\ ^2 + \frac{1}{2\sigma^2} \ \theta\ ^2$

Predictive distribution
$\bar{f}(x) = \int f(\theta)(x) p(\theta S) d\theta$
Bayesian ensembling: $\bar{f}^{(n)}(x)$ =
$\frac{\sum_{j=1}^n \exp[-E(\theta_j)] f(\theta_j)(x)}{\sum_{j=1}^n \exp[-E(\theta_j)]}$

9.3 GPs & Infinite Width

Gaussian processes
$(f(x_1), \dots, f(x_s)) \sim \mathcal{N}$; $\sum_{i=1}^s \alpha_i f(x_i) \sim \mathcal{N} \, \forall \alpha \in \mathbb{R}^s$ Mean $\mu(x) := \mathbb{E}_x[f(x)]$, covariance $k(x, \xi) := \mathbb{E}_{x, \xi}[f(x) f(\xi)] - \mu(x) \mu(\xi)$

GPs in DNNs
Linear layer: $w \sim \mathcal{N}(0, \frac{\sigma^2}{n} I)$, $\mathbb{E}[y_i y_j] = \sigma^2 x_i^\top x_j$ Deep layers: near-normal for high-dim inputs. $f \sim \mathcal{GP}(0, K^{l-1})$

Kernel recursion
$K_{\mu\nu}^l = \mathbb{E}[\varphi(x_{\mu}^{l-1}) \varphi(x_{\nu}^{l-1})] = \sigma^2 \mathbb{E}[\varphi(f_\mu) \varphi(f_\nu)]$ Example kernels: $k(x, \xi) = x^\top \xi$, $k(x, \xi) = e^{-\gamma \ x - \xi\ ^2}$

Kernel regression
$f^*(x) = k(x)^\top K^+ y$ (mean of Bayesian predictive) $\mathbb{E}[(f(x) - f^*(x))^2] = K(x, x) - k(x)^\top K^+ k(x)$

9.4 Statistical Learning Theory

Generalization gap
$\mathcal{R}(f) - \hat{\mathcal{R}}(f) = \mathbb{E}_{x, y}[\ell(f(x), y)] - \frac{1}{n} \sum_i \ell(f(x_i), y_i)$
PAC bound
$\mathbb{P}[\mathcal{R}(f) - \hat{\mathcal{R}}(f) \leq \epsilon] \geq 1 - \delta$
VC dimension
VC-dim(F) := max_s sup _{S =s} 1 [$F(S)$] = 2^s
VC inequality: $\mathbb{P}[\sup_F \hat{E}(f) - E(f) > \epsilon] \leq 8 F(s) e^{-s \epsilon^2/32}$
Rademacher complexity
$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{\sigma, S}[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_i \sigma_i f(x_i)]$
Generalization bound
$\mathcal{R}(f) \leq \hat{\mathcal{R}}(f) + 2 \mathcal{R}_n(\mathcal{F}) + \sqrt{\frac{\ln(1/\delta)}{2n}}$
Bias-variance tradeoff
$\mathbb{E}[(f(x) - y)^2] = \text{Bias}^2 + \text{Var} + \sigma^2$

Generalization gap
 Δ<!-- Δ --> := max (0 , E −<!-- − --> Ê<!-- Ê -->) , E : expected population error, Ê<!-- Ê --> : empirical
Double descent
Beyond interpolation point, models may level out at lower generalization error.
KL divergence
 D KL (p ∥<!-- ∥ --> q) = ∫<!-- ∫ --> p (x) log ⁡<!-- ⁡ --> p (x) q (x) d x = E x ∼<!-- ∼ --> p [ln ⁡<!-- ⁡ --> p (x) q (x)]

9.5 Loss Landscape

Critical points
 ∇<!-- ∇ --> ℳ<!-- ℳ --> (θ<!-- θ --> ∗<!-- ∗ -->) = 0 . Local min : H ⪰<!-- ⪰ --> 0 . Saddle : H indefinite.
Sharpness
 λ<!-- λ --> max (H) . Flat minima →<!-- → --> better generalization.
Mode connectivity
Local minima connected by paths of low loss.
Lottery ticket hypothesis
Sparse subnetworks can match full network performance if initialized correctly.

