

- Closed form of KL for Gaussians:

$$D_{KL} = \frac{1}{2} \left[\log \frac{\det(\Sigma_2)}{\det(\Sigma_1)} - d + \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^\top \Sigma_2^{-1} (\mu_2 - \mu_1) \right].$$

- Cosine theorem:

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle.$$

- Points of \mathcal{S} are in general position if any subset $\Xi \subseteq \mathcal{S}$ with $|\Xi| \leq d$ is linearly independent.

Connectionism

McCulloch-Pitts neuron

$$f[\sigma, \theta](\mathbf{x}) \doteq \mathbb{1}\{\sigma^\top \mathbf{x} \geq \theta\}, \quad \sigma \in \{-1, 1\}^n, \mathbf{x} \in \{0, 1\}^n, \theta \in \mathbb{R}.$$

Perceptron

$$f[\mathbf{w}, b](\mathbf{x}) \doteq \text{sgn}(\mathbf{w}^\top \mathbf{x} + b)$$

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

$$\gamma[\mathbf{w}, b](S) \doteq \min_{(\mathbf{x}, \mathbf{y}) \in S} \gamma[\mathbf{w}, b](\mathbf{x}, \mathbf{y})$$

$$\mathcal{V}(S) \doteq \{(\mathbf{w}, b) \mid \gamma[\mathbf{w}, b](S) > 0\}.$$

Decision boundary hyperplane: $\mathbf{w}^\top \mathbf{x} / \|\mathbf{w}\| + b / \|\mathbf{w}\| \stackrel{!}{=} 0$. If $\gamma[\mathbf{w}, b](S) > 0$, then \mathcal{S} is linearly separated. $\mathcal{V}(S) \neq \emptyset$ iff \mathcal{S} is linearly separable. Adding points to \mathcal{S} makes $\mathcal{V}(S)$ smaller.

The perceptron algorithm tries to find any $(\mathbf{w}, b) \in \mathcal{V}(S)$. It does not aim to find solution with smaller error if $\mathcal{V}(S) = \emptyset$. Iterative mistake-driven algorithm:

$$f[\mathbf{w}, b](\mathbf{x}) \neq y \implies \mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}, \quad b \leftarrow b + y.$$

For all iterates $\mathbf{w}_t \in \text{span}(\mathbf{x}_1, \dots, \mathbf{x}_n)$.

Convergence can be proven by using $\mathbf{w}^\top \mathbf{w}_t \geq t\gamma$ and $\|\mathbf{w}_t\| \leq R\sqrt{t}$, where $\|\mathbf{w}\| = 1$, $\gamma[\mathbf{w}](S) > 0$, and $R = \max_{\mathbf{x} \in S} \|\mathbf{x}\|$. Then, bound $1 \geq \cos \angle(\mathbf{w}, \mathbf{w}_t) = \mathbf{w}^\top \mathbf{w}_t / \|\mathbf{w}_t\|$. We convergence within $\lceil R^2 / \gamma^2 \rceil$ iterations.

Number of unique classifications:

$$\mathcal{C}(S, d) = \left| \left\{ \mathbf{y} \in \{-1, 1\}^n \mid \exists \mathbf{w} : \forall i : y_i(\mathbf{w}^\top \mathbf{x}_i) > 0 \right\} \right|$$

Assume that points are in general position. Cover's theorem:

$$\mathcal{C}(n+1, d) = 2 \sum_{i=0}^{d-1} \binom{n}{i}.$$

Proof: Base cases are both 2 and adding a point has two cases \rightarrow Recurrence:

$$\mathcal{C}(n+1, d) = \mathcal{C}(n, d) + \mathcal{C}(n, d-1).$$

For $n \leq d$, we have $\mathcal{C}(n, d) = 2^n$. After $n = 2d$, there is a steep decrease in number of linear classification, quickly moving toward 0.

Parallel distributed processing

(1) A set of processing units with states of activation; (2) Output functions for each unit; (3) A pattern of connectivity between units; (4) Propagation rules for propagating patterns of activity; (5) Activation functions for units; (6) A learning rule to modify connectivity based on experience; (7) An environment within which the system must operate.

Hopfield networks

Models an associative memory, which aims to reconstruct a memory from an input that has been subjected to noise. Energy function via second-order interactions between n binary neurons:

$$H(\mathbf{x}) \doteq \frac{1}{2} \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j + \sum_{i=1}^n b_i x_i, \quad \mathbf{x} \in \{-1, 1\}^d.$$

We have $w_{ii} = 0$ and $w_{ij} = w_{ji}$. Simple dynamics:

$$x_i \leftarrow \begin{cases} +1 & H([, x_{i-1}, +1, x_{i+1},]) \leq H([, x_{i-1}, -1, x_{i+1},]) \\ -1 & \text{otherwise.} \end{cases}$$

Or: $x_i \leftarrow \text{sgn}(H_i)$, where $H_i \doteq \sum_{j=1}^d w_{ij} x_j - b_i$.

Hebbian learning (neurons frequently in the same state reinforce):

$$\mathbf{W} = \frac{1}{d} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top.$$

Pattern \mathbf{x}_i is memorized if meta-stable: update rule makes no updates to it:

$$x_{ii} = \text{sgn}(x_{ii} + C_{ii}), \quad C_{ii} \doteq \frac{1}{d} \sum_{j=1}^d \sum_{r \neq i} x_{ri} x_{rj} x_{ij}.$$

If cross-talk $|C_{ii}| < 1$ for all i , then \mathbf{x}_i is meta-stable.

Feedforward networks

Regression models

Mean-squared error

$$\ell[\theta](S) = \frac{1}{2} \sum_{i=1}^n (f[\theta](\mathbf{x}_i) - y_i)^2.$$

Linear model: $\ell[\mathbf{w}](S) = \frac{1}{2} \|\mathbf{X}^\top \mathbf{w} - \mathbf{y}\|^2 \rightarrow$ Closed form

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

Logistic regression (binary outputs) use sigmoid:

$\sigma(z) \doteq 1/(1 + \exp(-z))$. Binary cross-entropy loss:

$$\ell[\theta](S) = -\frac{1}{n} \sum_{i=1}^n y_i \log \sigma(f[\theta](\mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(f[\theta](\mathbf{x}_i))).$$

(Multi-class) Cross-entropy loss:

$$\ell[\theta](\mathbf{x}, y) = -\log \text{softmax}(f[\theta](\mathbf{x}))_y.$$

Layers and units

Mapping: $f[\mathbf{W}, \mathbf{b}](\mathbf{x}) = \phi(\mathbf{W}\mathbf{x} + \mathbf{b})$, where ϕ is a pointwise activation function. DNNs compose: $G = F_L[\theta_L] \circ \dots \circ F_1[\theta_1]$.

Intermediate layers are permutation symmetric: $\mathbf{F}[\mathbf{W}, \mathbf{b}](\mathbf{x}) = \mathbf{P}^{-1} \phi(\mathbf{P}\mathbf{W}\mathbf{x} + \mathbf{P}\mathbf{b}) = \mathbf{P}^{-1} \mathbf{F}[\mathbf{P}\mathbf{W}, \mathbf{P}\mathbf{b}](\mathbf{x}) \rightarrow$ Parameters are not unique.

Linear networks

Linear layers are closed under composition \rightarrow We do not gain representational power from composing them.

Residual networks

Residual layers:

$$\mathbf{F}[\mathbf{W}, \mathbf{b}](\mathbf{x}) = \mathbf{x} + \phi(\mathbf{W}\mathbf{x} + \mathbf{b}) - \phi(\mathbf{0}).$$

Then, $\mathbf{F}[\mathbf{0}, \mathbf{0}] = \text{Id}$. This makes it such that the model learns how to incrementally learn a better representation, rather than having to learn it at every layer. If input and output have different dimensionality, linearly project \mathbf{x} . Using this architecture, the depth can be increased significantly, because gradients propagate better.

