

# 1 Connectionism

## 1.1 Perceptron

**Threshold Unit**  $f[w, b](x) = \text{sign}(x \cdot w + b)$  with **Dec. Boundary**  $x \cdot w + b = 0 \Leftrightarrow \frac{x \cdot w}{\|w\|} + \frac{b}{\|w\|} = 0$ .

**Geometric Margin**  $\gamma[w, b](x, y) = \frac{y(x \cdot w + b)}{\|w\|}$ .

**Maximum Margin Classifier**

$(w^*, b^*) \in \arg\max_{w, b} \gamma[w, b](\mathcal{S})$ , with  $\gamma[w, b](\mathcal{S}) := \min_{(x, y) \in \mathcal{S}} \gamma[w, b](x, y)$ .

### Perception Learning

If  $f[w, b](x) \neq y$ : update  $w += yx$ , and  $b += y$ .

$w_t \in \text{span}(x_1, \dots, x_s) \Rightarrow w_t \in \text{span}(x_1, \dots, x_s) \forall t$ .

### Convergence

1. If  $\exists w^*, \|w^*\| = 1$ , s.t.  $\gamma[w^*](\mathcal{S}) = \gamma > 0 \Rightarrow w_t \cdot w^* \geq t\gamma$ .
2. Let  $R = \max_{x \in \mathcal{S}} \|x\|$ . Then  $\|w_t\| \leq R\sqrt{t}$ .

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

### Cover's Theorem for $\mathcal{S} \subset \mathbb{R}^n$ , $|\mathcal{S}| = s$

$C(\mathcal{S}, n)$ : # of ways to separate  $\mathcal{S}$  in  $n$  dimensions. Position of pts does not matter (general position).

$$C(s+1, n) = 2 \sum_{i=0}^{n-1} \binom{s}{i}, C(s, n) = 2^s \text{ for } s \leq n.$$

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

## 1.2 Hopfield Networks

**Hopfield Model**  $E(X) = -\frac{1}{2} \sum_{i \neq j} w_{ij} X_i X_j + \sum_i b_i X_i$ , where  $X_i \in \{\pm 1\}$ .  $w_{ij} = w_{ji}$ ,  $w_{ii} = 0$ .

### Hebbian Learning

Choose patterns  $\{x\}_{t=1}^s \in \{\pm 1\}^n$ , build weights once using them:  $w_{ij} = \frac{1}{n} \sum_{t=1}^s x_i^t x_j^t$ ,  $w_{ii} = 0$ . For inference, update  $X$  iteratively:  $X_i^{t+1} = \text{sign}\left(\sum_j w_{ij} X_j^t + b_i\right)$  asynchronously. Capacity for random, uncorrelated patterns:  $s_{\max} \approx 0.138n$ .

# 2 Feedforward Networks

## 2.1 Linear Models

### Linear regression (MSE)

$$L[w](X, y) = \frac{\|Xw - y\|^2}{2n}, \nabla L = \frac{X^\top Xw - X^\top y}{n}.$$

### Moore-Penrose inverse solution

$w^* = X^*y \in \arg\min_w L[w](X, y)$ , where  $X^* = \lim_{\delta \rightarrow 0} (X^\top X + \delta I)^{-1} X^\top$  Moore-Penrose inverse.

### Stochastic gradient descent update

$$w_{t+1} = w_t + \eta(y_{i_t} - w_t^\top x_{i_t})x_{i_t}, i_t \sim \mathcal{U}([1, n]).$$

### Gaussian noise model

$y_i = w^\top x_i + \varepsilon_i$ ,  $\varepsilon_i \sim N(0, \sigma^2)$ , LSQ equivalent to NLL of gaussian noise model.

### 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 for $y \in \{0, 1\}$

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

**Logistic regression with CE loss:**  $L[w] = \frac{1}{n} \sum_{i=1}^n \ell_i(y_i, w^\top x_i), \nabla \ell_i = [\sigma(w^\top x_i) - y_i] x_i$ .

## 2.2 Feedforward Networks

### Generic feedforward layers

$$F : \underbrace{\mathbb{R}^{m(n+1)}}_{\text{parameters}} \times \underbrace{\mathbb{R}^n}_{\text{input}} \rightarrow \underbrace{\mathbb{R}^m}_{\text{output}}, F[\theta](x) = \varphi(Wx + b).$$

**Composition of layers**  $G = F^L[\theta^L] \circ \dots \circ F^1[\theta^1]$ .

### 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**  $(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$ , CE-loss in terms of logits  $\ell(y, z) = -z_y + \log \sum_j e^{z_j}$ .

**Residual layer**  $F[W, b](x) = x + (\varphi(Wx + b) - \varphi(0))$ , therefore  $F[0, 0] = \text{id}$ . Link that propagates  $x$  forward is called a **skip connection**.

## 2.3 Sigmoid Networks

**Sigmoid activation**  $\sigma(z) = \frac{1}{1+e^{-z}} = 1 - \sigma(-z)$ .  $\sigma'(z) = \sigma(z)(1-\sigma(z))$ .

### Hyperbolic tangent activation

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

### Smooth function approximation

Polynomials, ridge functions ( $\varphi(a^\top x + b)$ ) and MLPs with  $C^\infty$  activations are universal approximators.

### Barron's Theorem: Approximation error

For  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  with  $C_f = \int \|\omega\| |\hat{f}(\omega)| d\omega < \infty$ ,  $\exists$  width- $m$  MLP  $g_m$  s.t.:  $\int_B |f - g_m|^2 dx \leq O(\frac{1}{m})$

## 2.4 ReLU( $z$ ) = max(0, $z$ ) networks

ReLU networks are universal approximators.

### Zalavsky's Thoerem: Activation patterns

$m$  ReLU neurons in  $\mathbb{R}^n$ . Each neuron's hyperplane  $\{w_i^\top x = 0\}$  partitions  $\mathbb{R}^n$  into  $R(m)$  connected regions of constant activation pattern.  $R(m) \leq \sum_{i=0}^{\min(n, m)} \binom{m}{i} \ll 2^m$ .

### Montufar: Connected regions in ReLU network

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

# 3 Gradient-Based Learning

## 3.1 Backpropagation

$$x^\ell = \varphi(W^\ell x^{\ell-1} + b^\ell), \frac{\partial \mathcal{L}}{\partial W^\ell} = \delta^\ell(x^{\ell-1})^\top, \frac{\partial \mathcal{L}}{\partial b^\ell} = \delta^\ell, \\ \delta^\ell = \frac{\partial \mathcal{L}}{\partial x^\ell} \odot \varphi'(W^\ell x^{\ell-1} + b^\ell), \frac{\partial \mathcal{L}}{\partial x^\ell} = (W^{\ell+1})^\top \delta^{\ell+1}.$$

## 3.2 Gradient Descent

Update:  $x_{t+1} = x_t - \eta \nabla f(x_t)$ .