10. Generative Models

10.1 Variational Auto Encoders

Linear autoencoder
 x ↦<!-- ↦ --> z = C x , z ↦<!-- ↦ --> x ̂<!-- ̂ --> = D z , E (C , D) (x) = 1 2 ∥<!-- ∥ --> x −<!-- − --> D C x ∥<!-- ∥ --> 2 D C X = X ̂<!-- ̂ --> = U Σ<!-- Σ --> m V ⊤<!-- ⊤ --> , for centered data ≡<!-- ≡ --> PCA
Linear factor analysis
 x = μ<!-- μ --> + W z + η<!-- η --> , η<!-- η --> ∼<!-- ∼ --> N (0 , Ψ<!-- Ψ -->) , x ∼<!-- ∼ --> N (μ<!-- μ --> , W W ⊤<!-- ⊤ --> + Ψ<!-- Ψ -->) for z ∼<!-- ∼ --> N (0 , I)
 μ<!-- μ --> z x = W ⊤<!-- ⊤ --> (W W ⊤<!-- ⊤ --> + Ψ<!-- Ψ -->) −<!-- − --> 1 (x −<!-- − --> μ<!-- μ -->)
Generative model
 p θ<!-- θ --> (x , z) = p θ<!-- θ --> (x z) p (z) , p (z) = N (0 , I)

ELBO
 log ⁡<!-- ⁡ --> p (θ<!-- θ -->) (x) = log ⁡<!-- ⁡ --> ∫<!-- ∫ --> q (z) p (θ<!-- θ -->) (x z) p (z) q (z) d z ≥<!-- ≥ --> ∫<!-- ∫ --> q (z) log ⁡<!-- ⁡ --> p (θ<!-- θ -->) (x z) d z −<!-- − --> D KL (q ∥<!-- ∥ --> p) =: L (θ<!-- θ --> , q) (x)
Inference network
 z ∼<!-- ∼ --> N (μ<!-- μ --> (x) , Σ<!-- Σ --> (x))
Encoder
 q ϕ<!-- ϕ --> (z x) = N (μ<!-- μ --> ϕ<!-- ϕ --> (x) , σ<!-- σ --> ϕ<!-- ϕ --> 2 (x))
Decoder
 p θ<!-- θ --> (x z) = N (μ<!-- μ --> θ<!-- θ --> (z) , σ<!-- σ --> 2 I) or Bernoulli
Reparameterization trick
 z = μ<!-- μ --> ϕ<!-- ϕ --> (x) + σ<!-- σ --> ϕ<!-- ϕ --> (x) ⊙<!-- ⊙ --> ϵ<!-- ϵ --> , ϵ<!-- ϵ --> ∼<!-- ∼ --> N (0 , I)
KL divergence (Gaussian)
 D KL = 1 2 ∑<!-- ∑ --> j (μ<!-- μ --> j 2 + σ<!-- σ --> j 2 −<!-- − --> ln ⁡<!-- ⁡ --> σ<!-- σ --> j 2 −<!-- − --> 1)
 β<!-- β --> -VAE
 ℳ<!-- ℳ --> = E q [ln ⁡<!-- ⁡ --> p (x z)] −<!-- − --> β<!-- β --> D KL (q ∥<!-- ∥ --> p) . β<!-- β --> > 1 for disentangle-ment.

10.2 Generative Adversarial Networks

Generator
 G : z ↦<!-- ↦ --> x , z ∼<!-- ∼ --> p (z)
Discriminator
 D : x ↦<!-- ↦ --> [0 , 1] , probability that x is real
GAN objective
 V (G , D) = E x r ∼<!-- ∼ --> p data [D (x r)] + E z ∼<!-- ∼ --> p z [1 −<!-- − --> D (G (z))]
Bayes-optimal classifier
 q θ<!-- θ --> (x) := P { y = 1 x } = p (x) p (x) + p θ<!-- θ --> (x)
Jensen-Shannon objective
 ℓ<!-- ℓ --> ∗<!-- ∗ --> = JS (p , p θ<!-- θ -->) −<!-- − --> ln ⁡<!-- ⁡ --> 2 ℓ<!-- ℓ --> ∗<!-- ∗ --> (θ<!-- θ -->) ≥<!-- ≥ --> sup ϕ<!-- ϕ --> ℓ<!-- ℓ --> (θ<!-- θ --> , ϕ<!-- ϕ -->) where ϕ<!-- ϕ --> : discriminator, θ<!-- θ --> : genera-tor