Sigmoid networks

MLP with sigmoid activation:

$$g[\mathbf{v}, \mathbf{W}, \mathbf{b}](\mathbf{x}) \doteq \sigma^\top \mathbf{v}(\mathbf{W}\mathbf{x} + \mathbf{b}), \quad \mathbf{v}, \mathbf{b} \in \mathbb{R}^m, \mathbf{W} \in \mathbb{R}^{m \times n}.$$

Function class: $\mathcal{G}_n \doteq \bigcup_{m=1}^n \mathcal{G}_{n,m}$, where

$$\mathcal{G}_{n,m} \doteq \{g[\mathbf{v}, \mathbf{W}, \mathbf{b}] \mid \mathbf{b} \in \mathbb{R}^m, \mathbf{W} \in \mathbb{R}^{m \times n}\}.$$

Or, as a linear span of units,

$$\mathcal{G}_n = \text{span}\{\sigma(\mathbf{w}^\top \mathbf{x} + b) \mid \mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}\}.$$

This set universally approximates $\mathcal{C}(\mathbb{R}^n)$. But, it does not provide insight into how depth affects performance \rightarrow Barron: Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with finite \mathcal{C}_f and any $r > 0$, there is a sequence of one hidden layer MLPs $(g_m)_{m \in \mathbb{N}}$ such that

$$\int_{r\mathbb{B}} (f(\mathbf{x}) - g_m(\mathbf{x}))^2 \mu(d\mathbf{x}) \leq \mathcal{O}(1/m),$$

where $r\mathbb{B}$ is a r -radius ball and μ is any probability measure. Relaxing the notion of approximation to squared error over ball with radius r gives a decay of $1/m$ on the approximation error.

ReLU networks

ReLU: $(z)_+ \doteq \max_{0,z}$. Consider a layer of m ReLU units, then each unit is active or inactive: $\mathbb{1}\{\mathbf{W}\mathbf{x} + \mathbf{b} > 0\} \in \{0, 1\}^m$. As such, we can partition the input space into cells that have the same activation pattern:

$$\mathcal{X}_\kappa \doteq \{\mathbf{x} \mid \mathbb{1}\{\mathbf{W}\mathbf{x} + \mathbf{b} > 0\} = \kappa\}.$$

The number of cells is a proxy for the complexity of a network. Consider a ReLU network with L layers of $m > n$ width. The number of linear regions is lower bounded by

$$R(m, L) \geq R(m) \lfloor m/n \rfloor^{n(L-1)}.$$

Thus, complexity is related to depth.

Piecewise linear functions are approximators of $\mathcal{C}(\mathbb{R})$ and a piecewise linear function g with m pieces can be written as

$$g(x) = ax + b + \sum_{i=1}^{m-1} c_i(x - x_i)_+.$$

Using the lifting lemma, ReLU networks are universal approximators of $\mathcal{C}(\mathbb{R})$.

Gradient-based learning

Backpropagation

Backpropagation computes gradient in linear time if we know the gradient of all basic blocks in the function. Assume we have the following function,

$$\mathbf{F}[\theta](\mathbf{x}) \doteq (F_L \circ \dots \circ F_1)(\mathbf{x}), \quad \mathbf{h}_\ell \doteq F_\ell[\theta_\ell](\mathbf{h}_{\ell-1}), \quad \mathbf{h}_0 = \mathbf{x}.$$

We need the following,

$$\delta_\ell = \frac{\partial h(\theta)}{\partial \mathbf{h}_\ell}.$$

We have a recurrence,

$$\delta_\ell = \left[\frac{\partial \mathbf{h}_{\ell+1}}{\partial \mathbf{h}_\ell} \right]^\top \delta_{\ell+1}, \quad \delta_\ell = \frac{\partial h(\theta)}{\partial \mathbf{h}_\ell}.$$

Then, to compute parameter gradients,

$$\frac{\partial h(\theta)}{\partial \theta_\ell} = \delta_\ell \frac{\partial \mathbf{h}_\ell}{\partial \theta_\ell}.$$

Gradient descent

Update rule:

$$\theta^{t+1} = \theta^t - \eta \nabla_\theta h(\theta^t), \quad \eta > 0, \quad h(\theta) \doteq \ell \circ F[\theta].$$

Discretization of ODE:

$$d\theta = -\nabla h(\theta) dt.$$

Trajectory outcome (point where $\nabla h(\theta) = \mathbf{0}$) depends on initial conditions.

Gradient descent can only be successful if gradients change slowly \rightarrow Smoothness: h is L -smooth if

$$\|\nabla h(\theta) - \nabla h(\theta')\| \leq L \|\theta - \theta'\|, \quad \forall \theta, \theta'.$$

Or: $\|\nabla^2 h(\theta)\|_2 \leq L$ for all θ . If $\eta = 1/L$, then we have sufficient decrease,

$$h(\theta') \leq h(\theta) - \frac{1}{2L} \|\nabla h(\theta)\|^2, \quad \forall \theta, \theta'.$$

Let h be L -smooth and $\eta = 1/L$, then an ϵ -critical point will be found in at most

$$T = \frac{2L}{\epsilon^2} (h(\theta^0) - h(\theta^*)).$$

Proof: Sufficient decrease \rightarrow Telescoping sum and min $\leq \Sigma$.

h satisfies PL-inequality with $\mu > 0$ if

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

Intuition: If θ has small gradient, then it is near-optimal. Let h be L -smooth and μ -PL and $\eta = 1/L$, then

$$h(\theta_T) - h(\theta^*) \leq \left(1 - \frac{\mu}{L}\right)^T (h(\theta_0) - h(\theta^*)).$$

Proof: Sufficient decrease \rightarrow PL \rightarrow Subtract $h(\theta^*)$ both sides.

Acceleration, adaptivity, and momentum

Nesterov acceleration achieves better theoretical guarantees than GD:

$$\mathbf{x}^{t+1} = \theta^t + \gamma(\theta^t - \theta^{t-1}), \quad \theta^{t+1} = \mathbf{x}^{t+1} - \eta \nabla h(\mathbf{x}^{t+1}).$$

Momentum intuition: If gradient is stable, we can make bolder steps. Heavy Ball:

$$\theta^{t+1} = \theta^t - \eta \nabla h(\theta^t) + \beta(\theta^t - \theta^{t-1}), \quad \beta \in [0, 1].$$

Assuming constant gradient δ , we have

$$\theta^{t+1} = \theta^t - \eta \left(\sum_{\tau=1}^{t-1} \beta^\tau \right) \delta.$$

Thus, learning rate increases in case of a stable gradient.

Adaptivity intuition: Different parameters behave differently \rightarrow Parameter-specific learning rates:

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

$$\eta_i^t \doteq \frac{\eta}{\sqrt{\gamma_i^t + \delta}}, \quad \gamma_i^t \doteq \gamma_i^{t-1} + [\partial_i h(\theta^t)]^2.$$

Parameters with small gradient magnitude will have a larger step size.

Adam combines adaptivity and momentum:

$$\mathbf{g}_t = \beta \mathbf{g}_{t-1} + (1 - \beta) \nabla h(\theta_t), \quad \beta \in [0, 1]$$

$$\gamma_t = \alpha \gamma_{t-1} + (1 - \alpha) \nabla h(\theta_t)^{\odot 2}, \quad \alpha \in [0, 1]$$

$$\theta_{t+1} = \theta_t - \eta_t \odot \mathbf{g}_t, \quad \eta_t = 1 \odot (\sqrt{\gamma_t} + \delta).$$

\mathbf{g}_t is a smooth gradient estimator and γ_t measures the stability of the optimization landscape.

Stochastic gradient descent

When the dataset is too large, computing the full gradient is infeasible \rightarrow Estimate gradient with a mini-batch (SGD). SGD outperforms GD in practice, because it has a lower chance of getting stuck in a local optimum due to variance in the gradient estimator.

Convolutional networks

Convolution

Integral operator:

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

Fourier transform:

$$(\mathcal{F}f)(u) \doteq \int_{-\infty}^{\infty} e^{-2\pi i t u} f(t) dt.$$

Convolution:

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

Commutative: $f * h = h * f$, Shift-equivariant: $f_\Delta * h = (f * h)_\Delta$, Convolution as Fourier: $f * h = \mathcal{F}^{-1}(\mathcal{F}f \cdot \mathcal{F}h)$. All proofs are done by defining new variables dependent on existing ones. Linear shift-equivariant operator \iff Convolution.

Discrete convolution:

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

Cross-correlation:

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

Toeplitz matrix $\mathbf{H}_n^h \in \mathbb{R}^{(n+m-1) \times n}$

$$\mathbf{H}_n^h \doteq \begin{bmatrix} h_1 & 0 & \dots & 0 & 0 \\ h_2 & h_1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & h_m & h_{m-1} \\ 0 & 0 & \dots & 0 & h_m \end{bmatrix}.$$

This can then be applied to a vectorized $\mathbf{f} \in \mathbb{R}^n$. Effectively a proof that convolution is linear with increased statistical efficiency.

