

Deep Learning Cheat Sheet

Authors

January 2026

1 History

1.1 Perceptron

Threshold Unit

$$f[w, b](x) = \text{sign}(x^\top w + b)$$

Decision Boundary

$$x^\top w + b \neq 0 \Leftrightarrow \frac{x \cdot w}{\|w\|} + \frac{b}{\|w\|} \neq 0$$
$$x^\top w - b = [x \ -1] \cdot \begin{bmatrix} w \\ b \end{bmatrix} =: \tilde{x}^\top \tilde{w}, \quad \tilde{x}, \tilde{w} \in \mathbb{R}^{n+1}.$$

Geometric Margin

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

Maximum Margin Classifier

$$(w^*, b^*) \in \arg\max_{w, b} \gamma[w, b](S),$$

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 $b \leftarrow b + y$

$$w_0 \in \text{span}(x_1, \dots, x_s) \Rightarrow w_t \in \text{span}(x_1, \dots, x_s)$$

(forall t)

Convergence

If exists w , $\|w\| = 1$, such that $\gamma[w](S) = \gamma > 0$ then $w_t * w \geq t\gamma$.

$$R = \max_{x \in S} \|x\| \Rightarrow \|w_t\| \leq R\sqrt{t}$$
$$\cos(u, w_t) = \frac{u \cdot w_t}{\|w_t\|} \geq \frac{t\gamma}{\sqrt{t}R} = \sqrt{t} \frac{\gamma}{R} \leq 1 \Rightarrow t \leq \frac{R^2}{\gamma^2}$$

Cover's 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

$$C(s, n) = 2^s \text{ for } s \leq n$$

Phase transition at $s = 2n$. For $s > 2n$ empty version space is the exception, otherwise the rule.

1.1.1 Hopfield Networks

Hopfield Model

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

where $X_i \in \{-1, +1\}$

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

Hebbian Learning

$$x^t \in \{+1, -1\}^n \quad (1 \leq t \leq s)$$

$$w_{ij} = \frac{1}{n} \sum_{t=1}^s x_i^t x_j^t$$

$$W = \frac{1}{n} \sum_{t=1}^s x^t (x^t)^\top$$

1.2 Feedforward Networks

1.2.1 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

$$w^* = X^* y \in \arg\min_w h[w]$$

where $X^* := \lim_{\delta \rightarrow 0} (X^\top X + \delta I)^{-1} X^\top$

Stochastic gradient descent update

$$w_{t+1} := w_t + \eta(y_{i_t} - w_t^\top x_{i_t})x_{i_t}$$

with $i_t \sim \text{Uniform}(1, \dots, s)$

Gaussian noise model

$$y_i = w^\top x_i + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

Least squares equivalent to negative log likelihood of gaussian noise model

Ridge regression

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

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 with cross entropy loss

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

1.2.2 Feedforward Networks

Generic feedforward layer definition

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

Composition of layers

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

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

Layer activations

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

identifying $x^0 = x$, $x^L = F(x)$

Softmax function

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}, \text{softmax}(A)_{ij} = \frac{e^{A_{ij}}}{\sum_k e^{A_{ik}}}$$
$$\ell(y, z) = \frac{-z_y + \log \sum_j e^{z_j}}{\ln 2}$$

Residual layer definition

$$F[W, b](x) = x + (\varphi(Wx + b) - \varphi(0))$$

therefore $F[0, 0] = \text{id}$

Skip connection: Concatenate previous layer back in

1.2.3 Sigmoid Networks

Sigmoid activation

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

Hyperbolic tangent activation

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

Baron's Theorem: Approximation error

For f with finite $C_f := \int \|\omega\| |\hat{f}'(\omega)| d\omega < \infty$ there exists MLP g with one hidden layer of width m that:

$$\int_B (f(x) - g_m(x))^2 \mu(dx) \leq O\left(\frac{1}{m}\right)$$

1.2.4 ReLU Networks

ReLU activation

$$\varphi(z) := (z)_+ := \max\{0, z\}$$

ReLU networks are universal function approximators

Zalabsky: Connected regions

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

Montufar: Connected regions in ReLU network

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

1.3 Gradient-Based Learning

1.3.1 Backpropagation

Parameter derivatives for ridge function layers

$$\frac{\partial x_i^l}{\partial w_{ij}^{l-1}} = \dot{\varphi}_i^l x_j^{l-1},$$
$$\dot{\varphi}_i^l := \dot{\varphi}^l \left((w_i^l)^\top x^{l-1} + b_i^l \right)$$
$$\frac{\partial x_i^l}{\partial b_i^l} = \dot{\varphi}_i^l$$

Loss derivatives

$$\frac{\partial h[\theta](x, y)}{\partial w_{ij}^{l-1}} = \frac{\partial h^{l[\theta]}(x^l, y)}{\partial w_{ij}^{l-1}} \frac{\partial x_i^l}{\partial w_{ij}^{l-1}} = \delta_i^l \dot{\varphi}_i^l x_j^{l-1},$$
$$\frac{\partial h[\theta](x, y)}{\partial b_i^l} = \frac{\partial h^{l[\theta]}(x^l, y)}{\partial x_i^l} \frac{\partial x_i^l}{\partial b_i^l} = \delta_i^l \dot{\varphi}_i^l$$

with $\delta_i^l = \frac{\partial h}{\partial x_i^l} \dot{\varphi}_i^l$

1.3.2 Gradient Descent

Gradient descent update

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

Gradient flow ODE

$$d_\alpha^\theta t = -\nabla h(\theta)$$

L-smoothness

$$\|\nabla h(\theta_1) - \nabla h(\theta_2)\| \leq L \|\theta_1 - \theta_2\| \text{ (forall } \theta_1, \theta_2)$$
$$\lambda_{\max}(\nabla^2 h) \leq L$$

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

Polyak-Lojasiewicz condition

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

Convergence rate

$$\eta = \frac{1}{L}$$
$$\Delta t = \frac{2L}{\varepsilon} (h(\theta_0) - \min h) \text{ for } \varepsilon\text{-critical point}$$
$$\Delta h(\theta_t) - \min h \leq \left(1 - \frac{\mu}{L}\right)^t (h(\theta_0) - \min h)$$