Alternating gradient descent
 θ<!-- θ --> t + 1 = θ<!-- θ --> t −<!-- − --> η<!-- η --> ∇<!-- ∇ --> θ<!-- θ --> ℓ<!-- ℓ --> (θ<!-- θ --> t , ϕ<!-- ϕ --> t) ϕ<!-- ϕ --> t + 1 = ϕ<!-- ϕ --> t + η<!-- η --> ∇<!-- ∇ --> ϕ<!-- ϕ --> ℓ<!-- ℓ --> (θ<!-- θ --> t + 1 , ϕ<!-- ϕ --> t)
Wasserstein GAN
 min G max D ∈<!-- ∈ --> 1 - Lip E x [D (x)] −<!-- − --> E z [D (G (z))]
Gradient penalty (WGAN-GP)
 λ<!-- λ --> E x̂<!-- x̂ --> [(∥<!-- ∥ --> ∇<!-- ∇ --> x̂<!-- x̂ --> D (x̂<!-- x̂ -->) ∥<!-- ∥ --> 2 −<!-- − --> 1) 2] , x̂<!-- x̂ --> = α<!-- α --> x + (1 −<!-- − --> α<!-- α -->) G (z)

10.3 Denoising Diffusion

Forward process
 q (x t x t −<!-- − --> 1) = N (x t ; √<!-- √ --> 1 −<!-- − --> β<!-- β --> t x t −<!-- − --> 1 , β<!-- β --> t I)
Marginal
 q (x t x 0) = N (√<!-- √ --> α<!-- α --> t x 0 , β<!-- β --> t I) , α<!-- α --> t = ∏<!-- ∏ --> τ<!-- τ --> = 1 t (1 −<!-- − --> β<!-- β --> τ<!-- τ -->) , β<!-- β --> t = 1 −<!-- − --> α<!-- α --> t
 ν<!-- ν --> t ≈<!-- ≈ --> N (√<!-- √ --> α<!-- α --> t x 0 , β<!-- β --> t I) t →<!-- → --> ∞<!-- ∞ --> →<!-- → --> N (0 , I)
Reverse process
 p θ<!-- θ --> (x t −<!-- − --> 1 x t) = N (m θ<!-- θ --> (x t , t) , Σ<!-- Σ --> (x t , t))
Forward: π<!-- π --> ∗<!-- ∗ --> →<!-- → --> ν<!-- ν --> T = π<!-- π --> . Backward: π<!-- π --> →<!-- → --> μ<!-- μ --> 0 θ<!-- θ --> ≈<!-- ≈ --> π<!-- π --> ∗<!-- ∗ -->
Training objective
 L t = ∥<!-- ∥ --> m (x t , x 0 , t) −<!-- − --> m θ<!-- θ --> (x t , t) ∥<!-- ∥ --> 2 2 σ<!-- σ --> t 2 + const
Reparameterization: x t = √<!-- √ --> α<!-- α --> t x 0 + √<!-- √ --> 1 −<!-- − --> α<!-- α --> t ϵ<!-- ϵ -->
Simplified criterion
 h (θ<!-- θ -->) (x) = 1 T ∑<!-- ∑ --> t = 1 T E [∥<!-- ∥ --> ϵ<!-- ϵ --> −<!-- − --> ϵ<!-- ϵ --> θ<!-- θ --> (√<!-- √ --> α<!-- α --> t x + √<!-- √ --> 1 −<!-- − --> α<!-- α --> t ϵ<!-- ϵ --> , t) ∥<!-- ∥ --> 2]
Forward trajectory target
 x t −<!-- − --> 1 x t , x 0 = N (m (x t , x 0 , t) , β<!-- β --> t I)
 m (x t , x 0 , t) = √<!-- √ --> α<!-- α --> t −<!-- − --> 1 β<!-- β --> t x 0 + (1 −<!-- − --> α<!-- α --> t −<!-- − --> 1) √<!-- √ --> 1 −<!-- − --> β<!-- β --> t 1 −<!-- − --> α<!-- α --> t x t
Sampling
 m (x t , x 0 , t) = 1 √<!-- √ --> α<!-- α --> t ⎡<!-- ⎡ --> x t −<!-- − --> β<!-- β --> t √<!-- √ --> 1 −<!-- − --> α<!-- α --> t ϵ<!-- ϵ --> ⎤<!-- ⎤ -->