Convolutional networks

Images are 2D, so use 2D definition of convolution:

$$(\mathbf{X} * \mathbf{W})[i, j] = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} x_{i-k, j-l} w_{k, l}.$$

Let \mathbf{X}^ℓ be output of the ℓ -th convolutional layer and $\Delta^\ell \doteq \frac{\partial \mathcal{L}}{\partial \mathbf{X}^\ell}$, then

$$\Delta^{\ell-1} = \Delta^\ell * \mathbf{W}^\ell, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}^\ell} = \Delta^\ell * \mathbf{X}^{\ell-1}.$$

Max-pooling layer:

$$x_{ij}^\ell = \max\{x_{i+k, j+l}^{\ell-1} \mid k, l \in [0, r)\}$$

$$\frac{\partial x_{ij}^\ell}{\partial x_{m, n}^{\ell-1}} = \mathbb{1}\{(m, n) = (i^*, j^*)\},$$

where (i^*, j^*) are the indices of maximum value in the forward pass.

Data generally has multiple channels:

$$(\mathbf{X} * \mathbf{W})[c, i, j] = \sum_{r=1}^{C_{in}} \sum_{k=-K}^K \sum_{l=-K}^K w_{c, r, k, l} x_{r, i-k, j-l}.$$

Fully connected in channels and local in spatial dimensions.

Recurrent neural networks

Activations are computed recursively to handle variable-length data,

$$\mathbf{z}_t = F[\boldsymbol{\theta}](\mathbf{z}_{t-1}, \mathbf{x}_t).$$

Dependent on application, compute output variables:

$$\mathbf{y}_t = G[\boldsymbol{\varphi}](\mathbf{z}_t).$$

Elman RNN:
 $F[\mathbf{U}, \mathbf{V}](\mathbf{z}, \mathbf{x}) = \phi(\mathbf{U}\mathbf{z} + \mathbf{V}\mathbf{x}), \quad G[\mathbf{W}](\mathbf{z}) = \psi(\mathbf{W}\mathbf{z}).$
Bidirectional RNNs do both ways and concatenate. Stacked RNNs connect layers horizontally:

$$\mathbf{z}_{t,l} = \phi(\mathbf{U}_l \mathbf{z}_{t-1,l} + \mathbf{V}_l \mathbf{z}_{t,l-1}), \quad \mathbf{z}_{t,0} = \mathbf{x}_t.$$

Let $L = \sum_{i=1}^T \ell(\hat{\mathbf{y}}_i, \mathbf{y}_i)$, then we have gradients:

$$\frac{\partial L}{\partial \mathbf{U}} = \sum_{i=1}^T \frac{\partial L}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{U}}, \quad \frac{\partial L}{\partial \mathbf{V}} = \sum_{i=1}^T \frac{\partial L}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{V}},$$

where

$$\frac{\partial L}{\partial \mathbf{z}_i} = \sum_{j=i}^T \frac{\partial \ell(\hat{\mathbf{y}}_j, \mathbf{y}_j)}{\partial \hat{\mathbf{y}}_i} \frac{\partial \hat{\mathbf{y}}_j}{\partial \mathbf{z}_i} \prod_{j=i+1}^i \Phi_j \mathbf{U}_j,$$

where $\Phi_j = \text{diag}(\phi'(\mathbf{U}\mathbf{z}_{j-1} + \mathbf{V}\mathbf{x}_j))$. This is only stable if $\|\Phi_j \mathbf{U}\|_2 = 1$, which is almost never the case \rightarrow exploding/vanishing gradient.

Gated memory

Solve vanishing gradient by gating.

LSTM (\mathbf{z}_t : cell state, $\boldsymbol{\zeta}_t$: hidden state):

$$\mathbf{z}_t = \sigma(\mathbf{F}\hat{\mathbf{x}}_t) \odot \mathbf{z}_{t-1} + \sigma(\mathbf{G}\hat{\mathbf{x}}_t) \odot \tanh(\mathbf{V}\hat{\mathbf{x}}_t)$$
$$\hat{\mathbf{x}}_t = [\boldsymbol{\zeta}_{t-1}, \mathbf{x}_t]$$
$$\boldsymbol{\zeta}_t = \sigma(\mathbf{H}\hat{\mathbf{x}}_t) \odot \tanh(\mathbf{U}\mathbf{z}_t).$$

GRU simplifies the LSTM to only 3 weight matrices:

$$\mathbf{z}_t = \boldsymbol{\sigma} \odot \mathbf{z}_{t-1} + (1 - \boldsymbol{\sigma}) \odot \hat{\mathbf{z}}_t, \quad \boldsymbol{\sigma} = \sigma(\mathbf{G}\hat{\mathbf{x}}_t)$$
$$\hat{\mathbf{x}}_t = [\mathbf{z}_{t-1}, \mathbf{x}_t]$$
$$\hat{\mathbf{z}}_t = \tanh(\mathbf{V}[\boldsymbol{\zeta}_t \odot \mathbf{z}_{t-1}, \mathbf{x}_t]), \quad \boldsymbol{\zeta}_t = \sigma(\mathbf{H}[\mathbf{z}_{t-1}, \mathbf{x}_t]).$$

Linear recurrent model

Simplify GRU to be linear such that it can exploit prefix scan parallelism which has $\mathcal{O}(\log T)$ runtime:

$$\mathbf{z}_t = \boldsymbol{\sigma} \odot \mathbf{z}_{t-1} + (1 - \boldsymbol{\sigma}) \odot \hat{\mathbf{z}}_t, \quad \boldsymbol{\sigma} = \sigma(\mathbf{G}\mathbf{x}_t), \quad \hat{\mathbf{z}}_t = \mathbf{V}\mathbf{x}_t.$$

We can ensure that the gradients do not explode by parametrizing the GRU smartly. The LRU has hidden state evolution in a discrete time linear system:

$$\mathbf{z}_{t+1} = \mathbf{A}\mathbf{z}_t + \mathbf{B}\mathbf{x}_t.$$

Diagonalize $\mathbf{A} = \mathbf{P}\boldsymbol{\Lambda}\mathbf{P}^{-1}$, where $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m), \lambda_i \in \mathbb{C}$. This linear system is stable if the modulus of the eigenvalues is bounded by 1. Thus, we parameterize them such that the moduli can only be in $(0, 1)$ in the following way,

$$\lambda_i = \exp(-\exp(v_i))(\cos(\phi_i) + \sin(\phi_i)i).$$

This uses $\exp(i\theta) = \cos(\theta) + \sin(\theta)i$ and that we can represent complex numbers in polar coordinate form via modulus r and phase ϕ :

$$\mathbf{z} = r(\cos(\phi)\sin(\phi)i), \quad r = |\mathbf{z}|.$$

Thus, $r_i = \exp(-\exp(v_i)) \in (0, 1)$. Thus, at initialization we sample

$$\phi_i \sim \text{Unif}([0, 2\pi]), \quad r_i \sim \text{Unif}(I), \quad I \subseteq [0, 1].$$

And compute $v_i = \log(-\log r_i)$. The modulus always remains upper bounded by 1.

We do not lose any representational power, because we can put all representational power into the output map:

$$\mathbf{y}_t = \text{MLP}(\text{Re}(\mathbf{G}\boldsymbol{\zeta}_t)), \quad \mathbf{G} \in \mathbb{C}^{k \times m}.$$

Sequence learning

Generate sequence step-by-step, given another sequence:

$$p(\mathbf{y}_{1:m} \mid \mathbf{x}_{1:m}) = \prod_{i=1}^n p(\mathbf{y}_i \mid \mathbf{y}_{1:i-1}, \mathbf{x}_{1:m}).$$

This is done by mapping input sequence to latent representation (encoder RNN):

$$\mathbf{x}_1, \dots, \mathbf{x}_m \mapsto \boldsymbol{\zeta}_t.$$

Decoder RNN uses this along with history to predict next token:

$$\boldsymbol{\zeta}_t, \mathbf{y}_1, \dots, \mathbf{y}_{t-1} \mapsto \mathbf{z}_{t-1}.$$

This is mapped to a distribution over next tokens,

$$\mathbf{z}_{t-1} \mapsto \boldsymbol{\mu}_t, \quad \mathbf{y}_t \sim p(\boldsymbol{\mu}_t).$$

Problem: Lossy compression of input sequence. Solution: Use attention such that decoder can look at full input sequence at every step:

$$a_{ij} = \text{softmax}_j(\text{MLP}([z_{i-1}, \boldsymbol{\zeta}_j])), \quad \mathbf{c}_t = \sum_{i=1}^m a_{it} \boldsymbol{\zeta}_i.$$

This allows for alignment between input and output sequence.