1.3.3 Acceleration and Adaptivity

Heavy ball momentum update

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

Nesterov acceleration

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

More theoretical grounding than heavy ball

AdaGrad updates

$$\theta_{i,t+1} = \theta_{i,t} - \eta_i^t \frac{\partial h}{\partial \theta_i}(\theta_t),$$
$$\gamma_i^t = \gamma_i^{t-1} + \left(\frac{\partial h}{\partial \theta_i}(\theta_t) \right)^2,$$
$$\eta_i^t = \frac{\eta}{\sqrt{\gamma_i^t + \delta}}$$

Adam updates

$$g_i^t = \beta g_i^{t-1} + (1 - \beta) \frac{\partial h}{\partial \theta_i}(\theta_t)$$
$$\gamma_i^t = \alpha \gamma_i^{t-1} + (1 - \alpha) \left(\frac{\partial h}{\partial \theta_i}(\theta_t) \right)^2$$
$$\theta_{i,t+1} = \theta_{i,t} - \eta_i^t g_i^t, \eta_i^t := \frac{\eta}{\sqrt{\gamma_i^t + \delta}}$$

RMSprop

Adam without momentum term

1.3.4 Stochastic Gradient Descent

Stochastic gradient descent update

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

SGD variance

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

SGD convergence rate

$$E[h(|(\theta)_t|)] - \min h \leq O\left(\frac{1}{\sqrt{t}}\right) \text{ (general)}$$
$$E[h(|(\theta)_t|)] - \min h \leq O\left(\log \frac{t}{t}\right) \text{ (strongly convex)}$$
$$E[h(|(\theta)_t|)] - \min h \leq O\left(\frac{1}{t}\right) \text{ (additionally smooth)}$$

1.3.5 Function Properties

Convexity

$$\ell(\lambda w + (1 - \lambda)w') \leq \lambda \ell(w) + (1 - \lambda) \ell(w')$$
$$\ell''(x) \geq 0 \text{ forall } x$$

Convexity and differentiability

$$\ell(w) \geq \ell(w') + \nabla \ell(w')^\top (w - w')$$

Implies convexity for differentiable functions and vice versa

Strong convexity and differentiability

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

1.4 Convolutional Networks

1.4.1 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 transform convolution property

$$F(f * g) = F(f) * F(g)$$

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

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

1.4.2 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

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

r : output channel, u : input channel

Number of parameters of a convolutional layer

$$D = (|r| * |u|) * (|\Delta s| * |\Delta t|)$$

fully connected · window size

1.4.3 Natural Language Processing with ConvNets

Word embedding

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

Conditional log-bilinear model

Prediction of output word μ given word w in neighborhood

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

$$h(\{x_w\}, \{y_\mu\}) = \sum_{(w, \mu)} \ell_{w, \mu}$$

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

Negative sampling

$$\ell_{w, \mu} = -\ln \sigma(x_w^\top y_\mu) - \beta E_{\mu \sim D} \ln(1 - \sigma(x_w^\top y_\mu))$$

1.5 Recurrent Networks

1.5.1 Simple Recurrent Networks

Time evolution equation

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

Output map

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

RNN parameterization

$$F[U, V](z, x) := \varphi(Uz + Vx)$$

$$G[W](z) := \psi(Wz), W \in \mathbb{R}^{q \times m}$$

Backpropagation through time

$$\frac{\partial h}{\partial z_t^T} = \sum_{s=t}^T \delta_k^s \sum_{j=1}^m \frac{\partial y_k^s}{\partial z_j^s} \frac{\partial z_j^s}{\partial z_t^T},$$

$$\frac{\partial y_k^s}{\partial z_j^s} = \psi_k^s w_{kj}$$

$$\frac{\partial h}{\partial v_{ij}} = \sum_{t=1}^T \frac{\partial h}{\partial z_t^T} \dot{\varphi}_i^t x_j^t$$

$$\frac{\partial h}{\partial u_{ij}} = \sum_{t=1}^T \frac{\partial h}{\partial z_t^T} \dot{\varphi}_i^t z_j^{t-1}$$

Spectral norm

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

Gradient norms

$$\frac{\partial z^T}{\partial z^0} = \dot{\Phi}^T U * \dots * \dot{\Phi}^1 U$$

The norm of gradients either:

1. Vanishes exponentially if $\sigma_1(U) < \frac{1}{|(\alpha)|}$:

$$\left\| \frac{\partial z^t}{\partial z^0} \right\|_2 \leq (|(\alpha) \sigma_1(U)|)^t \rightarrow 0$$

2. Explodes if $\sigma_1(U)$ is too large

Bidirectional RNNs

$$\hat{y}_t = \psi(W z_t + \tilde{W} \tilde{z}_t)$$

1.5.2 Gated Memory

LSTM

$$z_t := \sigma(F \tilde{x}_t) * z_{t-1} + \sigma(G \tilde{x}_t) * \tanh(V \tilde{x}_t)$$

$$\tilde{x}_t := \text{mat}(x_t; h_t), h_{t+1} = \sigma(H \tilde{x}_t) * \tanh(U z_t)$$

GRU

$$z_t = (1 - \sigma) * z_{t-1} + \sigma * \tilde{z}_t,$$

$$\sigma := \sigma(G[x_t, z_{t-1}])$$

$$\tilde{z}_t := \tanh(V[r_t * z_{t-1}, x_t])$$

$$r_t := \sigma(H[z_{t-1}, x_t])$$

1.5.3 Linear Recurrent Models

Linear state evolution

$$z_{t+1} = A z_t + B x_t$$

Diagonal form

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

Stability condition

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

Initialization

$$\lambda_i = \exp(-\exp(\kappa_i) + i \varphi_i),$$

$$e^{\kappa_i} = -\ln r_i$$

$$\varphi_i \sim \text{Uni}[0; 2\pi], r_i \sim \text{Uni}[I], I \subset [0; 1]$$

Advantages

(i) clear modeling of long/short range dependencies

(ii) no channel mixing required

(iii) parallelizable training

1.6 Attention and Transformers

1.6.1 Attention

Attention mixing

$$\xi_s := \sum_t \alpha_{st} W x_t, \alpha_{st} \geq 0, \sum_t \alpha_{st} = 1$$

$$A = (a_{st}) \in \mathbb{R}^{T \times T}, \text{ s.t. } \Xi = W X A^\top$$

Query-key matching

$$Q = U_Q X, K = U_K X$$

$$(U_Q, U_K \in \mathbb{R}^{q \times n})$$

$$Q^\top K = X^\top U_Q^\top U_K X \text{ rank} \leq q$$

$$(Q^\top K \in \mathbb{R}^{T \times T})$$

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 $\beta = \frac{1}{\sqrt{q}}$