**Gradient flow ODE**  $\frac{dx}{dt} = -\nabla f(x)$  gives ideal trajectory to be approximated by gradient descent.

**Newton's method** gives optimal step for quadratic model:  $\Delta x = -[\nabla^2 f(x)]^{-1} \nabla f(x)$ .

### Optimal LR for Convex Quadratics

For  $f(x) = \frac{1}{2} x^\top Qx$ ,  $\eta^* = \frac{2}{\lambda_{\max}(Q) + \lambda_{\min}(Q)}$ . Stability requires  $\eta \leq \frac{2}{\lambda_{\max}(Q)}$ . Quadratic approx. of  $f$ :  $f(x + \Delta x) \approx f(x) + \nabla f(x)^\top \Delta x + \frac{1}{2} \Delta x^\top \nabla^2 f(x) \Delta x$ .

**L-smoothness:**  $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ .

Equivalently (if  $f$  twice differentiable):

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2} \|y - x\|^2.$$

Implies  $\lambda_i(\nabla^2 f(x)) \leq L$  for all EVs  $\lambda_i$  of  $\nabla^2 f(x)$ .

**Convexity:**  $f(\lambda w + (1-\lambda)w') \leq \lambda f(w) + (1-\lambda)f(w')$

**$\mu$ -Strong convexity:** ( $\mu = 0 \Leftrightarrow$  convex + diff)

$$\Leftrightarrow f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|^2.$$

Implies  $\lambda_i(\nabla^2 f(x)) \geq \mu$  for all EVs  $\lambda_i$  of  $\nabla^2 f(x)$ .

For  $f : \mathbb{R} \rightarrow \mathbb{R}$ , these become:  $L \geq f''(x) \geq \mu \quad \forall x$ .

If  $f$  diff and  $L$ -smooth:  $f(x) - f(x^*) \geq \frac{1}{2L} \|\nabla f(x)\|^2$ .

### Polyak-Lojasiewicz condition

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \mu(f(x) - \min f) \text{ (forall } x\text{).}$$

$\mu$ -strong convex  $\Rightarrow \mu$ -PL.

### GD Convergence Rates & Learning Rates

**L-smooth only:**  $\eta^* = \frac{1}{L}$ . To reach  $\epsilon$ -stationary point ( $\|\nabla f\| \leq \epsilon$ ) needs at most  $\frac{2L}{\epsilon^2} (f(x_0) - \min f)$  steps.

**$\mu$ -PL + L-smooth:** Use  $\eta^* = \frac{2}{L+\mu}$ . Convergence:  $f(x_t) - f(x^*) \leq (1 - \frac{\mu}{L})^t (f(x_0) - f(x^*))$ .

## 3.3 Stochastic Gradient Descent

### SGD variance

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

### SGD convergence rate with Polyak averaging

$$\mathbb{E}[f(\bar{\theta}_t)] - \min f \leq O\left(\frac{1}{\sqrt{t}}\right) \text{ (general)}$$

$$\mathbb{E}[f(\bar{\theta}_t)] - \min f \leq O\left(\frac{\log t}{t}\right) \text{ (strongly convex)}$$

$$\mathbb{E}[f(\bar{\theta}_t)] - \min f \leq O\left(\frac{1}{t}\right) \text{ (additionally smooth)}$$

**Minibatch SGD:** Variance  $\downarrow$  by  $\propto r$ . Can  $\uparrow \eta \propto r$ .

### Var. Reduction with SVRG w/ occasional snapshot

$$\bar{\theta}: \theta_{t+1} = \theta_t - \eta [\nabla f_i(\theta_t) - \nabla f_i(\bar{\theta}) + \nabla f(\bar{\theta})].$$

# 3.4 Acceleration and Adaptivity

### Heavy ball momentum update

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

### Nesterov acceleration

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

More theoretical grounding than heavy ball.

### AdaGrad updates

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \odot \nabla f(\theta_t), \\ \gamma^t = \gamma^{t-1} + \nabla f(\theta_t) \nabla f(\theta_t)$$

## 4 Convolutional Networks

### 4.1 Convolutions

$$(f * g)(u) = (g * f)(u) = \int_{-\infty}^{\infty} f(t)g(u-t) dt = \int_{-\infty}^{\infty} f(u-t)g(t) dt. \\ (f * g) = \text{Toeplitz-Matrix}(g)f.$$

**Fourier transform convolution property**  
 $\mathcal{F}(f * g) = \mathcal{F}(f) \cdot \mathcal{F}(g)$ .

**Cross-correlation:**  $(g * f)[u] = \sum_i g[t]f[u+t]$ .

## 4.2 Convolutional Networks

**Conventions for Padding:** Add zeros around input.

**ConvNets for Images** ( $r$  out channel,  $u$  in channel)  
 $y[r][s, t] = \sum_u \sum_{i,j} w[r, u][i, j] * x[u][s+i, t+j]$ .

**Number of parameters of a convolutional layer**  
 $(|r| \times |u|) \cdot (|i| \times |j|)$ : fully connected  $\times$  patchsize.

## 4.3 Word2Vec

Per word  $\omega$ , have input embedding  $x_\omega$  and output embedding  $y_\omega$ .

Predict context word  $\nu$  given center word  $\omega$ :  
 $P(\nu | \omega) = \frac{\exp(x_\omega^\top y_\nu)}{\sum_\mu \exp(x_\omega^\top y_\mu)}$ .

NLL loss:  $\ell_{\omega, \nu} = -x_\omega^\top y_\nu + \ln \sum_\mu \exp(x_\omega^\top y_\mu)$ . Total:  $h(\{\omega\}, \{\nu\}) = \sum_{(\omega, \nu)} \ell_{\omega, \nu}$  over observed pairs. Use only input embeddings after training.

# 5 Geometric Deep Learning

**Group** is set  $G$  with a binary operation s.t.:  
1)  $(gh)f = g(hf)$ , 2)  $\exists e \in G$  s.t.  $ef = fe = f$ , 3)  
 $\forall g \exists g^{-1} \in G$  s.t.  $gg^{-1} = g^{-1}g = e$ , 4)  $gh \in G \forall g, h$ .