Transformers

Self-attention

Let $\mathbf{X} \in \mathbb{R}^{T \times d}$ denote a sequence of input embeddings. Problem: They are non-contextual. Self-attention:
 $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V,$
where $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_k}$ and $\mathbf{W}_V \in \mathbb{R}^{d \times d_v}$. Attention:

$$\mathbf{\Xi} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}.$$

Multi-headed self-attention performs attention h times in parallel.

Cross-attention

Two sequences $\mathbf{A} \in \mathbb{R}^{T_A \times d_A}, \mathbf{B} \in \mathbb{R}^{T_b \times d_b}$ and we want to give information of \mathbf{B} to \mathbf{A} :

$$\mathbf{Q} = \mathbf{A}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{B}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{B}\mathbf{W}_V,$$

where $\mathbf{W}_Q \in \mathbb{R}^{d_A \times d_k}, \mathbf{W}_K \in \mathbb{R}^{d_b \times d_k}$, and $\mathbf{W}_V \in \mathbb{R}^{d_b \times d_v}$.

Positional encoding

Attention is permutation equivariant \rightarrow Add positional encoding matrix $\mathbf{P} \in \mathbb{R}^{T \times d}$ with

$$p_{tk} = \begin{cases} \sin(t\omega_k) & k \bmod 2 = 0 \\ \cos(t\omega_k) & k \bmod 2 = 1, \end{cases} \quad \omega_k \doteq C^{k/d}.$$

Machine translation

Encoder-decoder architecture where the encoder applies MHSA to input sequence and decoder applies masked MHSA to history and then cross-attention with contextualized input sequence. Furthermore, MLP, Layer normalization, and residual layers are used.

BERT

BERT is a transformer-based pretrained language model that is used for finetuning on downstream NLP tasks. BERT tokenizes uses WordPiece tokenization and prepends a [CLS] token. When the weights of the encoders are pretrained, we can place additional layers on top that operate on the contextualized BERT tokens.

Two pre-training stages:

1. Predicting masked out tokens using its left and right context as input (Cloze test; trains BERT's understanding of language);
2. Binary next sentence classification, where the model must classify two sentences as being consecutive or not (trains BERT to infer relationships between sentences).

Vision transformer

ViT adapts the transformer to images by treating projected 16×16 image patches as tokens. A possible reason for this model's effectiveness is that this architecture carries less inductive bias than CNN-based models. In general, this seems to be beneficial for very large datasets.

Geometric deep learning

GDL models neural networks that satisfy invariances by design.

Invariance and equivariance

f (arbitrary number of inputs) is order-invariant iff

$$f(\mathbf{X}) = f(\mathbf{P}\mathbf{X}), \quad \mathbf{X} \in \mathbb{R}^{M \times d},$$

where \mathbf{P} is a permutation matrix.

f (arbitrary number of inputs and same number of outputs) is equivariant iff

$$f(\mathbf{X}) = \mathbf{P}f(\mathbf{P}\mathbf{X}).$$

We want models that have these properties by design.

Deep sets

Let $\phi: \mathbb{R} \rightarrow \mathbb{R}^d$ be a pointwise feature extractor network. Deep Sets obtains an order-invariant representation of the input set by summing their features up. This representation can be given to any network $\rho: \mathbb{R}^d \rightarrow \mathcal{Y}$:

$$f(\mathbf{x}_1, \dots, \mathbf{x}_M) = \rho\left(\sum_{m=1}^M \phi(\mathbf{x}_m)\right).$$

We can easily turn this into an equivariant map by providing \mathbf{x}_m to $\rho: \mathbb{R} \times \mathbb{R}^d \rightarrow \mathcal{Y}$:

$$f(\mathbf{x}_1, \dots, \mathbf{x}_M)_i = \rho\left(\mathbf{x}_i, \sum_{m=1}^M \phi(\mathbf{x}_m)\right).$$

This architecture is universal for a fixed d , but it requires mapping that are highly discontinuous for $M \rightarrow \infty$.

PointNet

PointNet is a Deep Sets architecture on three-dimensional point clouds. The model employs T-net blocks, which apply rigid transformations to the point cloud, which is permutation invariant. These are applied twice alternately with MLPs to form ϕ . This gives a 64-dim intermediate feature vector and 1024-dim final feature vector. The features are aggregated by a max-pool operator.

Object classification: ρ is an MLP with a softmax head that takes the global feature vector as input.

Object segmentation: ρ concatenates intermediate local feature and final global feature, which is given to MLP with softmax head.

Graph neural networks

Let \mathbf{A} be the adjacency matrix of an undirected graph. f is order-invariant on a graph if

$$f(\mathbf{X}, \mathbf{A}) = f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^\top).$$

And equivariant if

$$f(\mathbf{X}, \mathbf{A}) = \mathbf{P}f(\mathbf{P}\mathbf{X}, \mathbf{P}\mathbf{A}\mathbf{P}^\top).$$

Let $\mathbf{X}_m = \{\{\mathbf{x}_n \mid a_{nm} = 1\}\}$ (multiset of neighbors' features). ϕ takes \mathbf{x}_m and \mathbf{X}_m as input (any pair of isomorphic graphs result in same feature representations):

$$\phi(\mathbf{x}_m, \mathbf{X}_m) = \phi\left(\mathbf{x}_m \bigoplus_{\mathbf{x} \in \mathbf{X}_m} \psi(\mathbf{x})\right).$$

This is a message-passing scheme.

Graph convolutional networks (GCNs) aggregates local neighborhoods with a fixed set of weights (coupling matrix),

$$\tilde{\mathbf{A}} \doteq \mathbf{D}^{-1/2}(\mathbf{A} + \mathbf{I})\mathbf{D}^{-1/2}, \quad \mathbf{D} = \text{diag}(\mathbf{d}), \quad d_m = 1 + \sum_{n=1}^M a_{mn}.$$

Now, $\tilde{\mathbf{A}}\mathbf{X}$ computes the average feature over neighbors and the node itself. GCNs introduce learnable parameters \mathbf{W} :

$$f[\mathbf{W}](\mathbf{X}, \mathbf{A}) = \sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}).$$

This is an equivariant function, which can be stacked. Then, we can use the final representations to do graph classification or node classification.

Limitations: Requires a depth equal to diameter of graph to exchange information between all nodes; In very deep GCNs, node features become indistinguishable due to smoothing of $\tilde{\mathbf{A}}$; Bottleneck effect of how much information can be stored in fixed-size representations. There are no canonical solutions to these problems.

GATs introduce attention (which is equivariant) in the neighborhood function and replaces $\tilde{\mathbf{A}}$. It does so by parametrizing the coupling matrix \mathbf{Q} ,

$$q_{mn} = \text{softmax}_n(\rho(\mathbf{u}^\top [\mathbf{V}\mathbf{x}_m, \mathbf{V}\mathbf{x}_n, \mathbf{x}_{mn}]))$$
$$\sum_{n=1}^M a_{mn} q_{mn} = 1, \quad \forall m \in [M].$$

Here, \mathbf{x}_{ij} is an edge feature.

Despite better adaptivity, GATs are still message-passing algorithms. Such algorithms have inherent limitations in the type of graphs they can distinguish. The Weisfeiler-Lehman graph isomorphism test computes whether there exists an isomorphism between two graphs. It can be shown that GCNs and GATs cannot distinguish graphs beyond the WL-test.