Feature transformation

$$X \mapsto \Xi \mapsto F(\Xi),$$

$$F(\theta)(\Xi) = (F(\xi_1), \dots, F(\xi_T))$$

Positional encoding

$$p_{tk} = \text{cases}(\sin(\omega_k t), k \text{ even}; \cos(\omega_k t), k \text{ odd}),$$

$$\omega_k = C^{\frac{k}{K}}$$

Transformer architecture

Self-attention: attend to its own values in the past

Cross-attention: E.g. decoder attends to encoder output (query from decoder, key and value from encoder)

Vision transformer patch embedding

$$\mathbb{R}^{p \times p \times q} \ni \text{patch}_t \mapsto x_t := V \text{vec}(\text{patch}_t) \in \mathbb{R}^n$$

with $V \in \mathbb{R}^{n \times (qp^2)}$

GELU activation

$$\varphi(z) = z \text{Prob}(z \leq Z), Z \sim \mathcal{N}(0, 1)$$

1.7 Geometric Deep Learning

1.7.1 Sets and Points

Function over sets

$$\{x_1, \dots, x_M\} \subset \mathbb{R}, f : 2^{\mathbb{R}} \rightarrow Y$$

Order-invariance property

$$f(x_1, \dots, x_M) = f(x_{\pi(1)}, \dots, x_{\pi(M)}) \text{ for all } \pi \in S_M$$

Equivariance property

$$f(x_1, \dots, x_M) = (y_1, \dots, y_M) \Rightarrow$$

$$f(x_{\pi(1)}, \dots, x_{\pi(M)}) = (y_{\pi(1)}, \dots, y_{\pi(M)})$$

Permutation invariant sum

$$\sum_{m=1}^M x_m = \sum_{m=1}^M x_{\pi(m)}, \text{ for all } M, \text{ for all } \pi \in S_M$$

Deep Sets model

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

Max pooling variant

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

Equivariant map construction

$$\rho : \mathbb{R} \times \mathbb{R}^N \rightarrow Y,$$

$$(x_m, \sum_{k=1}^M \varphi(x_k)) \mapsto y_m$$

1.7.2 Graph Convolutional Networks

Feature and adjacency matrices

$$X = \text{mat}(x_1^\top; \dots; x_M^\top), A = (a_{nm})$$

with

cases(1, if $\{v_n, v_m\} \in E$; 0, otherwise)

$$a_{nm} =$$

Permutation matrix constraints

$$P \in \{0, 1\}^{M \times M} \text{ s.t.}$$

$$\sum_{n=1}^M p_{nm} = \sum_{n=1}^M p_{mn} = 1 \text{ (for all } m)$$

Graph invariance definition

$$f(X, A) \neq f(PX, PAP^\top), \text{ for all } P \in \Pi_M$$

Graph equivariance definition

$$f(X, A) \neq Pf(PX, PAP^\top), \text{ for all } P \in \Pi_M$$

Node neighborhood features

$$X_m := \{\{x_n : \{v_n, v_m\} \in E\}\},$$

$$\{\{\cdot\}\} =$$

multiset

Message passing scheme

$$\varphi(x_m, X_m) = \varphi(x_m, m_{X_m} \psi(x))$$

m is a permutation-invariant operation

Normalized adjacency matrix

$$|(A) = D^{-\frac{1}{2}}(A + I)D^{-\frac{1}{2}}$$

$$D = \text{diag}(d_1, \dots, d_M), d_m = 1 + \sum_{n=1}^M a_{nm}$$

GCN layer

$$X^+ = \sigma(|(A)XW), W \in \mathbb{R}^{M \times N}$$

Two-layer GCN

$$Y = \text{softmax}(|(A)(|(A)XW^0)W^1)$$

1.7.3 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

$$L = D - A, (Lx)_n = \sum_{m=1}^M a_{nm}(x_n - x_m)$$

Normalized Laplacian

$$\tilde{L} = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}} = D^{-\frac{1}{2}} (D - A) D^{-\frac{1}{2}}$$

Graph Fourier transform

$$L = D - A = U \Lambda U^\top, \\ \Lambda := \text{diag}(\lambda_1, \dots, \lambda_M), \lambda_i \geq \lambda_{i+1}$$

Convolution

$$x * y = U((U^\top x) \odot (U^\top y))$$

Filtering operation

$$G_{\theta(L)} x = U G_{\theta(\Lambda)} U^\top x$$

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

$$x_i^{l+1} = \sum_j p_{ij}(L) x_j^l + b_i,$$

$$p_{ij}(L) = \sum_{k=0}^K \alpha_{ijk} L^k$$

1.7.4 Attention GNNs

Attention coupling matrix

$$Q = (q_{ij}), \\ q_{ij} = \text{softmax}(\rho(u^\top (V x_i; V x_j; x_{ij}))) \\ \text{s.t. } \sum_j A_{ij} q_{ij} = 1$$

Attention propagation

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

Weisfeiler-Lehman test

1.8 Tricks of the Trade

1.8.1 Initialization

Random initialization

$$\theta_i^0 \sim \mathcal{N}(0, \sigma_i^2), \text{ or } \\ \theta_i^0 \sim \text{Uniform}(-\sqrt{3}\sigma_i; \sqrt{3}\sigma_i)$$

LeCun initialization

$$w_{ij} \sim \text{Uniform}[-a; a], a := \frac{1}{\sqrt{n}}, b_i = 0$$

Stabilizes variance

Glorot initialization

$$w_{ij} \sim \text{Uniform}[-\sqrt{3}\gamma; \sqrt{3}\gamma],$$

$$\gamma := \frac{2}{n+m} \\ \text{Stabilizes variance of gradients in backpropagation}$$

He initialization

$$w_{ij} \sim \mathcal{N}(0, \gamma) \text{ or } w_{ij} \sim \text{Uniform}[-\sqrt{3}\gamma; \sqrt{3}\gamma],$$

$$\gamma := \frac{2}{n} \\ \text{In ReLU networks typically only } \frac{n}{2} \text{ units active}$$