## 5.1 Sets and Points

**Order-invariance property:**

$$f(x_1, \dots, x_M) = f(x_{\pi(1)}, \dots, x_{\pi(M)}) \quad (\text{perturbations}).$$

**Equivariance property:**

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

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

**Deep Sets model** (invariant layer):

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

**Equivariant map construction:**

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

## 5.2 Graph Convolutional Networks

**Feature and adjacency matrices**

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

with  $a_{nm} = 1 \Leftrightarrow \{v_n, v_m\} \in E$ .

**Permutation matrix constraints**

$$P \in \{0, 1\}^{M \times M} \text{ with single 1 in each row and col.}$$

**Graph invariance definition**

$$f(X, A) \stackrel{!}{=} f(PX, PAP^\top), \forall P \in \Pi_M.$$

**Graph equivariance definition**

$$f(X, A) \stackrel{!}{=} Pf(PX, PAP^\top), \forall P \in \Pi_M.$$

**Node neighborhood features**

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

**Message passing scheme**

$$\varphi(x_m, X_m) = \varphi\left(x_m, \bigoplus_{x \in X_m} \psi(x)\right),$$

⊕ is some permutation-invariant operation.

**Normalized adjacency matrix**

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

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

**One GCN layer**

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

### 5.2.1 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**

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

**Graph Fourier transform**

$$L = D - A = U\Lambda U^\top, \hat{f} = U^\top f, f = U\hat{f}.$$

$$\Lambda := \text{diag}(\lambda_1, \dots, \lambda_M), \lambda_i \geq \lambda_{i+1}.$$

$$\text{Convolution: } x * y = U((U^\top x) \odot (U^\top y)).$$

**Filtering operation**

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

**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**

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

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

## 6 Theory of DNNs

### 6.1 Statistical Learning Theory

**Risk decomposition**

$$f^*: \text{optimal predictor over all functions,}$$

$$f_H^* = \underset{\text{excess risk}}{\text{argmin}}_{f \in H} \mathcal{R}(f), \hat{f}_H: \text{learned from finite data.}$$

$$\underbrace{\mathcal{R}(\hat{f}_H) - \mathcal{R}(f^*)}_{\text{estimation error}} = \underbrace{\mathcal{R}(\hat{f}_H) - \mathcal{R}(f_H^*)}_{\text{approximation error}} + \underbrace{\mathcal{R}(f_H^*) - \mathcal{R}(f^*)}_{\text{excess risk}}$$

**Rademacher Complexity**  $G = \{g_h \mid h \in H\}$ :

$$\text{For } \sigma \in \{-1, 1\}, \text{ measures how well } G \text{ can fit random noise: } \hat{\mathcal{R}}_{D_n}(G) = \mathbb{E}_\sigma \left[ \sup_{g \in G} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right].$$

$$\mathbb{E} \left[ \sup_{h \in H} \mathcal{R}(h) - \hat{\mathcal{R}}_{D_n}(h) \right] \leq 2\hat{\mathcal{R}}_{D_n}(G),$$

$$\mathbb{E} \left[ \mathcal{R}(\hat{h}_H) \right] \leq \mathcal{R}(h_H^*) + 2\hat{\mathcal{R}}_{D_n}(G).$$

**Double descent:**

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

**Implicit bias towards min norm solutions:**

Any convergent algorithm with iterates in  $\text{span}\{x_1, \dots, x_n\}$  finds the minimum norm solution.

#### 6.1.1 A PAC-Bayesian result

**PAC-Bayesian theorem**

Bounds generalization gap for stochastic classifiers ( $f \sim Q$ ):  $E_Q[\mathcal{R}(f)] - E_Q[\hat{\mathcal{R}}_n(f)] \leq \sqrt{\frac{2}{n} [\text{KL}(Q||P) + \ln(2\sqrt{n}/\varepsilon)]}$

- $P$ : prior,  $Q$ : posterior (learned). Rate  $\tilde{O}(1/\sqrt{n})$
- $\text{KL}(Q||P)$ : “information cost” of moving  $P \rightarrow Q$
- Insight: generalization depends on **distance moved**, not parameter count

**PAC-Bayesian for DNNs**

$$P = \mathcal{N}(0, \lambda I), Q = \mathcal{N}(\theta, \text{diag}(\sigma_i^2))$$

$$\text{KL}(Q||P) = \sum_i \left[ \log \frac{\lambda}{\sigma_i^2} + \frac{\sigma_i^2 + \theta_i^2 - 1}{2\lambda^2} - \frac{1}{2} \right]$$

$$\text{Minimize directly: } E_Q[\hat{\mathcal{R}}] + \sqrt{\frac{2}{n} [\text{KL}(Q||P) + \dots]}$$

⇒ encourages wide/flat minima (perturbations  $\theta + \varepsilon$  must also perform well)

**Implementation:** Reparameterization:  $\tilde{\theta} = \theta + \text{diag}(\sigma_i)\eta, \eta \sim \mathcal{N}(0, I)$ , Backprop through  $\theta$  and  $\sigma$ .

### 6.2 Linearized DNNs and NTK

Training neural network  $f(\theta)(x)$  can be approximated by **linearizing** around initialization  $\theta_0$  when parameters change slowly.

**Linearization → Kernel Regression:**

Taylor approximation:  $h(\beta)(x) = f(\theta_0)(x) + \beta \cdot \nabla f(\theta_0)(x), \beta = \theta - \theta_0$ .

With residuals  $\tilde{y}_i = y_i - f(\theta_0)(x_i)$ , training becomes **linear regression** with features  $\nabla f(\theta_0)(x_i)$ .

**Neural Tangent Kernel (NTK):**

Definition:  $k_\theta(x, x') := \nabla f(\theta)(x) \cdot \nabla f(\theta)(x')$ .

**Dual representation:**

$$h(\alpha)(x) = f(\theta_0)(x) + \sum_{i=1}^n \alpha_i k_{\theta_0}(x_i, x).$$

**Optimization problem:**  $\min_{\alpha} \frac{1}{2} \|K_{\theta_0}\alpha - \tilde{y}\|^2$

**Optimal solution (kernel regression):**

$$\alpha^* = K_{\theta_0}^\dagger(y - f(\theta_0)), h^*(x) = k_{\theta_0}(x)^\top \alpha^*$$

**Functional Gradient Flow**

Training dynamics in function space:

$$\dot{f}(\theta) = K(\theta)(y - f(\theta))$$

- If  $K(\theta)$  constant → **linear ODE** with closed-form solution
- If  $K(\theta)$  evolves → **nonlinear dynamics**, feature learning

**Infinite-Width Limit**

Initialization:  $w_{ij}^{(\ell)} \sim \frac{\sigma_w}{\sqrt{m_\ell}} \mathcal{N}(0, 1)$ . Result: As width  $m \rightarrow \infty$ :  $k_{\theta(t)} \rightarrow k_\infty$  (constant during training).

