

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 $b \leftarrow b + y$
 $w^t \in \text{span}(x_1, \dots, x_s) \Rightarrow w^t \in \text{span}(x_1, \dots, x_s) (\forall t)$

Convergence

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

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

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

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

$$E(X) = -\frac{1}{2} \sum_{i \neq j} w_{ij} X_i X_j + \sum_i b_i X_i \text{ where } X_i \in \{-1, +1\}$$

$w_{ij} = w_{ji}$ ($\forall i, j$), $w_{ii} = 0$ ($\forall i$): Interaction strengths

Hebbian Learning

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

2. Feedforward Networks

2.1 Linear Models

Linear regression

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

Moore-Penrose inverse solution

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

SGD update

$$w_{t+1} := w_t + \eta \underbrace{(y_{it} - w_t^\top x_{it})}_{\text{residual}} x_{it}$$

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 of noise model.

Ridge regression

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

Logistic function

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

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

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

$$F : \underbrace{\mathbb{R}^{m(n+1)}}_{\text{parameters}} \times \underbrace{\mathbb{R}^n}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^m}_{\text{output}}$$

$$F[\theta](x) := \phi(Wx + b), \quad \theta := \text{vec}(W, b)$$

Composition of layers

$$G = F^L[\theta^L] \circ \dots \circ F^1[\theta^1]$$

where $F^l[W^l, b^l](x) := \phi^l(W^l x + b^l)$

Layer activations

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

Softmax function

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

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

Residual layer

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

Skip connection

Concatenate previous layer back in

2.3 Sigmoid Networks

Sigmoid/Tanh activations

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

Barron's Theorem

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

2.4 ReLU Networks

ReLU activation

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

Zaslavsky: Connected regions

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

Montufar: Connected regions

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

3. Gradient-Based Learning

3.1 Backpropagation

Parameter derivatives

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

Loss derivatives

$$\frac{\partial h}{\partial w_{ij}^l} = \delta_i^l \dot{\phi}_i^l x_j^{l-1}, \quad \frac{\partial h}{\partial b_i^l} = \delta_i^l \dot{\phi}_i^l \text{ with } \delta_i^l = \frac{\partial h}{\partial x_i^l} \dot{\phi}_i^l$$

3.2 Gradient Descent

GD update & flow

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

L-smoothness

$$\|\nabla h(\theta_1) - \nabla h(\theta_2)\| \leq L \|\theta_1 - \theta_2\| \quad (\forall \theta_1, \theta_2)$$

$$\lambda_{\max}(\nabla^2 h) \leq L, \quad \ell''(x) \leq L$$

$$\ell(w) - \ell(w') \leq \nabla \ell(w')^\top (w - w') + \frac{L}{2} \|w - w'\|_2^2$$

Polyak-Łojasiewicz

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

Convergence rate

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

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

$$\theta^{t+1} = \theta^t + \beta(\theta^t - \theta^{t-1}), \quad \theta^{t+1} = \tilde{\theta}^{t+1} - \eta \nabla h(\tilde{\theta}^{t+1})$$

More theoretical grounding than heavy ball.

AdaGrad

$$\nu_i^t = \nu_i^{t-1} + [\partial_i h(\theta^t)]^2, \quad \eta_i^t = \frac{\eta}{\sqrt{\nu_i^t + \epsilon}}$$

$$\theta_i^{t+1} = \theta_i^t - \eta_i^t \partial_i h(\theta^t)$$

Adam

$$g_i^t = \beta g_i^{t-1} + (1 - \beta) \partial_i h(\theta^t)$$

$$\nu_i^t = \alpha \nu_i^{t-1} + (1 - \alpha) [\partial_i h(\theta^t)]^2$$

$$\theta_i^{t+1} = \theta_i^t - \frac{\eta}{\sqrt{\nu_i^t + \epsilon}} g_i^t$$

RMSprop

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

Muon

$$M^t = \mu_1 M^{t-1} + \nabla h(\theta^t)$$

$$P^t = \text{orthogonalize}(M^t)$$

$$\theta^{t+1} = \theta^t - \eta \sqrt{d_{\text{out}}/d_{\text{in}}} P^t \text{ orthogonalize}(A) = \operatorname{argmin}_{B: B^\top B = I} \|A - B\|_F = UV^\top \text{ where } A = U\Sigma V^\top \text{ (truncated SVD). Muon solves constrained opt. problem:}$$