Orthogonal initialization

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

1.8.2 Weight Decay

L2 regularization

$$\Omega_{\mu(\theta)} = \frac{\mu}{2} \|\theta\|^2, \mu \geq 0$$

Gradient descent with weight decay

$$\Delta\theta = -\eta \nabla E(\theta) - \eta \nabla \Omega_{\mu(\theta)} = -\eta \nabla E(\theta) - \eta \mu \theta$$

Weight decay for multiple layers

$$\theta = (\text{vec}(W^1), \text{vec}(W^2), \dots, \text{vec}(W^L)),$$

$$\Omega_{\mu(\theta)} = \sum_{l=1}^L \mu_l \|W^l\|_F^2$$

Local loss landscape

$$\theta_\mu^* = (H + \mu I)^{-1} H \theta^*, H = Q^\top \Lambda Q$$

$$(\Lambda + I)^{-1} \Lambda = \text{diag}\left(\frac{\lambda_i}{\lambda_i + \mu}\right)$$

The minimum θ^* is shrunk along directions with small eigenvalues

Generalization

$$\mu = \frac{\sigma^2}{u^2}, u: \text{teacher "sign" al}$$

Optimal weight decay inverse proportional to the "sign" al-to-noise ratio

1.8.3 Dropout

Probability φ_i of keeping a unit

Dropout as Ensembling

$$p(y | x) = \sum_{b \in \{0,1\}^R} p(b) p(y | x; b)$$

$$\text{with } p(b) = \prod_{i=1}^R \varphi_i^{b_i} (1 - \varphi_i)^{1-b_i}$$

Weight scaling for inference

$$\tilde{w}_{ij} \leftarrow \varphi_j w_{ij}$$

1.8.4 Normalization

Batch normalization

E and V from minibatches or population statistics

$$|(f) = \frac{f - E[f]}{\sqrt{V[f]}}, E[|(f)] = 0, V[|(f)] = 1$$

$$|(f)[\mu, \gamma] = \mu + \gamma |(f)$$

Weight normalization

$$f(v, \gamma)(x) = \varphi(w^\top x), w := \frac{\gamma}{\|v\|_2} v$$

Gradient descent with respect to decoupled γ and v :

$$\frac{\partial E}{\partial \gamma} = \nabla_w E * \frac{v}{\|v\|_2}$$

$$\nabla_v E = \frac{\gamma}{\|v\|} \left(I - \frac{w w^\top}{\|w\|^2} \right) \nabla_w E$$

Layer normalization

$$\hat{f}_i = \frac{f_i - E[f]}{\sqrt{V[f]}},$$

$$E[f] = \frac{1}{m} \sum_{i=1}^m f_i$$

$$V[f] = \frac{1}{m} \sum_{i=1}^m (f_i - E[f])^2$$

Using population averages across units in a layer

1.8.5 Model Distillation

Tempered cross entropy loss for distillation

$$\ell(x) = \sum_{y=1}^K \frac{\exp\left[\frac{F_y(x)}{T}\right]}{\sum_{\mu=1}^K \exp\left[\frac{F_\mu(x)}{T}\right]} \left[\frac{1}{T} G_{y(x)} - \right.$$

$$\ln \sum_{\mu=1}^K \exp\left[\frac{G_\mu(x)}{T}\right]$$

$T > 0, F_y$: teacher logits, G_y : student logits

Gradient of distillation loss

$$\frac{\partial \ell}{\partial G_y} = \frac{1}{T} \left[\frac{e^{\frac{F_y}{T}}}{\sum_\mu e^{\frac{F_\mu}{T}}} - \frac{e^{\frac{G_y}{T}}}{\sum_\mu e^{\frac{G_\mu}{T}}} \right]$$

1.9 Theory

1.9.1 Neural Tangent Kernel

Linearized DNN taylor approximation

$$h(\beta)(x) = f(x) + \beta * \nabla f(x)$$

$$\text{with } \beta \approx \theta - \theta_0, f(x) := f(\theta_0)(x)$$

Kernel of gradient feature maps

$$k(x, \xi) = \nabla f(x) * \nabla f(\xi), \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

Dual representation

$$h(\alpha)(x) = f(x) + \sum_{i=1}^s \alpha_i \nabla f(x_i) * \nabla f(x)$$

Squared loss

$$E(\alpha) = \frac{1}{2s} \sum_{i=1}^s \left(\sum_{j=1}^s \alpha_j \nabla f(x_j) * \nabla f(x_i) + f(x_i) - \frac{1}{2} \log p(\theta | S) = E(\theta) + \text{const}, \right. \\ \left. \tilde{E}(\theta) = \frac{1}{2\gamma^2} \|y - f\|^2 + \frac{1}{2\sigma^2} \|\theta\|^2 \right)$$

Optimal solution of linearized DNN

$$K = [k(x_i, x_j)]_{i,j=1}^n \in \mathbb{R}^{n \times n}$$

$$\alpha^* = K^+(y - f),$$

$$h^{*(x)} = k(x) K^+(y - f)$$

Neural Tangent Kernel NTK

$$k(\theta)(x, \xi) := \nabla f(\theta)(x) * \nabla f(\theta)(\xi)$$

Quadratic loss

$$E(\theta) = \frac{1}{2} \|f(\theta) - y\|^2, y := (y_1, \dots, y_s)^\top$$

Gradient flow ODE

$$\dot{\theta} := \text{d}_\theta t = \sum_{i=1}^s (y_i - f_i(\theta)) \nabla f_i(\theta)$$

Functional gradient flow

$$\dot{f}_j = \nabla f_j * \theta = \sum_{i=1}^s (y_i - f_i) k(\theta)(x_i, x_j)$$

$$f = K(\theta)(y - f)$$

Infinite width limit

$$w_{ij}^l = \frac{\sigma_w}{\sqrt{m_i}} \varepsilon_{ij}^l,$$

$$b_i^l = \frac{\sigma_b}{\sqrt{m_i}} \beta_i^l,$$

$$\varepsilon_{ij}^l, \beta_i^l \sim \mathcal{N}(0, 1)$$

$$k(\theta) \rightarrow k_\infty \text{ for } m_l \rightarrow \infty$$

Initial NTK converges to deterministic limit

NTK constancy

$$\text{d} \frac{k(\theta(t))}{\text{d} t} = 0$$

$$f_{\infty(x)} = k(x) K^+(y - f), k = k_\infty$$

NTK remains constant when training in infinite width limit

Vanishing curvature

$$\frac{\|\nabla^2 f(\theta_0)\|_2}{\|\nabla f(\theta_0)\|_2^2} \ll 1$$

Near-constancy

$$\|k(\theta_0) - k(\theta_t)\|_F \in O\left(\frac{1}{m}\right), m = m_1 = \dots = m_L$$