- Kernel becomes **deterministic** (depends only on architecture/init scheme)
- Training = kernel regression with frozen  $k_\infty$
- No feature learning

**Finite-width:**  $\|K(\theta_0) - K(\theta(t))\| = \mathcal{O}(\frac{1}{m})$

**Why Kernel Stays Constant:**

Kernel grad  $\nabla K = \nabla^2 f(x) \nabla f(z) + \nabla^2 f(z) \nabla f(x)$ , at  $m = \infty$ :  $\nabla^2 f \rightarrow 0 \Rightarrow \nabla K \rightarrow 0 \rightarrow$  kernel frozen.

Lazy Training (NTK Regime)	Feature Learning (Rich Regime)
$m \rightarrow \infty$ , small LR	Finite width, normal LR
$K$ const → linear dynamics	$K$ evolves → nonlinear
No feature learning	Learns representations
Theoretically tractable	SOTA performance

**Takeaways:**

**Linearization** turns NN training into kernel regression with features  $\nabla f(\theta_0)(x)$ .

**NTK**  $k_\theta = \nabla f(\theta)(x) \cdot \nabla f(\theta)(x')$  governs training dynamics via  $\dot{f} = K(y - f)$ .

**Infinite width** → kernel constant → NN = kernel machine (no feature learning).

**Finite width** → kernel evolves  $\mathcal{O}(\frac{1}{m})$  → enables feature learning.

**NTK explains lazy regime but NOT why deep learning works** → real power is feature learning when kernel changes.

## 6.3 Random NNs and GPs

**Marginals and Conditionals of MV Gaussians**

Let  $X \in \mathbb{R}^d \sim \mathcal{N}(\mu, \Sigma)$  with partition:

$$X = \begin{pmatrix} X_A \\ X_B \end{pmatrix}, \mu = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

**Marginal:**  $X_A \sim \mathcal{N}(\mu_A, \Sigma_{AA})$

**Conditional:**  $X_B \mid X_A \sim \mathcal{N}(\mu_{B|A}, \Sigma_{B|A})$  where

$$\mu_{B|A} = \mu_B + \Sigma_{BA}\Sigma_{AA}^{-1}(X_A - \mu_A)$$

$$\Sigma_{B|A} = \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB}$$

### 6.3.1 Bayesian Linear Regression

**Least-squares:**  $\hat{w} = \arg \min_w \frac{1}{2n\sigma^2} \|y - Xw\|^2$ .

**Closed-form solution:**  $\hat{w} = (X^\top X)^{-1} X^\top y$ .

**MLE interpretation:**  $y_i = x_i^\top w + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,

$$\mathcal{L}(w) = \log \prod_{i=1}^n p(y_i \mid x_i, w),$$

$$y_i \mid x_i, w \sim \mathcal{N}(x_i^\top w, \sigma^2).$$

**Prior:**  $p(w) = \mathcal{N}(\mathbf{0}, I_d)$ .

**Posterior:**  $p(w \mid y, X) = \frac{p(y \mid X, w)p(w)}{p(y \mid X)}$ .

**Predictive distribution:**  $p(y_{n+1}) = \int p(y_{n+1} \mid w)p(w \mid y) dw$ .

$$\Sigma_{ww} = I_d, \Sigma_{yw} = X, \Sigma_{yy} = XX^\top + \sigma^2 I_n.$$

$$\mu_{w \mid y} = X^\top \Sigma_{yy}^{-1} y = (X^\top X + \sigma^2 I_d)^{-1} X^\top y,$$

$$\Sigma_{w \mid y} = I_d - X^\top \Sigma_{yy}^{-1} X = \sigma^2 (X^\top X + \sigma^2 I_d)^{-1}.$$

Same result as ridge:  $\hat{w} = (X^\top X + \sigma^2 I_d)^{-1} X^\top y$ .

Equiv. to GP with linear kernel:  $f(\cdot) \sim \mathcal{GP}(0, k(\cdot, \cdot))$ :

$$k(x, x') = \phi(x)^\top \phi(x'), y = f + \varepsilon, f \sim \mathcal{N}(0, K),$$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n), \text{ posterior } p(f \mid y) = \mathcal{N}(\mu_{f \mid y}, \Sigma_{f \mid y}),$$

$$\mu_{f \mid y} = K(K + \sigma^2 I_n)^{-1} y.$$

### 6.3.2 NNGPs

**Setup:** Random 1-hidden-layer NN with  $m$  units:

$$f(x) = v_0 + \frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \varphi(\theta_j^\top x). \text{ Random init: } v_0 \sim \mathcal{N}(0, \sigma_0^2), \mathbb{E}[v_j^2] = \sigma_v^2, \text{Cov}(\theta_j) = \Sigma_\theta.$$

**Result:** As  $m \rightarrow \infty$ ,  $f(\cdot) \rightarrow \text{GP}(0, k)$  where  $k(x, x') = \sigma_0^2 + \sigma_v^2 \mathbb{E}_\theta[\varphi(\theta^\top x)\varphi(\theta^\top x')]$

**Monte Carlo approximation:** Sample  $B$  random NNs  $\{f_b\}_{b=1}^B$ , define features:

$$\cdot \varphi(x) = \frac{1}{\sqrt{B}} (f_1(x), \dots, f_B(x))^\top$$

$$\cdot \Phi = [f_1 \dots f_B] \in \mathbb{R}^{n \times B}$$
 (feature matrix)

$$\cdot \hat{K} = \Phi \Phi^\top$$
 (approximate kernel matrix)

**GP regression:** Posterior mean and variance:

$$\mathbb{E}[f(x) \mid y] = \varphi(x)^\top (\Phi^\top \Phi + \sigma^2 I_B)^{-1} \Phi^\top y$$

$$\text{Var}[f(x) \mid y] = \sigma^2 \varphi(x)^\top (\Phi^\top \Phi + \sigma^2 I_B)^{-1} \varphi(x)$$

**Key advantage:** Inverts  $B \times B$  matrix instead of  $n \times n$  when  $B \ll n$ .

# 7 Generative Models

## 7.0.1 Linear Autoencoders

**Setup.** Encoder  $C \in \mathbb{R}^{k \times d}$ , decoder  $D \in \mathbb{R}^{d \times k}$ , data  $X \in \mathbb{R}^{d \times n}$  (centered cols):  $\min_{C,D} \|X - DCX\|_F^2$ .

**Optimal Solution (PCA).** Let  $S = XX^\top$  with eigendecomposition  $S = Q\Lambda^2Q^\top$ ,  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$ . Optimal reconstruction via rank- $k$  projection:  $\hat{X} = U_k^*U_k^{\top}X$  where  $U_k^* = Q[:,1:k]$  are top- $k$  eigenvectors of  $S$  (equiv, top- $k$  left singular vectors of  $X$ ).