$$\min_{\theta} \langle \nabla \mathcal{L}(\theta), \Delta \theta \rangle_F + \frac{1}{2} \frac{d_{\text{in}}}{d_{\text{out}}} \|\Delta \theta\|_F^2 \text{ s.t. } \|\Delta \theta\|_2 \leq 1$$

whose solution is an update in direction of the top singular vectors: $\Delta W \propto -UV^\top$

3.4 SGD

SGD update & variance

$$\theta_{t+1} = \theta_t - \eta \nabla h(\theta_t)(x_{it}, y_{it})$$

$$V[\theta] = \frac{1}{s} \sum_{i=1}^s \|\nabla h(\theta_t) - \nabla h(x_i, y_i)\|^2$$

SGD convergence

Convergence in expectation: $\mathbb{E}[h(\theta^t)] - \min h$

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

3.5 Function properties

Convexity

$$\ell(\lambda w + (1 - \lambda)w') \leq \lambda \ell(w) + (1 - \lambda) \ell(w'), \quad \ell''(x) \geq 0 \quad \forall x$$

$$\ell(w) \geq \ell(w') + \nabla \ell(w')^\top (w - w') \text{ (differentiable case)}$$

Strong convexity

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

$$\ell''(x) \geq \mu$$

4.1 Convolutions

Convolution definition

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

Cross-correlation

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

Fourier property

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

Discrete convolution

$$1\text{D: } (f * g)[u] := \sum_{t=-\infty}^{\infty} f[t] g[u - t]$$

$$2\text{D: } (f * g)[u, v] := \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} f[m, v - n] g[m, n]$$

Discrete cross-correlation

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

Toeplitz matrices

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

$$G = \text{Toeplitz-Matrix}(g):$$

$$\begin{bmatrix} g[0] & g[-1] & g[-2] & \cdots & g[-(n-1)] \\ g[1] & g[0] & g[-1] & \ddots & g[-(n-2)] \\ g[2] & g[1] & g[0] & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ g[n-1] & g[n-2] & \cdots & g[1] & g[0] \end{bmatrix}$$

Then $G_{ij} = g[i - j]$

Properties of Convolutions

$$(f * g)(u) = (g * f)(u) \implies * \text{ is commutative}$$

$$(f * g)(u) = (g * f)(-u) \implies * \text{ is not}$$

$$(f * g)(y) = (\mathcal{T}(f) * g)(y) \text{ with } (\mathcal{T}(f))(u) = f(-u)$$

$$S_t(f * g)(u) = (S_t f) * g(u) = (f * S_t g)(u) \text{ with } S_t(f)(u) = f(u - t) \implies * \text{ is equivariant under } S_t$$

A is a continuous linear operator. Then A equivariant under translations $\implies A$ is a convolution: $\exists g$ s.t. $A(f) = f * g$

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

$$y[r][s, t] = \sum_u \sum_{\Delta s, \Delta t} w[r, u] [\Delta s, \Delta t] \cdot x[u][s + \Delta s, t + \Delta t]$$

Parameters count

$$D = \#r \cdot \#\underbrace{\Delta s}_{\text{channels}} \cdot \#\underbrace{\Delta t}_{\text{window size}} \text{ (in-channels fully connected to out-channels)}$$

Separable Kernels

If $w = uv^\top \in \mathbb{R}^{q \times q}$ and $x \in \mathbb{R}^{d \times d}$, then $x * w$ can be computed in $\mathcal{O}(d(d-q)q)$ instead of $\mathcal{O}((d-q)^2 q^2)$.