Spectral graph theory

Laplacian operator measures local deviation from the mean in vanishingly small neighborhoods:

$$\Delta f = \sum_{i=1}^d \frac{\partial^2 f}{\partial x_i^2}.$$

Graph Laplacian: $\mathbf{L} = \mathbf{D} - \mathbf{A}$. Degree-normalized Laplacian: $\tilde{\mathbf{L}} = \mathbf{D}^{-1/2}(\mathbf{D} - \mathbf{A})\mathbf{D}^{-1/2}$. We can generalize Fourier transform to graphs: Diagonalize $\mathbf{L} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$ and \mathbf{U} can be seen as graph Fourier basis and $\boldsymbol{\Lambda}$ as frequencies. Graph convolution can be computed as pointwise multiplication in Fourier domain:

$$\mathbf{X} * \mathbf{Y} = \mathbf{U}((\mathbf{U}^\top \mathbf{X}) \odot (\mathbf{U}^\top \mathbf{Y})).$$

This can be learned by

$$G[\boldsymbol{\theta}](\mathbf{L})\mathbf{X} = \mathbf{U}G[\boldsymbol{\theta}](\boldsymbol{\Lambda})\mathbf{U}^\top \mathbf{X}.$$

Problem: Eigendecomposition of \mathbf{L} takes $\mathcal{O}(M^3)$. Solution:

Use polynomial kernels:

$$\mathbf{U}\left(\sum_{k=0}^K \alpha_k \boldsymbol{\Lambda}\right)\mathbf{U}^\top \mathbf{X} = \sum_{k=1}^K \alpha_k \mathbf{L}^k \mathbf{X}.$$

Here, the polynomial order K defines the kernel size (or neighborhood size). $\boldsymbol{\alpha} \in \mathbb{R}^K$ are parameters.

Tricks of the trade

Initialization

Parameters are typically chosen with a fixed variance by sampling from

$$\boldsymbol{\theta} \sim \mathcal{N}(0, \sigma^2), \quad \boldsymbol{\theta} \sim \text{Unif}([- \sqrt{3}\sigma, \sqrt{3}\sigma]).$$

LeCun init, $\sigma = 1/\sqrt{n}$, preserves input variance. Glorot init, $\sigma = \sqrt{2/(n+m)}$, normalizes magnitude of gradient (intuition: backpropagation combines n -dim input and m -dim output).

Kaiming init, $\sigma = \sqrt{2/n}$, is designed to be used with ReLU by observing that only half of the units are active in expectation. Orthogonal init considers the weights holistically by initializing as orthogonal matrix (benefit: eigenvalues are ± 1).

Weight decay

$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta(\nabla h(\boldsymbol{\theta}_t) - \mu\boldsymbol{\theta}_t) = (1 - \eta\mu)\boldsymbol{\theta}_t - \nabla h(\boldsymbol{\theta}_t)$.
Equivalent to (1) gradient descent with ℓ_2 regularization, (2) Bayesian prior that weights are sampled from normal distribution, (3) Lagrangian with constraint $\|\boldsymbol{\theta}\| \leq \mu$.

Under the basis of the Hessian's eigenvectors, the optimum of ℓ_2 -regularized ℓ_μ is

$$\boldsymbol{\theta}^* = \text{diag}\left(\frac{\lambda_i}{\lambda_i + \mu}\right)\boldsymbol{\theta}^*.$$

Each axis is scaled based on sensitivity. If $\lambda_i \gg \mu$, then $\lambda_i/\lambda_i + \mu \approx 1$, so the solution does not change much.

Early stopping

Stop if validation performance does not improve for past p epochs. Crudely equivalent to weight decay with μ if stopped at $T \approx 1/\eta\mu$.

Dropout

Randomly disable subset of parameters during training \rightarrow Units become less dependent on one another \rightarrow Instead of units being specialized, they become generally useful.

Two view: (1) regularization and (2) ensemble of networks:

$$p[\boldsymbol{\theta}](\mathbf{y} \mid \mathbf{x}) = \sum_{\mathbf{b} \in \{0,1\}^P} p(\mathbf{b})p[\boldsymbol{\theta} \odot \mathbf{b}](\mathbf{y} \mid \mathbf{x}).$$

This can be approximated by scaling weights by their probability of being active.

Normalization

Goal: Make all units more similar. A unit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ can be normalized by

$$\tilde{f} = \frac{f - \mathbb{E}[f(\mathbf{x})]}{\sqrt{\text{Var}[f(\mathbf{x})]}}.$$

This removes 2 DOF (bias and variance) \rightarrow Explicitly parameterize:

$$\tilde{f}[\mu, \gamma](\mathbf{x}) = \mu + \gamma \tilde{f}(\mathbf{x}).$$

BatchNorm approximates $\mathbb{E}[f]$ and $\text{Var}[f]$ over a mini-batch.

Normalization is very effective and sometimes essential. We used to believe that it was because it helped combat covariance shift. However, a modern motivation shows that normalization is the same as weight normalization and scaling by $\|\mathbf{w}\|_{\mathbf{I}}/\|\mathbf{w}\|_{\Sigma}$, where $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}]$.

Layer normalization estimates mean and variance over the feature dimension instead of batch dimension.

Weight normalization

Normalize weights before applying them:

$$f[\mathbf{v}, \gamma](\mathbf{x}) = \phi(\mathbf{w}^{\top} \mathbf{x}), \quad \mathbf{w} = \frac{\gamma}{\|\mathbf{v}\|_2} \mathbf{v}.$$

Gradients:

$$\frac{\partial h}{\partial \gamma} = \frac{\partial h}{\partial \mathbf{w}} \frac{\partial \mathbf{w}}{\partial \gamma}, \quad \frac{\partial h}{\partial \mathbf{v}} = \frac{\gamma}{\|\mathbf{v}\|} \frac{\partial h}{\partial \mathbf{w}} \left(\mathbf{I} - \frac{\mathbf{w} \mathbf{w}^{\top}}{\|\mathbf{w}\|^2} \right).$$

Here $\mathbf{I} - \mathbf{w} \mathbf{w}^{\top} / \|\mathbf{w}\|^2$ is the projection matrix onto the complement of $\mathbf{w} \rightarrow$ The direction of \mathbf{v} is projected out.

Data augmentation

Transform input data to train invariances.

Label smoothing: Replace labels by noisy probability distributions, because classifiers are not good at dealing with mislabeled data.

Distillation

Let F be the teacher and G its student and we want the student to match its teacher's logits. Then, we have tempered cross-entropy loss:

$$\ell(\mathbf{x}) = \sum_{y \in \mathcal{Y}} \frac{\exp(F_y(\mathbf{x})/T)}{\sum_{y' \in \mathcal{Y}} \exp(F_{y'}(\mathbf{x})/T)} \cdot \left(\frac{1}{T} G_y(\mathbf{x}) - \log \sum_{y' \in \mathcal{Y}} \exp(G_{y'}(\mathbf{x})/T) \right).$$

Gradient:

$$\frac{\partial h}{\partial G_y} = \frac{1}{T} \left(\frac{\exp(F_y(\mathbf{x})/T)}{\sum_{y' \in \mathcal{Y}} \exp(F_{y'}(\mathbf{x})/T)} - \frac{\exp(G_y(\mathbf{x})/T)}{\sum_{y' \in \mathcal{Y}} \exp(G_{y'}(\mathbf{x})/T)} \right).$$

Neural tangent kernel

Linearized model

$$f[\boldsymbol{\theta}] \approx f[\boldsymbol{\theta}_0] + \langle \nabla f[\boldsymbol{\theta}_0], \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle.$$

f is non-linear w.r.t. \mathbf{x} , but linear w.r.t. $\boldsymbol{\theta} \rightarrow$ Define linear model with gradient feature map:

$$h[\boldsymbol{\beta}](\mathbf{x}) = f[\boldsymbol{\theta}_0](\mathbf{x}) + \boldsymbol{\beta}^{\top} \nabla f[\boldsymbol{\theta}_0](\mathbf{x}).$$

Kernel method with $k(\mathbf{x}, \mathbf{x}') = \langle \nabla f[\boldsymbol{\theta}_0](\mathbf{x}), \nabla f[\boldsymbol{\theta}_0](\mathbf{x}') \rangle$ and MSE loss,

$$\boldsymbol{\beta}^* = \boldsymbol{\Phi}^{\top} \mathbf{K}^{-1} (\mathbf{y} - \mathbf{f}), \quad \mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^{\top}, \quad \boldsymbol{\phi}_i = \nabla f[\boldsymbol{\theta}_0](\mathbf{x}_i).$$

Predictions:

$$h^*(\mathbf{x}) = \mathbf{k}(\mathbf{x})^{\top} \mathbf{K}^{-1} (\mathbf{y} - \mathbf{f}).$$

Linearized models are non-competitive with full networks. Also they may be intractable due to number of samples and parameters. Benefit: We can look at DNNs through the lens of kernel methods if the parameters do not evolve far away from $\boldsymbol{\theta}_0$.