- Any  $C = U_k^{\top}A$ ,  $D = A^{-1}(U_k^*)^\top$  is optimal ( $\forall A$ )
- Reduces to truncated SVD:  $\hat{X} = U^*\Lambda_k V^\top$  with  $\Lambda_k = \text{Diag}(\lambda_1, \dots, \lambda_k, 0, \dots, 0)$
- Convex objective with no spurious local minima (gradient descent finds global optimum)
- Singular vectors may not be uniquely identified

## 7.0.2 Factor analysis

**Latent Variable Models** are a generic way to describe generative models. Latent variable  $z \sim p(z)$ , conditional models for observables  $x$ ,  $p(x|z)$ , observed data model:  $p(x) = \int p(x|z)p(z) dz$ .

**Mixture models:** simple discrete models:  $z \in [K]$ ,  $p(z)$  mixing proportions,  $p(x|z)$  condit. densities.

$$x \sim \mathcal{N}(\mu, \Sigma): p(x; \mu, \Sigma) = \frac{\exp[-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)]}{\sqrt{(2\pi)^n \det(\Sigma)}}$$

## 7.0.3 Linear Factor Analysis

- Latent prior:  $z \sim \mathcal{N}(\mathbf{0}, I)$ ,  $z \in \mathbb{R}^m$
- Observation:  $x = \mu + Wz + \eta$ ,  $\eta \sim \mathcal{N}(\mathbf{0}, \Sigma)$
- Independence:  $\eta \perp z$
- Typically  $m < n$  (fewer factors than features)

**Marginal distribution:**  $x \sim \mathcal{N}(\mu, WW^\top + \Sigma)$

- $WW^\top$ : shared variance (low-rank, explained by latent factors)
- $\Sigma$ : unique variance (diagonal, observation-specific)

Non-identifiability:  $(WQ)(WQ)^\top = WQQ^\top W^\top = WW^\top$  for any orthogonal  $Q$ . Factors only identifiable up to rotations/reflections.  $\Rightarrow$  Use factor rotations (varimax, etc.) for interpretability.

**MLE estimation:**  $\theta = (\mu, W) \xleftarrow{\text{max}} \log p(X; \mu, W)$

- $\hat{\mu} = \frac{1}{s} \sum_{i=1}^s x_i$  (closed form)
- No closed form for  $W \rightarrow$  use GD or EM algorithm

**Posterior (encoder):**  $p(z | x) = \frac{p(x | z)p(z)}{p(x)}$ .

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

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

**Probabilistic PCA:** Special case  $\Sigma = \sigma^2 I$ . Optimal  $i$ -th column:  $w_i = \rho_i u_i$ ,  $\rho_i^2 = \max\{0, \lambda_i - \sigma^2\}$  where  $(\lambda_i, u_i)$  is  $i$ -th eigenpair of data covariance.

As  $\sigma \rightarrow 0$ :  $\mu_{z|x} \rightarrow W^\top(x - \mu)$  (standard PCA).

## 7.1 Variational Autoencoders

$z \in \mathbb{R}^d$  is learned embedding of  $x$ . For generation,  $z \sim \mathcal{N}(\mathbf{0}, I)$ , decoder  $p_\theta(x|z)$  maps latent to data.

**Problem:**  $p_\theta(x) = \int p(z)p_\theta(x|z) dz$  intractable.

**Solution:** Maximize ELBO instead:  $\log p_\theta(x) \geq \underbrace{\mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)]}_{\text{reconstruction}} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \parallel p(z))}_{q_\phi \text{ regularization}}$

- **Reconstruction:** Encode  $x \rightarrow z$ , decode back ( $p_\theta$ )
- **KL term:** Keep encoder output close to prior  $p(z) \sim \mathcal{N}(\mathbf{0}, I)$ , ensures generation using latents

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

$$\text{KL closed form: } D_{\text{KL}}(q_\phi(z|x) \parallel p(z)) = \frac{1}{2} \sum_{j=1}^d \left( \frac{\sigma_j^2}{1} + \frac{(\mu_j - 0)^2}{1} - 1 - \log\left(\frac{\sigma_j^2}{1}\right) \right).$$

$$\log p_\theta(x|z) = -\frac{1}{2\sigma^2} \|x - \mu_\theta(z)\|_2^2 - \frac{d}{2} \log(2\pi\sigma^2).$$

**Reparameterization trick:**  $z = \mu_\phi(x) + \sigma_\phi(x) \odot \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, I)$ , enables backprop through sampling.

$$\log p_\theta(x) - \text{ELBO} = D_{\text{KL}}(q_\phi(z|x) \parallel p_\theta(z|x)), \text{ tight when } q_\phi = \text{true posterior.}$$

$$\text{Monte Carlo estimation: } E_{q_\phi(z|x)}[\log p_\theta(x|z)] \approx -\frac{1}{2\sigma^2 K} \sum_{k=1}^K \|x - \mu_\theta(z_k)\|_2^2 - \frac{d}{2} \log(2\pi\sigma^2).$$

## 7.2 Normalizing Flows

Transform simple distribution  $z \sim \mathcal{N}(\mathbf{0}, I)$  through invertible map  $T$  to get complex  $x = T(z)$ . Exact likelihood (no ELBO like VAEs), easy sampling.

**Change of Variables Formula:**

$$p_x(x) = p_z(T^{-1}(x)) \cdot |\det J_{T^{-1}}(x)|, \quad |\det J_{T^{-1}}(x)| = \frac{1}{|\det J_T(T^{-1}(x))|}.$$

**Diffeomorphism:**  $T$  is bijective, differentiable, with differentiable inverse. Guarantees  $\det J_T \neq 0$ .

**Computational problem:** Computing  $\det J$  is  $O(d^3)$  for dense Jacobian. **Solution:** Design  $T$  s.t. Jacobian is **triangular**, then only  $O(d)$ !