4.3 NLP with ConvNets

Word embedding

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

Conditional log-bilinear model

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

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

Negative sampling		Weight scaling for inference	
$\hat{\ell}_{\omega,\nu} = -\ln \sigma(x_{\omega}^T y_{\nu}) - \beta \mathbb{E}_{\mu \sim D} \ln(1 - \sigma(x_{\omega}^T y_{\mu}))$		$\tilde{w}_{ij} \leftarrow \pi_j w_{ij}$	
<h2>5. Recurrent Networks</h2>		<h2>8.4 Normalization</h2>	
<h3>5.1 Simple RNNs</h3>		Batch normalization	
Time evolution $z_t := F[\theta](z_{t-1}, x_t), z_0 := 0 \ (\forall t)$		$\bar{f} = \frac{f - \mathbb{E}[f]}{\sqrt{\mathbb{V}[f]}}, \mathbb{E}[\bar{f}] = 0, \mathbb{V}[\bar{f}] = 1$	
Output map $\hat{y}_t := G[\psi](z_t)$		$\bar{f}[\gamma, \beta] = \gamma + \beta \bar{f}$ (learnable)	
RNN parameterization $F[U, V](z, x) := \varphi(Uz + Vx), G[W](z) := \Phi(Wz), W \in \mathbb{R}^{q \times m}$		Weight normalization	
BPTT		$f(v, \epsilon)(x) = \varphi(w^T x), w := \frac{\epsilon}{\ v\ _2} v$	
$\frac{\partial h}{\partial z_i^t} = \sum_{s=t}^T \sum_{k=1}^m \sum_j \frac{\partial \hat{y}_k^s}{\partial z_j^t} \frac{\partial z_j^s}{\partial z_i^t}, \frac{\partial \hat{y}_k^s}{\partial z_j^t} = \dot{\Phi}_k^s w_{kj}$		$\partial_{\epsilon} E = \nabla_w E \cdot \frac{v}{\ v\ _2}, \nabla_v E = \frac{\epsilon}{\ v\ _2} \left(I - \frac{w w^T}{\ w\ _2^2} \right) \nabla_w E$	
$\frac{\partial h}{\partial v_{ij}} = \sum_{t=1}^T \frac{\partial h}{\partial z_i^t} \dot{\varphi}_i^t x_j^t, \frac{\partial h}{\partial u_{ij}} = \sum_{t=1}^T \frac{\partial h}{\partial z_i^t} \dot{\varphi}_i^t z_j^t$		Layer normalization	
Spectral norm $\ A\ _2 = \max_{x: \ x\ =1} \ Ax\ _2 = \sigma_1(A)$		$\tilde{f}_i = \frac{f_i - \mathbb{E}[f]}{\sqrt{\mathbb{V}[f]}}, \mathbb{E}[f] = \frac{1}{m} \sum_i f_i, \mathbb{V}[f] = \frac{1}{m} \sum_i (f_i - \mathbb{E}[f])^2$	
Gradient norms $\frac{\partial z_T}{\partial z_0} = \dot{\Phi}^T U \dots \dot{\Phi}^1 U.$ Vanishes if $\sigma_1(U) < 1/\kappa$, explodes if $\sigma_1(U)$ too large.		Reparameterization trick	
Bidirectional RNNs $\hat{y}_t = \Phi(Wz_t + \tilde{W}\tilde{z}_t)$		$z = \mu + \Sigma^{1/2} \eta, \eta \sim \mathcal{N}(0, I)$ $\nabla_{\mu} \mathbb{E}[f(z)] = \mathbb{E}[\nabla_z f(z)], \nabla_{\Sigma} \mathbb{E}[f(z)] = \frac{1}{2} \mathbb{E}[\nabla^2 f(z)]$	
<h2>5.2 Gated Memory</h2>		<h2>8.5 Model Distillation</h2>	
LSTM $z_t := \sigma(F\tilde{x}_t) \odot z_{t-1} + \sigma(G\tilde{x}_t) \odot \tanh(V\tilde{x}_t), \tilde{x}_t := [x_t, \ell_t], \ell_{t+1} = \sigma(H\tilde{x}_t) \odot \tanh(Uz_t)$		Tempered cross entropy	