1.9.2 Bayesian DNNs

Bayesian predictive distribution

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

Bayes rule

$$p(\theta | S) = \frac{p(\theta) p(S | \theta)}{p(S)},$$

$$p(S) = \int p(\theta) p(S | \theta) d\theta$$

Parameter priors (Gaussian)

$$p(\theta) = \prod_{i=1}^d p(\theta_i), \theta_i \sim \mathcal{N}(0, \sigma_i^2)$$

$$- \log p(\theta) = \frac{1}{2\sigma^2} \|\theta\|^2 + \text{const}$$

Essentially a weight decay term

Likelihood (Gaussian noise)

$$- \log p(S | \theta) = \frac{1}{2\gamma^2} \|y - f(\theta)\|^2 + \text{const.}$$

$$\text{with } y_i = f^{*(x_i)} + \nu_i, \nu_i \sim \mathcal{N}(0, \gamma^2)$$

Posterior

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

Bayesian ensembling (post hoc)

$$f(\Theta)(x) = \sum_{i=1}^n \frac{\exp\left[\frac{-E(\theta_i)}{\gamma}\right]}{\sum_{j=1}^n \exp\left[\frac{-E(\theta_j)}{\gamma}\right]} f(\theta_i)(x)$$

Relative posterior weighting

Markov chain monte carlo (MCMC)

$$\theta_0, \theta_1, \theta_2, \dots,$$

$$\theta_{t+1} | \theta_t \sim \Pi$$

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

Metropolis-Hastings

$$\Pi(\theta_1 | \theta_2) = \tilde{\Pi}(\theta_1 | \theta_2) A(\theta_1 | \theta_2)$$

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

Modified transition probability with acceptance step A

Hamiltonian monte carlo

$$E(\theta) = - \sum_{x,y} \log p(y | x; \theta) - \log p(\theta)$$

$$H(\theta, v) = E(\theta) + \frac{1}{2} v^\top M^{-1} v$$

$$\text{with } p(\theta, v) \text{ propto } \exp[-H(\theta, v)]$$

$$\dot{v} = -E'(\theta), \theta = v$$

$$\theta_{t+1} = \theta_t + \eta v_t$$

$$v_{t+1} = v_t - \eta \nabla E(\theta_t)$$

Langevin dynamics

$$\dot{\theta} = v$$

$$dv = -\nabla E(\theta) dt - B v dt + N(0, 2B dt)$$

$$\theta_{t+1} = \theta_t + \eta v_t$$

$$v_{t+1} = (1 - \eta\gamma) v_t - \eta \int \nabla \tilde{E}(\theta) + \sqrt{2\gamma\eta} N(0, I)$$

1.9.3 Gaussian Processes

Gaussian process

$(f(x_1), \dots, f(x_s)) \sim N$
 $\sum_{i=1}^s \alpha_i f(x_i) \sim N$, for all $\alpha \in \mathbb{R}^s$

Mean and covariance functions

GPs are completely defined by first and second order statistics

$$\begin{aligned}\mu(x) &:= E_{x[f(x)]} \\ k(x, \xi) &:= E_{x, \xi}[f(x)f(\xi)] - \mu(x)\mu(\xi) \\ K_{\mu\nu} &= k(x_\mu, x_\nu), K \in \mathbb{R}^{s \times s}\end{aligned}$$

Example kernels

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

GPs in DNN

Treating parameters as random variables. Each unit in a DNN becomes a random function.

Linear Layer

$$\begin{aligned}w &\sim N\left(0, \frac{\sigma^2}{n} I_{n \times n}\right) \\ E[y_i y_j] &= \frac{\sigma^2}{n} x_i^\top x_j\end{aligned}$$

Deep layers

$$W^{l+1} X^l, l \geq 1$$

No longer normal as products break normality, but near-normal for high dimensional inputs.

Non-linear activations

$$\begin{aligned}\mu(x^{l+1}) &= E[\varphi(W^l x^l)] \\ \text{Kernel recursion} \\ K_{\mu\nu}^l &= E[\varphi(x_{i\mu}^{l-1}) \varphi(x_{i\nu}^{l-1})] \\ &= \sigma^2 E[\varphi(f_\mu) \varphi(f_\nu)] \\ f &\sim \text{GP}(0, K^{l-1})\end{aligned}$$

Kernel regression

$$\begin{aligned}\text{Mean of bayesian predictive distribution} \\ f^{*(x)} &= k(x)^\top K^+ y \\ E\left[(f(x) - f^{*(x)})^2\right] &= K(x, x) - k(x)^\top K^+ k(x)\end{aligned}$$

1.9.4 Statistical Learning Theory

VC learning theory

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

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

VC inequality

$$P(\sup_{f \in F} |\hat{E}(f) - E(f)| > \varepsilon) \leq 8 |F(s)| e^{-s \frac{\varepsilon^2}{32}}$$

Double descent

Beyond the interpolation point, models start to learn and eventually may level out at a lower generalization error.

Generalization gap

$$\Delta := \max(0, E - \hat{E})$$

E : expected population error, \hat{E} : empirical error

KL divergence

$$\begin{aligned}D_{\text{KL}}(p \parallel q) &= \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx = E_{x \sim p}\left[\ln\left(\frac{p(x)}{q(x)}\right)\right]\end{aligned}$$

PAC-Bayesian theorem

For fixed E and any Q over s samples:

$$\begin{aligned}E_{Q[E(f)]} - \\ E_{Q[\hat{E}(f)]} &\leq \sqrt{\frac{2}{s} \left[\text{KL}(Q \parallel P) + \ln\left(\frac{1}{2\sqrt{s\varepsilon}}\right) \right]}\end{aligned}$$

Ensures general rate $\tilde{O}\left(\frac{1}{\sqrt{s}}\right)$

PAC-Bayesian bound

$$\begin{aligned}Q &:= N(\theta, \text{diag}(\sigma_i^2)) \\ \text{KL}(Q \parallel P) &= \sum_i \log\left(\frac{\lambda}{\sigma_i} + \frac{\sigma_i^2 + \theta_i^2}{2\lambda^2} - \frac{1}{2}\right) \\ E_{\text{PAC}}(Q) &:= \\ E_{Q[\hat{E}]} &+ \sqrt{\frac{2}{s} \left[\text{KL}(Q \parallel P) + \ln\left(\frac{1}{2\sqrt{s\varepsilon}}\right) \right]}\end{aligned}$$