Training dynamics

Gradient flow ODE with MSE loss:

$$\dot{\boldsymbol{\theta}}_t = \sum_{i=1}^n (y_i - f[\boldsymbol{\theta}_t](\mathbf{x}_i)) \nabla f[\boldsymbol{\theta}_t](\mathbf{x}_i).$$

Functional gradient flow:

$$\dot{f}[\boldsymbol{\theta}_t] = \boldsymbol{\theta}_t^{\top} \nabla f[\boldsymbol{\theta}_t](\mathbf{x}_j) = \sum_{i=1}^n (y_i - f[\boldsymbol{\theta}_t](\mathbf{x}_i)) k[\boldsymbol{\theta}_t](\mathbf{x}_i, \mathbf{x}_j).$$

In matrix form:

$$\dot{\mathbf{f}}[\boldsymbol{\theta}_t] = \mathbf{K}[\boldsymbol{\theta}_t](\mathbf{y} - \mathbf{f}[\boldsymbol{\theta}_t]), \quad \dot{\mathbf{f}}[\boldsymbol{\theta}_t] = -\mathbf{K}[\boldsymbol{\theta}_t] \nabla_{\mathbf{f}[\boldsymbol{\theta}_t]} \ell(\boldsymbol{\theta}_t).$$

NTK $\mathbf{K}[\boldsymbol{\theta}_t]$ governs the evolution of the joint sample predictions. Problem: NTK has a dependence on the parameters.

Infinite width

In practice it has been found that as the width m of a model is scaled, the parameters stay more closely to their initialization during gradient descent. It can be shown (under basic assumptions) that the NTK converges in probability to a deterministic limit as the model is scaled to infinite width:

$$k[\boldsymbol{\theta}] \xrightarrow{m \rightarrow \infty} k_{\infty}.$$

The deterministic limit depends only on the law of initialization. Under these training dynamics, minimizing MSE equates to solving a kernel regression problem with k_{∞} . This provides insight into why overparameterization works so well in practice, despite having the ability to overfit.

NTK constancy

The NTK remains constant under gradient flow: $\frac{\partial \mathbf{K}[\boldsymbol{\theta}_t]}{\partial t} = \mathbf{0}$.

Bayesian learning

Goal is to compute the Bayesian predictive posterior:

$$f(\mathbf{x}) = \int p(\boldsymbol{\theta} | \mathcal{S}) f[\boldsymbol{\theta}](\mathbf{x}) d\boldsymbol{\theta},$$

where

$$p(\boldsymbol{\theta} | \mathcal{S}) = \frac{p(\mathcal{S} | \boldsymbol{\theta})}{p(\mathcal{S})}, \quad p(\mathcal{S}) = \int p(\boldsymbol{\theta}) p(\mathcal{S} | \boldsymbol{\theta}) d\boldsymbol{\theta}.$$

$p(\mathcal{S})$ is intractable but often not necessary. An isotropic Gaussian prior leads to MAP optimizing ℓ_2 regularization.

Markov chain Monte Carlo

However, we do not want MAP only, we want the full distribution. MCMC methods sample from the posterior by constructing a Markov chain, where the stationary distribution is the posterior.

Detailed Balance Equation:

$$q(\boldsymbol{\theta}) \Pi(\boldsymbol{\theta}' | \boldsymbol{\theta}) = q(\boldsymbol{\theta}') \Pi(\boldsymbol{\theta} | \boldsymbol{\theta}'), \quad \forall \boldsymbol{\theta}, \boldsymbol{\theta}'.$$

Then, the Markov chain is time reversible and has the unique stationary distribution q .

Metropolis-Hastings samples from an arbitrary Markov chain with kernel $\tilde{\Pi}$ and adjusts it such that DBE is satisfied for the posterior. It does so by constructing a new kernel:

$$\Pi(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \tilde{\Pi}(\boldsymbol{\theta}' | \boldsymbol{\theta}) \alpha(\boldsymbol{\theta}' | \boldsymbol{\theta})$$

$$\alpha(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \min \left\{ 1, \frac{p(\boldsymbol{\theta} | \mathcal{S}) \tilde{\Pi}(\boldsymbol{\theta}' | \boldsymbol{\theta})}{p(\boldsymbol{\theta}' | \mathcal{S}) \tilde{\Pi}(\boldsymbol{\theta} | \boldsymbol{\theta}')} \right\}.$$

This is the unique choice of acceptance function α that has a one-sided structure. If $\tilde{\Pi}$ is symmetric, we only need the ratio of posteriors, not $p(\mathcal{S})$. Problems: Burn-in period can be arbitrarily long due to poor $\tilde{\Pi}$ leading to high rejection probabilities.

Hamiltonian Monte Carlo obtains posterior averages. Energy function:

$$E(\boldsymbol{\theta}) = - \sum_{\mathbf{x}, y} \log p[\boldsymbol{\theta}](y | \mathbf{x}) - \log p(\boldsymbol{\theta}).$$

The Hamiltonian augments with momentum vector:

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

Hamiltonian dynamics:

$$\dot{\mathbf{v}} = -\nabla E(\boldsymbol{\theta}), \quad \dot{\boldsymbol{\theta}} = \mathbf{v}.$$

HMC discretizes with stepsize η :

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \mathbf{v}_t, \quad \mathbf{v}_{t+1} = \mathbf{v}_t - \eta \nabla E(\boldsymbol{\theta}_t).$$

We sample the posterior by following momentum-based gradient descent dynamics. (We can also view this as momentum-based gradient descent leading to a single sample approximation of the predictive distribution.) Problem: We need to compute the full gradient.

Langevin dynamics extends HMC by friction:

$$\dot{\boldsymbol{\theta}} = \mathbf{v}, \quad d\mathbf{v} = -\nabla E(\boldsymbol{\theta}) dt - \mathbf{B} \mathbf{v} dt + \mathcal{N}(\mathbf{0}, 2\mathbf{B} dt).$$

Friction reduces momentum and dissipates kinetic energy, while the Wiener noise injects stochasticity. Discretize:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \mathbf{v}_t$$

$$\mathbf{v}_{t+1} = (1 - \eta \gamma) \mathbf{v}_t - \eta \mathbf{s} \nabla \tilde{E}(\boldsymbol{\theta}) + \sqrt{2\gamma\eta} \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Here, \tilde{E} is a stochastic potential function which is the empirical loss over a random mini-batch of data.

Gaussian process

GPs is a fully tractable Bayesian method. f is a GP if for every finite subset $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathcal{X}$, the resulting finite marginal is jointly normally distributed,

$$[f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)] \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

In GPs, the mean can be computed by deterministic regression and the covariance matrix is evaluated by a kernel function:

$$\sigma_{ij} = k(\mathbf{x}_i, \mathbf{x}_j).$$

The kernel function can be seen as a prior over function space that describes how related the output values corresponding to the two input value should be. RBF kernel encodes that close input value should have close output values:

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2).$$

Linear networks assume a random Gaussian weight vector:

$$\mathbf{w} \sim \mathcal{N}\left(\mathbf{0}, \frac{\sigma^2}{d} \mathbf{I}_d\right).$$

Outputs are computed by $y_i = \mathbf{w}^{\top} \mathbf{x}_i$. Vectorized:

$$\mathbf{y} = \mathbf{X} \mathbf{w}.$$

This is a Gaussian vector:

$$\mathbf{y} \sim \mathcal{N}\left(\mathbf{0}, \frac{\sigma^2}{d} \mathbf{X}^{\top} \mathbf{X}\right).$$

In other words, it is a GP with the following kernel:

$$k(\mathbf{x}, \mathbf{x}') = \frac{\sigma^2}{d} \mathbf{x}^{\top} \mathbf{x}'.$$

We can do this for multiple units because the preactivations of units in the same layer are independent, conditioned on the input. In general, we do not get the same effect if we increase the depth, because there is randomness not only in the weights but also the preactivations. However, a deep preactivation process is "near normal" for high-dimensional inputs.

We can extend this to non-linear networks. But, the activations are no longer Gaussian due to the non-linearity. However, due to CLT, they are effectively shaped back into Gaussians when they propagate to the next layer. The mean function and kernels are computed by

$$\mu(\mathbf{x}^{\ell}) = \mathbb{E}[\phi(\mathbf{W}^{\ell-1} \mathbf{x}^{\ell-1})], \quad k^{\ell}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{E}[\phi(\mathbf{x}_i^{\ell-1}) \phi(\mathbf{x}_j^{\ell-1})].$$

We can now use kernel regression:

$$f^*(\mathbf{x}) = \mathbf{k}(\mathbf{x})^{\top} \mathbf{K}^{-1} \mathbf{y}, \quad k = k^L.$$

In conclusion, DNNs in the infinite-width limit can be thought of as GPs (because then all preactivations can be viewed as Gaussians). Benefit: Uncertainty quantification, No training, Leverage tricks form kernel machines. Problems: Computing f^* and storing \mathbf{K}^{ℓ} are not feasible, It is much less efficient than optimizing weights with gradient descent.