GRU $z_t = (1 - \sigma) \odot z_{t-1} + \sigma \odot \tilde{z}_t, \sigma := \sigma(G[x_t, z_{t-1}])$ $\tilde{z}_t := \tanh(V[\ell_t \odot z_{t-1}, x_t]), \ell_t := \sigma(H[z_{t-1}, x_t])$		$\ell(x) = \sum_y \frac{q \exp[F_y/T]}{\sum_{\nu} \exp[F_{\nu}/T]} [\frac{1}{T} G_y - \ln \sum_{\nu} \exp[G_{\nu}/T]]$	
<h2>5.3 Linear Recurrent Models</h2>		Distillation gradient	
Linear state evolution $z_{t+1} = Az_t + Bx_t$		$\frac{\partial \ell}{\partial G_y} = \frac{1}{T} \left[\frac{e^{qF_y/T}}{\sum_{\nu} e^{F_{\nu}/T}} - \frac{qe^{G_y/T}}{\sum_{\nu} e^{G_{\nu}/T}} \right]$	
Diagonal form $A = P \Lambda P^{-1}, \Lambda := \text{diag}(\lambda_1, \dots, \lambda_m), \lambda_i \in \mathbb{C}$		<h2>9. Theory</h2>	
Stability $\max_j \lambda_j \leq 1$		<h3>9.1 Infinite Width (NTK)</h3>	
Initialization $\lambda_i = \exp(-\exp(\nu_i) + i\theta_i), e^{\nu_i} = -\ln r_i$ $\theta_i \sim \text{Uni}[0; 2\pi], r_i \sim \text{Uni}[I], I \subseteq [0; 1]$		Neural tangent kernel	
Advantages (i) clear long/short range dependencies (ii) no channel mixing required (iii) parallelizable training		$k(x, \xi) = \nabla f(x) \cdot \nabla f(\xi), \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ Linearized: $h(\beta)(x) = f(x) + \beta \cdot \nabla f(x)$ with $\beta \approx \theta - \theta_0$	
<h2>6. Attention and Transformers</h2>		Gradient flow	
<h3>6.1 Attention</h3>		$\text{ODE: } \dot{\theta} = \sum_{i=1}^s (y_i - f_i(\theta)) \nabla f_i(\theta)$ Functional: $\dot{f}_j = \sum_{i=1}^s (y_i - f_i) k^{(\theta)}(x_i, x_j), \dot{f} = K^{(\theta)}(y - f)$	
Attention mixing $\xi_s := \sum_t a_{st} W x_t, a_{st} \geq 0, \sum_t a_{st} = 1$ $A = (a_{st}) \in \mathbb{R}^{T \times T}, \Xi = W X A^T$		Dual representation	
Query-key matching $Q = U_Q X, K = U_K X$ ($U_Q, U_K \in \mathbb{R}^{q \times n}$) $Q^T K = X^T U_Q^T U_K X$ ($Q^T K \in \mathbb{R}^{T \times T}$, rank $\leq q$)		$h(\alpha)(x) = f(x) + \sum_{i=1}^s \alpha_i \nabla f(x_i) \cdot \nabla f(x)$ Optimal: $\alpha^* = K^+(y - f), h^*(x) = k(x) K^+(y - f)$	
Softmax attention $A = \text{softmax}(\beta Q^T K), a_{st} = \frac{e^{\beta(Q^T K)_{st}}}{\sum_r e^{\beta(Q^T K)_{sr}}}, \text{ usually } \beta = 1/\sqrt{q}$		Infinite width limit	
Feature transformation $X \mapsto Y \mapsto F(Y), F(\theta)(Y) = (F(y_1), \dots, F(y_T))$		$w_{ij}^l = \frac{\sigma_w}{\sqrt{m}} \omega_{ij}^l, b_i^l = \frac{\sigma_b}{\sqrt{m}} \beta_i^l, \omega^l, \beta^l \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ $k^{(\theta)} \rightarrow k^{\infty}$ for $m_l \rightarrow \infty$	
Positional encoding $p_{tk} = \begin{cases} \sin(t\omega_k) & k \text{ even} \\ \cos(t\omega_k) & k \text{ odd} \end{cases}, \omega_k = C^{k/K}$		NTK constancy	