**Two architectures with triangular Jacobians:**

	MAF	IAF
Fast / parallel	Density eval	Sampling
Slow / sequential	Sampling	Density

**Coupling layers:** Trick that makes both directions fast, at the cost of being less expressive per layer.

## 7.3 Autoregressive Models

$$p(x) = \prod_{i=1}^d p(x_i | x_{<i}).$$

## 7.4 Generative Adversarial Networks

Likelihood-free generative model: train via adversarial game between two networks: **Generator**  $G_\theta$  maps latent  $z \sim p_z$  (typically Gaussian) to fake samples; **Discriminator**  $D_\varphi$ : outputs prob that input is **real**.

$$\text{GAN Objective: } \min_\theta \max_\varphi \underbrace{\mathbb{E}_{x \sim p_r} [\log D_\varphi(x)]}_{\text{fake samples}} + \underbrace{\mathbb{E}_{z \sim p_z} [\log(1 - D_\varphi(G_\theta(z)))]}_{\text{real samples}}$$

- **Discriminator** maximizes: correctly classify real (high  $D$ ) and fake (low  $D$ )
- **Generator** minimizes: fool discriminator (make  $D(G(z))$  high)

### 7.4.1 Theoretical Foundation

Binary classification with  $p(y=1) = p(y=0) = \frac{1}{2}$ :

- $y=1$ : sample from real  $p_{r(x)}$
- $y=0$ : sample from generator  $p_{\theta(x)}$

**Bayes Optimal Classifier** (prob that  $x$  is real):

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

**Generator Logistic Objective = JS Divergence:**

$$\ell^*(\theta) = \mathbb{E}_{\tilde{p}_\theta(x,y)}[y \ln q_\theta(x) + (1-y) \ln(1 - q_\theta(x))] = \text{JS}(p_r \parallel p_\theta) - \ln 2.$$

$$\text{Jensen-Shannon Divergence: } \text{JS}(p_r \parallel p_\theta) = \frac{1}{2} D_{\text{KL}}(p_r \parallel p_m) + \frac{1}{2} D_{\text{KL}}(p_\theta \parallel p_m), \quad p_m = \frac{p_r + p_\theta}{2}.$$

$$\text{Bounded: } 0 \leq \text{JS}(p_r \parallel p_\theta) \leq \log 2.$$

### 7.4.2 Training

Alternating SGD (heuristic, may diverge!). Training is **Saddle-point problem**, notoriously unstable!

**JS Divergence Saturates** when distributions don't overlap. If  $p_r$  and  $p_\theta$  have disjoint supports: discriminator perfect, no gradient for generator!

**Wasserstein Distance (Earth Mover's Distance):**

$$W(p_r, p_\theta) = \inf_{\gamma \in \Pi(p_r, p_\theta)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

Minimum total "work" to transport mass from  $p_r$  to  $p_\theta$ .

Provides meaningful gradients even without overlap.

**Dual (Kantorovich-Rubinstein):**

$$W(p_r, p_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p_r}[f(x)] - \mathbb{E}_{x \sim p_\theta}[f(x)].$$

Maximize gap between avg score of real vs fake samples w.r.t. Lipschitz constraint.

Max achievable gap = Wasserstein distance.

**WGAN** uses critic  $f_w$  (not classical discriminator!):  $\min_\theta \max_w \mathbb{E}_{x \sim p_r}[f_w(x)] - \mathbb{E}_{z \sim p_z}[f_w(G_\theta(z))]$ .

**Enforcing Lipschitz:**

- **Weight clipping** (original): crude, problematic
- **Gradient penalty**: add  $\lambda \mathbb{E}_{\hat{x}}[(\|\nabla_{\hat{x}} f_w(\hat{x})\|_2 - 1)^2]$

**Mode Collapse:** Generator produces only few samples that fool discriminator, ignoring full distribution diversity.

## 7.5 Diffusion Models

**Forward process (fixed):** Gradually add Gaussian noise over  $T$  steps until data becomes pure noise.

Fwd step:  $q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$ .  
Full fwd process:  $q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$ .

**Noise schedule:**  $\{\beta_t \in (0, 1)\}_{t=1}^T$  noise added at each step. **Define:**  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ .

**Direct sampling (reparameterization trick):**

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I). \\ x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I).$$

**Reverse process (learned):** Train NN to denoise step by step:  $p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I)$ . For small  $\beta_t$ , the reverse  $q(x_{t-1} | x_t)$  is also Gaussian.

$\log p_\theta(x_0)$  intractable, so derive **Variational Lower Bound (VLB)**:  $-\log p_\theta(x_0) \leq \mathcal{L}_{\text{VLB}} = \mathbb{E}_q[\log \frac{q(x_{1:T} | x_0)}{p_\theta(x_0, T)}]$ . **Decomposition into 3 terms:**

$$\mathcal{L}_{\text{VLB}} = \underbrace{\mathbb{E}_{x \sim p(x)}[D_{\text{KL}}(q(x_T | x_0) \parallel p(x_T))]}_{L_T} + \underbrace{\sum_{t=2}^T \mathbb{E}_{q(x_t | x_0)}[D_{\text{KL}}(q(x_{t-1} | x_t, x_0) \parallel p_\theta(x_{t-1} | x_t))]}_{L_{t-1}} - \underbrace{\mathbb{E}_{q(x_1 | x_0)}[\log p_\theta(x_0 | x_1)]}_{L_0}$$

•  $L_T$ : Is  $q(x_T | x_0) \approx \mathcal{N}(0, I)$ ? Not optimized.

•  $L_{t-1}$ : Match learned reverse to true reverse

•  $L_0$ : Reconstruction term

**Tractable Reverse Posterior**  $q(x_{t-1} | x_t, x_0)$  is Gaussian with closed form (product of Gaussians):  $q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_{q,t}(x_t, x_0), \sigma_t^2 I)$ , with:

$$\mu_{q,t}(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \varepsilon_0 \right), \quad \sigma_t^2 = \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \beta_t.$$

## 7.6 Noise Prediction Parameterization

Predict the **noise**  $\varepsilon_\theta(x_t, t)$  instead of mean directly. Parameterize learned mean to mirror true posterior:  $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \varepsilon_\theta(x_t, t) \right)$ .

Since both distributions are Gaussian with same variance:  $D_{\text{KL}}(\mathcal{N}(\mu_q, \sigma^2 I) \parallel \mathcal{N}(\mu_p, \sigma^2 I)) = \frac{1}{2\sigma^2} \|\mu_q - \mu_p\|^2$ . This simplifies  $L_{t-1}$  to comparing noise:  $L_{t-1} = \mathbb{E}_{x_0, \varepsilon_0} \left[ \frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \alpha_t)\sigma_t^2} \varepsilon_0 - \varepsilon_\theta(x_t, t) \right]^2$ .