Score function
 ∇<!-- ∇ --> x ln ⁡<!-- ⁡ --> p (x) ≈<!-- ≈ --> −<!-- − --> ϵ<!-- ϵ --> θ<!-- θ --> (x , t) √<!-- √ --> 1 −<!-- − --> α<!-- α --> t }
Classifier-free guidance
 ε<!-- ε --> ̃<!-- ̃ --> = (1 + w) ϵ<!-- ϵ --> θ<!-- θ --> (x t , t , c) −<!-- − --> w ϵ<!-- ϵ --> θ<!-- θ --> (x t , t)
DDIM (deterministic)
 x t −<!-- − --> 1 = √<!-- √ --> α<!-- α --> t −<!-- − --> 1 (x t −<!-- − --> √<!-- √ --> 1 −<!-- − --> α<!-- α --> t ϵ<!-- ϵ --> θ<!-- θ -->) + √<!-- √ --> 1 −<!-- − --> α<!-- α --> t −<!-- − --> 1 ϵ<!-- ϵ --> θ<!-- θ -->

11. Ethics

11.1 Adversarial Examples

Adversarial perturbation
 x ′<!-- ′ --> = x + δ<!-- δ --> , ∥<!-- ∥ --> δ<!-- δ --> ∥<!-- ∥ --> p ≤<!-- ≤ --> ϵ<!-- ϵ --> , F (x ′<!-- ′ -->) ≠<!-- ≠ --> F (x) ∥<!-- ∥ --> x ∥<!-- ∥ --> p = (∑<!-- ∑ --> i x i p) 1 / p , ∥<!-- ∥ --> x ∥<!-- ∥ --> ∞<!-- ∞ --> = max i x i , ∥<!-- ∥ --> x ∥<!-- ∥ --> 0 = { i : x i ≠<!-- ≠ --> 0 }
FGSM (Fast Gradient Sign Method)
 δ<!-- δ --> = ϵ<!-- ϵ --> ⋅<!-- ⋅ --> sign (∇<!-- ∇ --> x ℳ<!-- ℳ --> (f (x) , y))
PGD (Projected Gradient Descent)
 η<!-- η --> t + 1 = Π<!-- Π --> ϵ<!-- ϵ --> [η<!-- η --> t + α<!-- α --> ∇<!-- ∇ --> x ℳ<!-- ℳ --> (f (x + η<!-- η --> t) , y)] Π<!-- Π --> ϵ<!-- ϵ --> [z] := z / ∥<!-- ∥ --> z ∥<!-- ∥ --> 2 for p = 2 , Π<!-- Π --> ϵ<!-- ϵ --> [z] := z / ∥<!-- ∥ --> z ∥<!-- ∥ --> ∞<!-- ∞ --> for p = ∞<!-- ∞ -->
Adversarial training
 min θ<!-- θ --> E x , y [max δ<!-- δ --> ∥<!-- ∥ --> δ<!-- δ --> ∥<!-- ∥ --> ≤<!-- ≤ --> ϵ<!-- ϵ --> ℳ<!-- ℳ --> (f (x + δ<!-- δ -->) , y ; θ<!-- θ -->)]
DeepFool (binary)
Optimal perturbation: η<!-- η --> ∝<!-- ∝ --> sign (f 1 (x) −<!-- − --> f 2 (x)) (w 2 −<!-- − --> w 1) for f i = w i ⊤<!-- ⊤ --> x + b i . Iterate: arg min ∥<!-- ∥ --> η<!-- η --> ∥<!-- ∥ --> 2 s.t. (∇<!-- ∇ --> f 1 −<!-- − --> ∇<!-- ∇ --> f 2) ⊤<!-- ⊤ --> η<!-- η --> < f 1 (x) −<!-- − --> f 2 (x)
Certified robustness
Provable guarantees via randomized smoothing, interval bound propagation.
Transferability
Adversarial examples often transfer between models.