Transformer architecture Self-attention: attend to its own values in the past. Cross-attention: decoder attends to encoder output.		$\frac{dk^{(\theta(t))}}{dt} = 0, f^{\infty}(x) = k(x) K^+(y - f)$ Near-constancy: $\ k(\theta_0) - k(\theta_t)\ _F^2 \in O(1/m)$	
Vision transformer patch embedding $\mathbb{R}^{p \times p \times c} \ni \text{patch}_t \mapsto x_t := V(\text{patch}_t) \in \mathbb{R}^n$ $V \in \mathbb{R}^{n \times (cp)^2}$		Kernel regression solution	
GELU activation $\varphi(z) = z \Pr(z \leq Z), Z \sim \mathcal{N}(0, 1)$		$\bar{F}_{\infty} = K_0(K_0 + \lambda I)^{-1} Y$	
<h3>7. Geometric Deep Learning</h3>		Function space	
<h4>7.1 Sets and Points</h4>		$\bar{F} \in \mathcal{H}_K$ (RKHS), $\ \bar{F}\ _{\mathcal{H}_K}^2 = \theta^T \theta$	
Order-invariance $f(x_1, \dots, x_M) = f(x_{\pi_1}, \dots, x_{\pi_M}) \ \forall \pi \in S_M$		<h4>7.2 Graph Conv Networks</h4>	
Permutation invariant sum: $\sum_{m=1}^M x_m = \sum_{m=1}^M x_{\pi_m}, \forall M, \forall \pi \in S_M$		Feature & adjacency	
Equivariance $f(x_{\pi_1}, \dots, x_{\pi_M}) = (y_{\pi_1}, \dots, y_{\pi_M})$		$X = [x_1^T; \dots; x_M^T], A = (a_{nm})$ with $a_{nm} = 1$ if $\{v_n, v_m\} \in E$	
Deep Sets model $f(x_1, \dots, x_M) = \rho \left(\sum_{m=1}^M \varphi(x_m) \right)$		Graph invariance	
Max pooling variant $f(x_1, \dots, x_M) = \rho \left(\max_{m=1}^M \varphi(x_m) \right)$		$f(X, A) = f(PX, PAP^T) \ \forall P$	
Equivariant map $\rho: \mathbb{R} \times \mathbb{R}^N \rightarrow Y, (x_m, \sum_{k=1}^M \varphi(x_k)) \mapsto y_m$		Graph equivariance	
<h4>7.3 Spectral Graph Theory</h4>		$f(X, A) = Pf(X, PAP^T) \ \forall P$	
Laplacian operator $\Delta f := \sum_{n=1}^N \frac{\partial^2 f}{\partial x_n^2}, f: \mathbb{R}^N \rightarrow \mathbb{R}$		Message passing $\varphi(x_m, X_m) = \varphi(x_m, \bigoplus_{X_m} \Phi(x)), \bigoplus \text{ permutation-invariant}$	
Normalized adjacency $\bar{A} = D^{-1/2} (A + I) D^{-1/2}, D = \text{diag}(d_m), d_m = 1 + \sum_n a_{nm}$		Normalizing adjacency $\theta_{\mu}^* = (H + \mu I)^{-1} H \theta^*. \text{ Minimum shrunk along small eigenvalue directions.}$	
GCN layer $X^+ = \sigma(\bar{A} X W), W \in \mathbb{R}^{M \times N}$		Optimal weight decay $\mu = \sigma^2/u^2.$ Inverse proportional to signal-to-noise ratio.	
Two-layer GCN: $Y = \text{softmax} \left(\bar{A} \sigma \left(\bar{A} X W^{(0)} \right) W^{(1)} \right)$		<h4>8.2 Weight Decay</h4>	
<h4>8.3 Dropout</h4>		Dropout as Ensembling	
Dropout as Ensembling $p(y x) = \sum_{b \in \{0,1\}^R} p(b)p(y x;b), p(b) = \prod_{i=1}^R \pi_i^{b_i} (1 - \pi_i)^{1-b_i}$		$p(\theta) = \prod_{i=1}^q p(\theta_i), \theta_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ $-\log p(\theta) = \frac{1}{2\sigma^2} \ \theta\ ^2 + \text{const (weight decay)}$	