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t \sim [1, T], x_0, \varepsilon_0} [\|\varepsilon_0 - \varepsilon_\theta(x_t, t)\|^2].$$

Training	Sampling
<ol style="list-style-type: none"> <li>1. Sample real image <math>x_0 \sim q(x_0)</math></li> <li>2. Sample random timestep <math>t \sim \text{Uniform}\{1, \dots, T\}</math></li> <li>3. Sample noise <math>\varepsilon \sim \mathcal{N}(0, I)</math></li> <li>4. Compute noisy image: <math>x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon</math></li> <li>5. Grad step on <math>\nabla_\theta \ \varepsilon - \varepsilon_\theta(x_t, t)\ ^2</math></li> </ol>	<ol style="list-style-type: none"> <li>1. Sample <math>x_T \sim \mathcal{N}(0, I)</math></li> <li>2. For <math>t = T, \dots, 1</math>: <math>z \sim \mathcal{N}(0, I)</math> if <math>t &gt; 1</math>, else <math>z = 0</math></li> <li>3. Compute <math>\varepsilon_\theta</math>: <math>x_t = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}}\varepsilon_\theta(x_t, t)) + \sigma_t z</math></li> <li>4. Return <math>x_0</math></li> </ol>

Cosine noise schedule performs better than linear.

Used architecture is U-Net. **Input:** Noisy image  $x_t + \text{timestep } t$ ; **Output:** Predicted noise  $\varepsilon_\theta(x_t, t)$ .

Model **conditional distribution**  $p_\theta(x_{0:T} | y)$  where  $y$  is condition (class, text, image). Extend denoiser to take  $y$  as input.

**Latent Diffusion Models (LDM)** run diffusion in **compressed latent space** instead of pixel space.

## 8 Tricks

**Short Connections** in DNs: Add less deep paths to a very deep network. **Residual** connections: shortcut and add back in. **Skip** connections: concatenate.

### 8.1 Weight Decay & Early Stopping

#### L2 regularization

$\mathcal{R}_\Omega(\theta; \mathcal{S}) = \mathcal{R}(\theta; \mathcal{S}) + \Omega(\theta)$ ,  $\Omega_\mu(\theta) = \frac{\mu}{2} \|\theta\|^2$ ,  $\mu \geq 0$ . Only penalize weights, not biases. **GD upd w/ WD**:  $\Delta\theta = -\eta \nabla \mathcal{R}(\theta) - \eta \nabla \Omega_\mu(\theta) = -\eta \nabla \mathcal{R}(\theta) - \eta \mu \theta$ .

Geometric interpretation (local quadratic approx): Regularized optimum:  $\theta_\mu^* = (H + \mu I)^{-1} H \theta^*$ , where  $H = Q^\top \Lambda Q$  gives  $\theta_\mu^* = Q \operatorname{diag}\left(\frac{\lambda_i}{\lambda_i + \mu}\right) Q^\top \theta^*$ .

$\lambda_i \gg \mu: \frac{\lambda_i}{\lambda_i + \mu} \approx 1 \rightarrow$  weak shrinkage (important dirs).  $\lambda_i \ll \mu: \frac{\lambda_i}{\lambda_i + \mu} \approx 0 \rightarrow$  strong shrinkage (flat dirs).

Adaptively shrinks based on loss geometry, preserves important dirs, removes unnecessary complexity.

**Early stopping**: Rather than training to convergence, stop when validation performance plateaus. Analysis shows that this is approximately equivalent to L2 regularization. GD trajectories can be approximated as  $\theta(k) = [I - (I - \eta \Lambda)^k] \theta^*$ . For small step sizes, behaves like weight decay when  $k = \frac{1}{\eta \mu}$ .

## 8.2 Ensemble Methods

**Bagging**: Create  $K$  bootstrap samples of Data (sampling with replacement), train separate models, and average predictions:  $p(y|x) = \frac{1}{K} \sum_{k=1}^K p(y|x; \theta_k)$ .

**Dropout**: Randomly drop units during training with probability  $1 - \pi$ . Creates an exponential ensemble of sub-networks sharing weights. Test time: Scale weights by  $\pi$  to approximate the ensemble average.

## 8.3 Normalization

**Batch Norm**: Normalize activations across mini-batch:  $\tilde{z} = \frac{z - \mu_{\text{batch}}}{\sigma_{\text{batch}}}$ ,  $\hat{z} = \alpha \tilde{z} + \beta$ ,  $\mu_{\text{batch}} = \frac{1}{b} \sum_{i=1}^b z_i$ ,  $\sigma_{\text{batch}} = \sqrt{\frac{1}{b} \sum_{i=1}^b (z_i - \mu_{\text{batch}})^2}$ . **Layer Norm**: Normalize features in a layer instead; particularly effective for RNNs (batch statistics are less stable).

## 8.4 Data / Task Augmentation

Augment Data by applying valid transformations. Semi-supervised Learning: Train jointly on labeled and unlabeled data w/ combined loss. Pre-training & Fine-tuning. Multi-task Learning. Self-supervised Learning: Create free supervision from data.

## 9 Recurrent Neural Networks

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

Optional output:  $y_t = G[\theta](z_t)$ .

**Simple RNN**:  $z_t = \phi(Wz_{t-1} + Ux_t)$  where  $W \in \mathbb{R}^{m \times m}$ ,  $U \in \mathbb{R}^{m \times n}$

**Backpropagation Through Time** (param sharing):

$$\frac{\partial R}{\partial w_{ij}} = \sum_t \frac{\partial R}{\partial z_i^t} \cdot \dot{\phi}_i^t \cdot z_j^{t-1}.$$

$$\frac{\partial R}{\partial u_{ik}} = \sum_t \frac{\partial R}{\partial z_i^t} \cdot \dot{\phi}_i^t \cdot x_k^t; \dot{\phi}_i^t = \phi'(F_i(z^{t-1}, x^t)).$$

**Gradient flow backward through time**:

$$\nabla_{x_t} \mathcal{R} = \left[ \prod_{r=t+1}^s W^\top S(z^r) \right] \cdot J_G \cdot \nabla_y \mathcal{R}$$

**Spectral analysis**:  $\left\| \prod W^\top S(z^r) \right\|_2 \leq [\sigma_{\max}(W)]^{s-t}$

**Root cause**: Repeated matmul through time.

$\Rightarrow$  Simple RNNs cannot learn long dependencies.

**Deep RNNs** stack layers vertically:

$$z^{t,\ell} = \varphi(W_\ell z^{t-1,\ell} + U_\ell z^{t,\ell-1}) \text{ where } z^{t,0} = x_t.$$

### 9.1 Long Short-Term Memory (LSTM)

- $C_t$ : cell state (internal memory, protected highway)
- $h_t$ : hidden state (external output, filtered view)

$$C_t = \sigma(F \tilde{x}^t) \odot C_{t-1} + \sigma(G \tilde{x}^t) \odot \tanh(V \tilde{x}^t), \\ z_t = \sigma(H \tilde{x}^t) \odot \tanh(C_t), \text{ where } \tilde{x}^t = [x_t, z_{t-1}].$$

### 9.2 Gated Recurrent Unit (GRU)

**Single state  $z_t$  Input**:  $\tilde{x}^t = [x_t, z_{t-1}]$ .

$$u_t = \sigma(U \tilde{x}^t), \quad r_t = \sigma(R \tilde{x}^t), \\ z_t = u_t \odot z_{t-1} + (1 - u_t) \odot \tanh(W[r_t \odot z_{t-1}, x_t])$$

Often comparable to LSTM with fewer resources. Gating creates identity paths  $\rightarrow$  better gradient flow.