Favours minima robust to parameter perturbations

PAC-bayesian learning implementation

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \nabla E_{Q[\hat{E}]} = \theta_t - \eta \nabla \hat{E}(\tilde{\theta}), \\ \text{with } \tilde{\theta} &\sim Q(\theta, \sigma) \\ \text{Gradient loss on perturbed parameters}\end{aligned}$$

Reparameterization trick

$$\begin{aligned}\tilde{\theta} &= \theta + \text{diag}(\sigma_i) \varepsilon, \varepsilon \sim N(0, I) \\ \text{Backpropagation to } \theta &\text{ and } \sigma_i\end{aligned}$$

1.10 Generative Models

1.10.1 Variational Autoencoders

Linear autoencoder

$$\begin{aligned}x &\mapsto z = Cx, C \in \mathbb{R}^{m \times n} \\ z &\mapsto \hat{x} = Dz, D \in \mathbb{R}^{n \times m} \\ E(C, D)(x) &= \frac{1}{2} \|x - \hat{x}\|^2 = \frac{1}{2} \|x - DCx\|^2 \\ DCX &= \hat{X} = U \Sigma_m V^\top \\ \Sigma_m &= \text{diag}(\sigma_1, \dots, \sigma_m, 0, \dots, 0) \\ \text{For centered data equivalent to PCA, but generally} \\ &\text{has non-global minima}\end{aligned}$$

Linear factor analysis

Probability Model

$$\begin{aligned}p_X(x) &= \int p_{Z(z)} p_{X|Z}(x | z) dz \\ Z: &\text{latent variables, } X: \text{observed variables}\end{aligned}$$

Linear observation model

$$\begin{aligned}x &= \mu + Wz + \nu \text{ with } \nu \sim N(0, \Sigma) \\ x &\sim N(\mu, WW^\top + \Sigma) \text{ for } z \sim N(0, I)\end{aligned}$$

Posterior mean and covariance

$$\begin{aligned}\mu_{z|x} &= W^\top (WW^\top + \Sigma)^{-1} (x - \mu) \\ \Sigma_{z|x} &= I - W^\top (WW^\top + \Sigma)^{-1} W\end{aligned}$$

Pseudoinverse limit

$$\begin{aligned}W^\top (WW^\top + \sigma^2 I)^{-1} &\rightarrow W^+ \in \mathbb{R}^{m \times n} \\ \mu_{z|x} &\rightarrow W^+ (x - \mu), \Sigma_{z|x} \rightarrow 0\end{aligned}$$

Maximum likelihood estimation

$$\mu, W \max \rightarrow \log p_{\mu, W}(S)$$

Optimality condition for W

$$\begin{aligned}w_i &= \rho_i u_i, \rho_i = \max\{0, \sqrt{\lambda_i - \sigma^2}\} \\ \text{With } (\lambda_i, u_i) &\text{ eigenvalues and eigenvectors of co-} \\ &\text{variance matrix.} \\ \text{For } \sigma = 0 &\text{ equivalent to PCA.}\end{aligned}$$

Variational autoencoder (VAE)

$$\begin{aligned}z &\sim N(0, I) \\ x &= F(\theta)(z) = (F^L @ \dots @ F^1)(z)\end{aligned}$$

Evidence lower bound (ELBO)

$$\begin{aligned}\log p_\theta(x) &= \log \int p_\theta(x | z) p(z) dz \\ &= \log \int q(z) \left[\frac{p_\theta(x | z) p(z)}{q(z)} \right] dz \\ &\geq \int q(z) \log p_\theta(x | z) dz - \int q(z) \log \left(\frac{q(z)}{p(z)} \right) dz \\ &=: L(\theta, q)(x) \\ \theta \max \rightarrow L(\theta, q)(S) &= \sum_{i=1}^s L(\theta, q)(x_i)\end{aligned}$$

Inference network

$$\begin{aligned}z &\sim N(\mu(x), \Sigma(x)) \\ z &= \mu + \Sigma^{\frac{1}{2}} \varepsilon, \varepsilon \sim N(0, I) \\ \nabla_\mu E[f(z)] &= E[\nabla_z f(z)] \\ \nabla_\Sigma E[f(z)] &= \frac{1}{2} E[\nabla_z^2 f(z)] \\ \text{Integration by parts derivation}\end{aligned}$$

1.10.2 Generative Adversarial Networks

GAN objective

$$\begin{aligned}V(G, D) &= E_{x_r \sim p_{\text{data}}} D(x_r) + E_{z \sim p_z} (1 - \\ &D(G(z)))\end{aligned}$$

Discriminator Mixture Model

$$\begin{aligned}\hat{p}_{\theta(x, y)} &= \frac{1}{2} (yp(x) + (1 - y)p_{\theta(x)}), \\ y &\in \{0, 1\},\end{aligned}$$

p : true probability, p_θ : model probability

Bayes-optimal classifier

$$\begin{aligned}q_{\theta(x)} &:= P\{y = 1 | x\} = \frac{p(x)}{p(x) + p_{\theta(x)}} \\ \text{To detect fake samples, } y = 1 &\text{ for real samples, } y = \\ 0 &\text{ for fake samples}\end{aligned}$$

Logistic likelihood

$$\begin{aligned}\theta \min \rightarrow \ell^{*(\theta)} &:= E_{p_\theta} [y \ln q_{\theta(x)} + (1 - y) \ln (1 - \\ &q_{\theta(x)})]\end{aligned}$$

Jensen-Shannon as effective objective

$$\begin{aligned}\ell^* &= E_{p_\theta} [y \ln q_{\theta(x)} + (1 - y) \ln (1 - q_{\theta(x)})] \\ &= -\frac{1}{2} H(p) - \frac{1}{2} H(p_\theta) + H\left(\frac{1}{2}(p + p_\theta)\right) - \ln 2 \\ &= \text{JS}(p, p_\theta) - \ln 2.\end{aligned}$$