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

$$\bar{f}(x) = \int f(\theta)(x)p(\theta|S)d\theta$$

$$\text{Bayesian ensembling: } \bar{f}^{(n)}(x) = \sum_{i=1}^n \frac{\exp[-E(\theta_i)]f(\theta_i)(x)}{\sum_{j=1}^n \exp[-E(\theta_j)]}$$

9.3 GPs & Infinite Width

Gaussian processes

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

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

$$K_{\mu\nu}^l = \mathbb{E}[\varphi(x_{i\mu}^{l-1})\varphi(x_{i\nu}^{l-1})] = \sigma^2 \mathbb{E}[\varphi(f_\mu)\varphi(f_\nu)]$$

Example kernels: $k(x, \xi) = x^\top \xi$, $k(x, \xi) = e^{-\gamma \|x - \xi\|^2}$

Kernel regression

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

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

$$\text{VC-dim}(F) := \max_{S} \sup_{|S|=s} \mathbf{1}[\|F(S)\| = 2^s]$$

$$\text{VC inequality: } \mathbb{P}[\sup_F |\hat{E}(f) - E(f)| > \epsilon] \leq 8|F(s)|e^{-se^2/32}$$

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

$\Delta := \max(0, E - \hat{E})$, E : expected population error, \hat{E} : empirical

Double descent

Beyond interpolation point, models may level out at lower generalization error.

KL divergence

$$D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E}_{x \sim p} [\ln \frac{p(x)}{q(x)}]$$

9.5 Loss Landscape

Critical points

$\nabla \mathcal{L}(\theta^*) = 0$. Local min: $H \succeq 0$. Saddle: H indefinite.

Sharpness

$\lambda_{\max}(H)$. Flat minima \rightarrow 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

$$x \mapsto z = Cx, z \mapsto \hat{x} = Dz, E(C, D)(x) = \frac{1}{2} \|x - DCx\|^2$$
$$DCx = \hat{X} = U\Sigma_m V^\top, \text{ for centered data } \equiv \text{PCA}$$

Linear factor analysis

$$x = \mu + Wz + \eta, \eta \sim \mathcal{N}(0, \Psi), x \sim \mathcal{N}(\mu, WW^\top + \Psi)$$
$$\text{for } z \sim \mathcal{N}(0, I)$$
$$\mu_{z|x} = W^\top(WW^\top + \Psi)^{-1}(x - \mu)$$

Generative model

$$p_\theta(x, z) = p_\theta(x|z)p(z), p(z) = \mathcal{N}(0, I)$$

ELBO

$$\log p(\theta)(x) = \log \int q(z)p(\theta)(x|z) \frac{p(z)}{q(z)} dz$$
$$\geq \int q(z) \log p(\theta)(x|z) dz - D_{KL}(q||p) =: L(\theta, q)(x)$$

Inference network

$$z \sim \mathcal{N}(\mu(x), \Sigma(x))$$

Encoder

$$q_\phi(z|x) = \mathcal{N}(\mu_\phi(x), \sigma_\phi^2(x))$$

Decoder

$$p_\theta(x|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 I) \text{ or Bernoulli}$$

Reparameterization trick

$$z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon, \epsilon \sim \mathcal{N}(0, I)$$

KL divergence (Gaussian)

$$D_{KL} = \frac{1}{2} \sum_j (\mu_j^2 + \sigma_j^2 - \ln \sigma_j^2 - 1)$$

β -VAE

$\mathcal{L} = \mathbb{E}_q[\ln p(x|z)] - \beta D_{KL}(q||p)$. $\beta > 1$ for disentanglement.

10.2 Generative Adversarial Networks

Generator

$$G : z \mapsto x, z \sim p(z)$$

Discriminator

$D : x \mapsto [0, 1]$, probability that x is real

GAN objective

$$V(G, D) = \mathbb{E}_{x_r \sim p_{\text{data}}} [D(x_r)] + \mathbb{E}_{z \sim p_z} [1 - D(G(z))]$$

Bayes-optimal classifier

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

Jensen-Shannon objective

$$\ell^* = \text{JS}(p, p_\theta) - \ln 2$$

$\ell^*(\theta) \geq \sup_\phi \ell(\theta, \phi)$ where ϕ : discriminator, θ : generator