Statistical learning theory

VC theory

Let \mathcal{F} be a class of binary classifiers. Then the following is the set of possible classification outcomes over a dataset \mathcal{S} :

$$\mathcal{F}(\mathcal{S}) = \{[f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)] \in \{0, 1\}^n \mid f \in \mathcal{F}\}.$$

Further define maximum number:

$$\mathcal{F}(n) = \sup_{|\mathcal{S}|=n} |\mathcal{F}(\mathcal{S})|.$$

\mathcal{F} shatters \mathcal{S} if $|\mathcal{F}(\mathcal{S})| = 2^n$ (every possible labeling is realized by some function $f \in \mathcal{F}$). VC dimensionality is defined as

$$\text{VC}(\mathcal{F}) = \max_{n \in \mathbb{N}} \sup_{|\mathcal{S}|=n} \mathbb{1}[\mathcal{F}(n) = 2^n].$$

Under uniform convergence, VC inequality holds:

$$\mathbb{P}\left(\sup_{f \in \mathcal{F}} |\hat{\ell}(f) - \ell(f)| > \epsilon\right) \leq 8|\mathcal{F}(n)| \exp\left(-\frac{n\epsilon^2}{32}\right),$$

where ℓ is the expected loss and $\hat{\ell}$ is the empirical loss. Intuition: No generalization guarantees can be given if \mathcal{F} can be fit to any labeling.

Randomization experiments with CIFAR-10 observations:

1. DNNs can perfectly fit the training data;
2. When randomly replacing training labels, the models can still perfectly fit the data (memorization);
3. The training time does not increase much when labels are randomized;
4. When randomly shuffling pixels, the models can also perfectly fit the data \rightarrow Inductive bias of CNNs does not provide much benefit in this regard.

These findings are unexplainable by the above theory.

Overparameterization can lead to double descent phenomenon, where large models will eventually start generalizing better after overfitting.

The flatness of local minima are linked to their generalization ability, because at flat minima, small perturbations in the parameters will only have a small effect on performance. These can be found by small-batch SGD, weight averaging, or entropy SGD.

PAC-Bayesian

For any $p \gg q$ and p -measurable X ,

$$\mathbb{E}_q[X] \leq D_{\text{KL}}(q \| p) + \log \mathbb{E}_p[\exp(X)].$$

For a fixed p , any q , and $\epsilon \in (0, 1)$, we have the following with probability greater than or equal to ϵ ,

$$\mathbb{E}_q[\ell(f)] - \mathbb{E}_q[\hat{\ell}(f)] \leq \sqrt{\frac{2}{n} \left(D_{\text{KL}}(q \| p) + \log \frac{2\sqrt{n}}{\epsilon} \right)}.$$

Here, p is a prior over parameters, q the posterior, and we bound the expected generalization gap over stochastic classifiers $\rightarrow \mathcal{O}(1/\sqrt{n})$ bound on generalization error. However, it only applies to stochastic classifiers, not single classifiers.

This motivated the PAC Bayes loss function

$$\ell_{\text{PAC}}(q) \doteq \mathbb{E}_q[\hat{\ell}] + \sqrt{\frac{2}{n} \left(D_{\text{KL}}(q \| p) + \log \frac{2\sqrt{n}}{\epsilon} \right)}.$$

Effectively just a regularization term. The KL term can be computed in closed form if prior and posterior are Gaussian.

Generative models

Autoencoders

Encoder \mathcal{E} maps data to latents and decoder \mathcal{D} maps latents to data. We want $\mathcal{D}(\mathcal{E}(\mathbf{x})) = \mathbf{x}$. Linear autoencoder with MSE loss is PCA of covariance matrix $\frac{1}{n} \mathbf{X}^{\top} \mathbf{X}$ and taking m principal eigenvectors. Intuition: Retain as much variance as possible.

VAE optimizes log-likelihood of data \rightarrow Intractable \rightarrow ELBO:

$$\log p[\boldsymbol{\theta}](\mathbf{x}) \geq \mathbb{E}_{p[\boldsymbol{\theta}](\mathbf{z}|\mathbf{x})} [\log p[\boldsymbol{\theta}](\mathbf{x} | \mathbf{z})] - D_{\text{KL}}(p[\boldsymbol{\theta}](\mathbf{z} | \mathbf{x}) \| p(\mathbf{z})).$$

$p[\boldsymbol{\theta}]$ is the encoder distribution and $p[\boldsymbol{\theta}]$ is the decoder distribution. This is effectively a reconstruction loss with a regularization term, where the regularization term ensures a well-behaved latent space.

In general, the posterior $p[\boldsymbol{\theta}](\mathbf{z} | \mathbf{x})$ is intractable, so we restrict it to Gaussians,

$$\mathbf{z} | \mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}[\boldsymbol{\theta}](\mathbf{x}), \boldsymbol{\Sigma}[\boldsymbol{\theta}](\mathbf{x})).$$

And the prior is $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Then, the KL divergence can be computed in a closed form:

$$\frac{1}{2} (\|\boldsymbol{\mu}[\boldsymbol{\theta}](\mathbf{x})\|^2 + \text{tr}(\boldsymbol{\Sigma}[\boldsymbol{\theta}](\mathbf{x})) - \log \det(\boldsymbol{\Sigma}[\boldsymbol{\theta}](\mathbf{x})) - m).$$

This can be optimized using the reparameterization trick.

Generative adversarial networks

Log-likelihood is not the only way to optimize a model. GANs provide a training signal by introducing a binary classifier that distinguishes between samples from “nature” (1) and the generator (0): discriminator.

Augmented distribution over samples:

$$\tilde{p}(\boldsymbol{x}, y) = \frac{1}{2} (y p(\boldsymbol{x}) + (1 - y) p[\boldsymbol{\theta}](\boldsymbol{x})).$$

Bayes-optimal classifier:

$$\mathbb{P}(y = 1 \mid \boldsymbol{x}) = \frac{p(\boldsymbol{x})}{p(\boldsymbol{x}) + p[\boldsymbol{\theta}](\boldsymbol{x})}.$$

Minimizing the logistic log-likelihood w.r.t. this discriminator gives the following loss for the generator,

$$\ell^*(\boldsymbol{\theta}) = D_{\text{JS}}(p \parallel p[\boldsymbol{\theta}]) - \log 2,$$

where $D_{\text{JS}}(p \parallel q) = H(p + p[\boldsymbol{\theta}]/2) - H(p) + H(p[\boldsymbol{\theta}])/2$. But, the optimal classifier is intractable, so we train a parametrized one $q[\boldsymbol{\varphi}]$,

$$\boldsymbol{\theta}^*, \boldsymbol{\varphi}^* \in \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \underset{\boldsymbol{\varphi}}{\operatorname{argmax}} \ell(\boldsymbol{\theta}, \boldsymbol{\varphi}),$$

where

$$\ell(\boldsymbol{\theta}, \boldsymbol{\varphi}) = \mathbb{E}_{p[\boldsymbol{\theta}]} [\boldsymbol{y} \log q[\boldsymbol{\varphi}](\boldsymbol{x}) + (1 - \boldsymbol{y}) \log (1 - q[\boldsymbol{\varphi}](\boldsymbol{x}))].$$

We have the following bound:

$$\ell^*(\boldsymbol{\theta}) \geq \sup_{\boldsymbol{\varphi}} \ell(\boldsymbol{\theta}, \boldsymbol{\varphi}).$$

Problem: Gradient descent-ascent is not guaranteed to

converge. Solution: Extragradient optimization algorithm:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \boldsymbol{\nabla}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_{t+1/2}, \boldsymbol{\varphi}_t), \quad \boldsymbol{\theta}_{t+1/2} = \boldsymbol{\theta}_t - \eta \boldsymbol{\nabla}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_t, \boldsymbol{\varphi}_t)$$

$$\boldsymbol{\varphi}_{t+1} = \boldsymbol{\varphi}_t + \eta \boldsymbol{\nabla}_{\boldsymbol{\varphi}} \ell(\boldsymbol{\theta}_t, \boldsymbol{\varphi}_{t+1/2}), \quad \boldsymbol{\varphi}_{t+1/2} = \boldsymbol{\varphi}_t + \eta \boldsymbol{\nabla}_{\boldsymbol{\varphi}} \ell(\boldsymbol{\theta}_t, \boldsymbol{\varphi}_t).$$

In practice, we also need to use a different loss function for the generator:

$$\ell(\boldsymbol{\theta} \mid \boldsymbol{\varphi}) = \mathbb{E}_{p[\boldsymbol{\theta}]} [-\log q[\boldsymbol{\varphi}](\boldsymbol{x})],$$

because otherwise the gradient goes to infinity when

$q[\boldsymbol{\varphi}](\boldsymbol{x}) = 1$, which makes the generator saturate.