Discriminator model

$$q_\varphi: x \mapsto [0; 1], \varphi \in \Phi$$

Objective bounds

$$\begin{aligned}\ell^{*(\theta)} &\geq \sup_{\varphi \in \Phi} \ell(\theta, \varphi) \\ \ell(\theta, \varphi) &:= E_{p_\theta} [y \ln q_{\varphi(x)} + (1 - y) \ln (1 - q_{\varphi(x)})]\end{aligned}$$

Saddle point optimization

$$\begin{aligned}\theta^* &:= \text{argmin}_{\theta \in \Theta} \left(\sup_{\varphi \in \Phi} \ell(\theta, \varphi) \right) \\ \varphi: &\text{Generator, } \theta: \text{Discriminator}\end{aligned}$$

Alternating gradient descent/ascent

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \nabla_\theta \ell(\theta_t, \varphi_t) \\ \varphi_{t+1} &= \varphi_t + \eta \nabla_\varphi \ell(\theta_{t+1}, \varphi_t)\end{aligned}$$

Extra-gradient steps

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \nabla_\theta \ell(\theta_{t+0.5}, \varphi_t) \\ \text{with } \theta_{t+0.5} &:= \theta_t - \eta \nabla_\theta \ell(\theta_t, \varphi_t) \\ \varphi_{t+1} &= \varphi_t + \eta \nabla_\varphi \ell(\theta_t, \varphi_{t+0.5}) \\ \text{with } \varphi_{t+0.5} &:= \varphi_t + \eta \nabla_\varphi \ell(\theta_t, \varphi_t)\end{aligned}$$

Deconvolutional DNN

Upside-down ConvNet for image generation

1.10.3 Denoising Diffusion

Markov chains

$$\begin{aligned}x_{0:t-1} &\perp x_{t+1:\infty} | x_t \text{ (for all } t) \\ p(x_t | x_{t-1}) &= p(x_t | x_0) \text{ (for all } t), \\ p(x_{s:t}) &= p(x_t) \text{prod}_{\tau=s+1}^t p(x_{\tau-1} | x_\tau) \\ p(x_{s:t}) &= p(x_s) \text{prod}_{\tau=s+1}^t p(x_\tau | x_{\tau-1}), \\ \pi(x_{t+1}) &= \int \pi(x_t) p(x_{t+1} | x_t) dx_t\end{aligned}$$

Denoising diffusion

Forward (noise generation)

$$\pi^* = \nu_0 \mapsto \nu_1 \mapsto \dots \mapsto \nu_{T-1} \mapsto \nu_T = \pi$$

Backward (denoising)

$$\pi = \mu_0^\theta \mapsto \mu_{T-1}^\theta \mapsto \dots \mapsto \mu_1^\theta \mapsto \mu_0^\theta \approx \pi^*$$

Gaussian example

$$\begin{aligned}\pi &\approx N(0, I), \\ x_t | x_{t-1} &\sim N(\sqrt{1 - \beta_t} x_{t-1}, \beta_t I)\end{aligned}$$

Forward SDE

$$dx_t = -\frac{1}{2} \beta_t x_t dt + \sqrt{\beta_t} d\omega_t$$

Backward SDE

$$\begin{aligned}dx_t &= \left[-\frac{1}{2} \beta_t x_t - \beta_t \nabla_{x_t} \log q_t(x_t) \right] dt + \\ &\sqrt{\beta_t} d|(\omega)_t \\ \text{score} \cdot \text{wiener process}\end{aligned}$$

ELBO bound

$$\begin{aligned}x_t &= \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon_t, \varepsilon_t \sim N(0, I) \\ \ln p_{\theta(x_0)} &= \ln \int q(x_{1:T} | x_0) \left(\frac{p_{\theta(x_0, T)}}{q(x_{1:T} | x_0)} \right) dx_{1:T} \\ &\geq E \left[\ln \left(\frac{p_{\theta(x_0, T)}}{q(x_{1:T} | x_0)} \right) | x_0 \right] \\ &= \sum_{t=0}^T L_t \\ L_t &:= \text{cases} \left(E \left[\ln p_{\theta(x_0 | x_1)} \right] \right), t = \\ 0; &-D(q(x_T | x_0) \parallel \pi), t = \\ T; &-D(q(x_{t-1} | x_t, x_0) \parallel p_{\theta(x_{t-1} | x_t)}), \text{else} \end{aligned}$$

Backward model assumption

$$x_{t-1} | x_t \sim N(m(x_t, t), \Sigma(x_t, t))$$

Entropy bounds

$$H(x_t) \geq H(x_{t-1}) \Rightarrow$$

$$H(x_t \mid x_{t-1}) \geq H(x_{t-1} \mid x_t)$$

Noise schedules

$$|(\alpha)_t = \prod_{\tau=1}^t (1 - \beta_\tau), |(\beta)_t = 1 - |(\alpha)_t$$

$$x_t \approx N\left(\sqrt{|(\alpha)_t} x_0, |(\beta)_t I\right) \quad t \rightarrow \infty \rightarrow N(0, I)$$

Forward trajectory target

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

$$m(x_t, x_0, t) =$$

$$\left(\frac{\sqrt{|(\alpha)_{t-1} \beta_t}}{1 - |(\alpha)_t}\right) x_0 + \frac{(1 - |(\alpha)_{t-1}) \sqrt{1 - \beta_t}}{1 - |(\alpha)_t} x_t$$

$$\text{with } \tilde{\beta}_t = \frac{1 - |(\alpha)_{t-1} \beta_t}{1 - |(\alpha)_t}$$

Fixed isotropic covariance

$$\Sigma(x_t, t) = \sigma_t^2 I, \text{ where } \sigma_t^2 \in \{\beta_t, \tilde{\beta}_t\}$$

Simplified ELBO

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

Reparameterization

$$x_t = \sqrt{|(\alpha)_t} x_0 + \sqrt{1 - |(\alpha)_t} \varepsilon \Rightarrow x_0 =$$

$$\frac{1}{\sqrt{|(\alpha)_t}} x_t - \frac{\sqrt{1 - |(\alpha)_t}}{\sqrt{|(\alpha)_t}} \varepsilon$$

$$m(x_t, x_0, t) = \frac{1}{\sqrt{\alpha_t}} \left[x_{t(x_0, \varepsilon)} - \frac{\beta_t}{\sqrt{1 - |(\alpha)_t}} \varepsilon \right]$$

$$\text{with } \varepsilon \sim N(0, I)$$

Expected squared error

$$E_{q[L_t \mid x_0]} = E$$

$$\varepsilon \left[\rho_t \left\| \varepsilon - \varepsilon_{\theta(\sqrt{|(\alpha)_t} x_0 + \sqrt{1 - |(\alpha)_t} \varepsilon, t)} \right\|^2 \mid x_0 \right]$$

$$\text{with } \rho_t = \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - |(\alpha)_t)}$$