Alternating gradient descent

$$\theta_{t+1} = \theta_t - \eta \nabla_\theta \ell(\theta_t, \phi_t)$$

$$\phi_{t+1} = \phi_t + \eta \nabla_\phi \ell(\theta^{t+1}, \phi_t)$$

Wasserstein GAN

$$\min_G \max_{D \in 1\text{-Lip}} \mathbb{E}_x[D(x)] - \mathbb{E}_z[D(G(z))]$$

Gradient penalty (WGAN-GP)

$$\lambda \mathbb{E}_{\hat{x}}[(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2], \hat{x} = \alpha x + (1 - \alpha)G(z)$$

10.3 Denoising Diffusion

Forward process

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$$

Marginal

$$q(x_t|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, \bar{\beta}_t I), \bar{\alpha}_t = \prod_{\tau=1}^t (1 - \beta_\tau), \bar{\beta}_t = \frac{1}{1 - \bar{\alpha}_t}$$

$$\nu_t \approx \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, \bar{\beta}_t I) \xrightarrow{t \rightarrow \infty} \mathcal{N}(0, I)$$

Reverse process

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(m_\theta(x_t, t), \Sigma(x_t, t))$$

Forward: $\pi^* \rightarrow \nu_T = \pi$. Backward: $\pi \rightarrow \mu_0^\theta \approx \pi^*$

Training objective

$$L_t = \frac{\|m(x_t, x_0, t) - m_\theta(x_t, t)\|^2}{2\sigma_t^2} + \text{const}$$

Reparameterization: $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$

Simplified criterion

$$h(\theta)(x) = \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2]$$

Forward trajectory target

$$x_{t-1}|x_t, x_0 = \mathcal{N}(m(x_t, x_0, t), \tilde{\beta}_t I)$$

$$m(x_t, x_0, t) = \frac{\sqrt{\bar{\alpha}_{t-1}\beta_t}}{1 - \bar{\alpha}_t} x_0 + \frac{(1 - \bar{\alpha}_{t-1})\sqrt{1 - \beta_t}}{1 - \bar{\alpha}_t} x_t$$

Sampling

$$m(x_t, x_0, t) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

Score function

$$\nabla_x \ln p(x) \approx -\frac{\epsilon_\theta(x, t)}{\sqrt{1 - \bar{\alpha}_t}}$$

Classifier-free guidance

$$\tilde{\epsilon} = (1 + w)\epsilon_\theta(x_t, t, c) - w\epsilon_\theta(x_t, t)$$

DDIM (deterministic)

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_\theta}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_\theta$$

11. Ethics

11.1 Adversarial Examples

Adversarial perturbation

$$x' = x + \delta, \|\delta\|_p \leq \epsilon, F(x') \neq F(x)$$

$$\|x\|_p = (\sum_i |x_i|^p)^{1/p}, \|x\|_\infty = \max_i |x_i|, \|x\|_0 = |\{i : x_i \neq 0\}|$$

FGSM (Fast Gradient Sign Method)

$$\delta = \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(f(x), y))$$

PGD (Projected Gradient Descent)

$$\eta^{t+1} = \Pi_\epsilon[\eta^t + \alpha \nabla_x \mathcal{L}(f(x + \eta^t), y)]$$

$$\Pi_\epsilon[z] := z/\|z\|_2 \text{ for } p = 2, \Pi_\epsilon[z] := z/\|z\|_\infty \text{ for } p = \infty$$

Adversarial training

$$\min_\theta \mathbb{E}_{x,y} [\max_{\|\delta\| \leq \epsilon} \mathcal{L}(f(x + \delta), y; \theta)]$$

DeepFool (binary)

Optimal perturbation: $\eta \propto \text{sign}(f_1(x) - f_2(x))(w_2 - w_1)$ for $f_i = w_i^\top x + b_i$. Iterate: $\arg \min_{\|\eta\|_2} \text{s.t. } (\nabla f_1 - \nabla f_2)^\top \eta < f_1(x) - f_2(x)$

Certified robustness

Provable guarantees via randomized smoothing, interval bound propagation.

Transferability

Adversarial examples often transfer between models.