## 9.3 Linear Recurrent Models

RNNs not parallelizable during training. LRU has linear dynamics:  $z_{t+1} = Az_t + Bx_t$ . Diagonalize to  $A = P\Lambda P^{-1}$ ,  $\lambda_i \in \mathbb{C}$ , change basis  $\zeta_t = P^{-1}z_t$ . Then:  $\zeta_{t+1} = \Lambda\zeta_t + Cx_t$ . Each dimension evolves independently (no channel mixing). Compensate with expressive output:  $y_t = \text{MLP}(\text{Re}(Gz_t))$ .

**Stability**: Require  $\max |\lambda_j| \leq 1$  (spectral radius  $\leq 1$ ).

**Parameterization**:  $\lambda_i = \exp(-\exp(\nu_i) + i\varphi_i)$  ensures  $|\lambda_i| \in (0, 1)$  automatically,  $|\lambda_i| \approx 1$ : Long-term memory,  $|\lambda_i| \approx 0$ : Short-term patterns.

**Provably universal** as sequence-to-sequence map.

### 9.3.1 Connectionist Temporal Classification

**Problem**: Unsegmented sequences (e.g., speech).

**Solution**: RNN outputs prob distribution over vocabulary at each time step. Model all alignments with blank symbol “-”:  $p(\ell | x) = \sum_{\pi \in \mathcal{B}^{-1}(\ell)} \prod_t y_{\pi_t}$ .  $\mathcal{B}$  removes blanks and repeated symbols.

## 9.4 Sequence Learning

**Teacher Forcing**:  $p(y^t)$  depends on  $y^{1:t-1}$  only through  $z^t$ , means during autoregressive generation, model doesn't see its own predictions. **Solution**: Add feedback connections from  $y^{t-1}$  to  $z^t$ :  $z^t = \text{RNN}(z^{t-1}, x^t, y^{t-1})$ , now model conditions on its own previous predictions  $\rightarrow$  more coherent gen.

**Professor Forcing**: Train two networks (teacher-forced + free-running), discriminator matches hidden states  $\rightarrow$  improved generalization.

**Seq2Seq**: Input and output sequences have different lengths: Use encoder-decoder framework.

sentence B follows A). Has bidirectional encoder, and task-specific heads. Can do FT for various tasks.

**Vision Transformers** split images into patches, add pos embeddings and a [CLS] token, then process with a standard transformer encoder.

## 11 Ethics

**Adversarial examples** (given  $f(x) = y$  correctly):

- Untargeted:  $\|\delta\| \leq \varepsilon$  s.t.  $f(x + \delta) \neq y$ .
  - Optimize  $\max_{\|\delta\| \leq \varepsilon} L(f(x + \delta), y)$ .
- Targeted:  $\|\delta\| \leq \varepsilon$  s.t.  $f(x + \delta) = t \neq y$ .
  - Optimize  $\min_{\|\delta\| \leq \varepsilon} L(f(x + \delta), t)$ .

**Linear Binary** ( $y \in \{-1, 1\}$ ,  $f(x) = w^\top x + b$ ):

**Correct**:  $y(w^\top x + b) > 0$ .

Adv. flips when  $yw^\top \delta \leq -y(w^\top x + b)$  cross hyperplane. **L2 optimal**:  $\delta^* = \frac{-w^\top x + b}{\|w\|_2^2} w$ ,  $\|\delta^*\|_2 = \frac{\|w^\top x + b\|}{\|w\|_2}$ .

$L_\infty$  optimal:  $\delta = -\varepsilon \operatorname{sign}(yw)$ .

**Multiclass**:  $f_k(x) = w_k^\top x + b_k$ , use  $\operatorname{argmax}_k f_k(x)$ .

**Margin** to class  $j$ :  $m_j(x) = (w_j - w_j)^\top x + (b_j - b_j)$ .

$f_j(x) = f_j(x) \Leftrightarrow m_j(x) = 0$ . **Correct** if  $f_j(x) > f_j(x) \forall j \neq y$ , **adversarial** if  $\exists j \neq y$  s.t.  $f_j(x + \delta) < f_j(x + \delta)$ . Distance to boundary:  $\frac{m_j(x)}{\|w_j - w_j\|_2}$ .

**Adversarial attacks for NNs**: Approximate boundary by  $f(x + \delta) \approx f(x) + \nabla f(x)^\top \delta$ . **FGSM** is a one-step  $L_\infty$  attack:  $\delta = \varepsilon \operatorname{sign}(\nabla_x L(f(x), y))$ . **PGD** is multi-step  $\delta_{\{t+1\}} = \operatorname{Proj}_{\|\delta\| \leq \varepsilon} (\delta_t + \alpha \operatorname{sign}(g_t))$ .

**Distributionally Robust Optimization**:

$\min_f \sup_{Q \in U(P)} E_Q[L(f(x))]$ , where  $U$  means close. Can use upper bound on Wasserstein distance e.g.

**Robust training**  $\min_f \mathbb{E}[\max_{\delta \in S} L(f(x + \delta), y)]$ .

Adversarial training can be viewed as robustness to distribution shift measured by Wasserstein distance.

**Interpretability**: Local - explain pred for specific  $x$ , Global - explain model behaviour on avg over data.

**Local**: Ceteris paribus (vary  $x_j$ , fix  $x_{-j}$ ), Sensitivity ( $\partial_{x_j} f(x)$ ), missing info ( $f(x) - \mathbb{E}[f(X) | X_{-j} = x_{-j}]$ ). **Global**: Mutual info ( $I(X_j; Y | X_{-j})$ ), Predictive util (train  $f$  w/ and w/o  $x_j$ ). For log-loss predictive util  $\approx$  conditional mutual information.

SHAP attributes predictions, while SAGE attributes risk reduction.

A protected attribute,  $Y$  target outcome,  $\hat{Y}$  prediction. Demographic Parity:  $\hat{Y} \perp A$ ; Equalized Odds:  $\hat{Y} \perp A | Y$ , Equality of Opportunity:  $\hat{Y} \perp A | Y = 1$ .