Final simplified criterion

$$h(\theta)(x) =$$

$$\frac{1}{T} \sum_{t=1}^T E \left[\left\| \varepsilon - \varepsilon_{\theta(\sqrt{|(\alpha)_t} x + \sqrt{1 - |(\alpha)_t} \varepsilon, t)} \right\|^2 \right]$$

1.11 Ethics**1.11.1 Adversarial Examples****Adversarial perturbation**

$$f(x + \nu) \neq f(x) \text{ s.t. } \|\nu\|_p \leq \varepsilon$$

p-norm definitions

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}$$

$$\|x\|_\infty = \max_i |x_i|, \|x\|_0 = |\{i : x_i \neq 0\}|$$

Optimal perturbation (linear binary classification)

$$\nu \propto \text{sign}(f_1(x) - f_2(x))(w_2 - w_1)$$

$$\text{for } f_i = w_i^\top x + b_i$$

Optimal perturbation (multiclass)

$$\nu = \text{argmin}_{i>1} \frac{f_1(x) - f_i(x)}{\|w_1 - w_i\|_2^2} (w_i - w_1)$$

DeepFool iterative optimization

$$\text{Iterate: } \text{argmin}_{\Delta\nu} \|\Delta\nu\|_2 \text{ s.t.}$$

$$(\nabla f_1(x) - \nabla f_2(x))^\top \Delta\nu < f_1(x) - f_2(x)$$

Robust training

$$\ell(f(x), y) \rightarrow \max_{\nu: \|\nu\|_p \leq \varepsilon} \ell(f(x + \nu), y)$$

Projected gradient ascent (p = 2)

$$\nu_{t+1} = \varepsilon \Pi[\nu_t + \alpha \nabla_x \ell(f(x + \nu_t), y)]$$

$$\Pi[z] := \frac{z}{\|z\|_2}$$

Projected gradient ascent (p = oo)

$$\nu_{t+1} = \varepsilon \Pi[\nu_t + \alpha \text{sign}(\nabla_x \ell(f(x + \nu_t), y))]$$

$$\Pi[z] := \frac{z}{\|z\|_\infty}$$

Fast Gradient Sign Method (FGSM)

$$\nu = \varepsilon \text{sign}(\nabla_x \ell(f(x), y))$$

2 Computer Vision

2.1 The Digital Image & Sensors

Charge Coupled Device (CCD)

Photons

- **Blooming:** Oversaturated photosites cause vertical channels to “flood” (bright vertical line)

Image Noise

Additive Gaussian noise:

Color camera concepts:

1. Prism (split light, 3 sensors, needs good alignment, good color separation)

2.2 Image Segmentation

Pixel-wise classification problem, to group pixels in an image that share common properties.

Segmentation of I : Find R_1, \dots, R_n such that $I = \bigcup_{i=1}^N R_i$ with $R_i \cap R_j = \emptyset \quad \forall i \neq j$.

Thresholding Segment image into 2 classes.

$B(x, y) = 1$ if $I(x, y) \geq T$ else 0, finding T with trial and error, compare results with ground truth.

Important Kernels

Laplacian	Prewitt _x	Low-pass/ Mean / High-pass Box
-----------	----------------------	--------------------------------------

Gaussian	Sobel _x	Diff _x	Diff _y
----------	--------------------	-------------------	-------------------

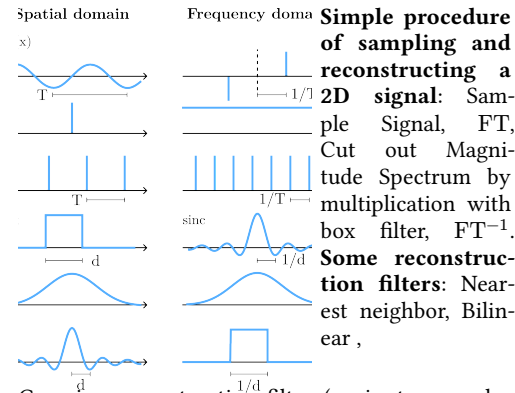
$\frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$ $\text{mat}(-1, 0, 1, \text{mat}(0, 2, -1), \text{mat}([-1], 1))^T$

Dirac delta: $\delta(x) = \begin{cases} 0 & \text{if } x \neq 0 \\ \text{undefined} & \text{else} \end{cases}$ with $\int_{-\infty}^{\infty} \delta(x) dx = 1$. $\mathcal{F}[\delta(x - x_0)](u) = e^{-i2\pi u x_0}$. $\delta(u) = \int_{\mathbb{R}} e^{-i2\pi x u} dx$.

Sampling f at points x_n : $f_s(x) = \sum_n f(x_n) \delta(x - x_n)$.

Property

Linearity	$f(x)$ $\alpha f_1(x) + \beta f_2(x)$	$F(u)$ $\alpha F_1(u) + \beta F_2(u)$
Duality	$F(x)$	$f(-u)$



Gaussian reconstruction filter (equiv. to convolving sampled signal w/ Gaussian kernel. $\text{sinc}(x) = \frac{\sin(\pi x)}{\pi x}$)

Image restoration

Image degradation is applying kernel h to some image I . The inverse \tilde{h} should compensate: $I \xrightarrow{\tilde{h}(x)} J \xrightarrow{h(x)} I$.

Determine with $\mathcal{F}[\tilde{h}](u, v) \cdot \mathcal{F}[h](u, v) = 1$. Or $\tilde{h} = \mathcal{F}^{-1}\left[\frac{1}{\mathcal{F}[h]}\right]$

Cancellation of frequencies & noise amplification \rightarrow Regularize using $\tilde{\mathcal{F}}[\tilde{h}](u, v) = \mathcal{F}[h] / (|\mathcal{F}[h]|^2 + \varepsilon)$.

Motion blur: $h(x, y) = \frac{1}{2l}[\theta(x + l) - \theta(x - l)]\delta(y)$