Diffusion models

Map a simple distribution to a complex one in many steps:

$$\pi = \pi_T \mapsto \pi_{T-1} \mapsto \cdots \mapsto \pi_0 \approx p.$$

SDE view:

$$\mathrm{d}\boldsymbol{x}_t = -\frac{1}{2} \beta_t \boldsymbol{x}_t \mathrm{d}t \sqrt{\beta_t} \mathrm{d}\boldsymbol{w}_t.$$

Time-reversed:

$$\mathrm{d}\boldsymbol{x}_t = \left(-\frac{1}{2} \beta_t \boldsymbol{x}_t - \beta_t \boldsymbol{\nabla}_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t) \right) \mathrm{d}t + \sqrt{\beta_t} \mathrm{d}\tilde{\boldsymbol{w}}_t.$$

Denoising amounts to approximating a vector field over the gradient of the probability distribution moving towards areas with high probability density. Score models approximate $\boldsymbol{\nabla}_{\boldsymbol{x}_t} \log q_t(\boldsymbol{x}_t)$.

ELBO view: Forward process:

$$\boldsymbol{x}_t = \sqrt{1 - \beta_t} \boldsymbol{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon}_t, \quad \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I}).$$

Energy of the stochastic process evolves as $\mathbb{E}[\|\boldsymbol{x}_t\|^2 \mid \boldsymbol{x}_{t-1}] = (1 - \beta_t) \|\boldsymbol{x}_{t-1}\|^2 + \beta_t \operatorname{tr}(\boldsymbol{I})$. If $\mathbb{E}[\|\boldsymbol{x}_0\|^2] = \operatorname{tr}(\boldsymbol{I}) = \dim(\boldsymbol{x}_0)$, then energy is conserved throughout the process.

Closed form ($\tilde{\boldsymbol{\alpha}}_t = \prod_{\tau=1}^t (1 - \beta_{\tau})$):

$$\boldsymbol{x}_t \sim \mathcal{N}(\sqrt{\tilde{\boldsymbol{\alpha}}_t}, (1 - \tilde{\boldsymbol{\alpha}}_t) \boldsymbol{I}).$$

Backward process:

$$\boldsymbol{x}_{t-1} \sim \mathcal{N}(\boldsymbol{\mu}[\boldsymbol{\theta}](\boldsymbol{x}_t, t), \boldsymbol{\Sigma}[\boldsymbol{\theta}](\boldsymbol{x}_t, t)).$$

ELBO:

$$\log p[\boldsymbol{\theta}](\boldsymbol{x}_0) \geq \sum_{t=0}^T \ell_t,$$

where

$$\ell_t = \begin{cases} \mathbb{E}_q[\log p[\boldsymbol{\theta}](\boldsymbol{x}_0 \mid \boldsymbol{x}_1)] & t = 0 \\ -D_{\text{KL}}(q(\boldsymbol{x}_t \mid \boldsymbol{x}_{t-1}, \boldsymbol{x}_0) \parallel p[\boldsymbol{\theta}](\boldsymbol{x}_t \mid \boldsymbol{x}_{t+1})) & 0 < t < T \\ -D_{\text{KL}}(q(\boldsymbol{x}_T \mid \boldsymbol{x}_0) \parallel \pi) & t = T. \end{cases}$$

The KL divergences can be analytically computed because all q_t are Gaussians an we parameterized the network as a Gaussian. The q targets are derived as

$$q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) = \mathcal{N}(\boldsymbol{\mu}(\boldsymbol{x}_t, \boldsymbol{x}_0, t), \tilde{\beta}_t),$$

where

$$\boldsymbol{\mu}(\boldsymbol{x}_t, \boldsymbol{x}_0, t) = \frac{\sqrt{\tilde{\boldsymbol{\alpha}}_{t-1}} \tilde{\beta}_t}{1 - \tilde{\boldsymbol{\alpha}}_t} \boldsymbol{x}_0 + \frac{1 - \tilde{\boldsymbol{\alpha}}_{t-1}}{1 - \tilde{\boldsymbol{\alpha}}_t} \sqrt{1 - \beta_t} \boldsymbol{x}_t$$

$$\tilde{\beta}_t = \frac{1 - \tilde{\boldsymbol{\alpha}}_{t-1}}{1 - \tilde{\boldsymbol{\alpha}}_t} \beta_t.$$

Thus ℓ_t simplify to

$$\ell_t = -\frac{1}{2\sigma_t^2} \|\boldsymbol{\mu}(\boldsymbol{x}_t, \boldsymbol{x}_0, t) - \boldsymbol{\mu}[\boldsymbol{\theta}](\boldsymbol{x}_t, t)\|^2,$$

where $\sigma_t^2 \in [\beta_t, \tilde{\beta}_t]$.

By noting the closed form forward process, we have

$$\boldsymbol{x}_0 = \frac{1}{\sqrt{\tilde{\boldsymbol{\alpha}}_t}} \boldsymbol{x}_t - \frac{\sqrt{1 - \tilde{\boldsymbol{\alpha}}_t}}{\tilde{\boldsymbol{\alpha}}_t} \boldsymbol{\epsilon}.$$

We can thus rewrite

$$\boldsymbol{\mu}(\boldsymbol{x}_t, \boldsymbol{x}_0, t) = \frac{1}{\sqrt{\tilde{\boldsymbol{\alpha}}_t}} \left(\boldsymbol{x}_t - \frac{\beta_t}{\sqrt{1 - \tilde{\boldsymbol{\alpha}}_t}} \boldsymbol{\epsilon} \right).$$

Note that $\boldsymbol{\epsilon}$ fully determines \boldsymbol{x}_t and \boldsymbol{x}_0 is constant. So, we only need to predict $\boldsymbol{\epsilon}$. Simplified loss:

$$\mathbb{E}_q[\ell_t \mid \boldsymbol{x}_0] = \mathbb{E}_{\boldsymbol{\epsilon}}[\lambda(t) \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}[\boldsymbol{\theta}](\boldsymbol{x}_t, t)\|^2]$$

$$\lambda(t) = \frac{\beta_t^2}{2\sigma_t^2 \tilde{\boldsymbol{\alpha}}_t (1 - \tilde{\boldsymbol{\alpha}}_t)}.$$

In practice, the loss is approximated by

$$\ell(\boldsymbol{\theta} \mid \boldsymbol{x}_0) = \frac{1}{T} \sum_{t=1}^T [\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}[\boldsymbol{\theta}](\boldsymbol{x}_t, t)\|^2 \mid \boldsymbol{x}_0].$$

Entropy bound:

$$H(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t) \leq H(\boldsymbol{x}_t \mid \boldsymbol{x}_{t-1}).$$

The entropy of the reverse process is bounded by the entropy of the forward process.

Adversarial attacks

The attacker wants to make small changes to the input such that the model gives a different result.

p-norm robustness

Consider a multi-class classifier $f : \mathbb{R}^d \rightarrow [m]$. The goal of an adversarial attack is to find a perturbation $\boldsymbol{\eta}$ such that

$$f(\boldsymbol{x} + \boldsymbol{\eta}) \neq f(\boldsymbol{x}), \quad \|\boldsymbol{\eta}\|_p \leq \epsilon.$$

Consider a binary affine classifier and $p = 2$,

$$f(\boldsymbol{x}) = \operatorname{argmax}\{\boldsymbol{w}_1^\top \boldsymbol{x} + b_1, \boldsymbol{w}_2^\top \boldsymbol{x} + b_2\}.$$

Assume \boldsymbol{x} is classified as 1 and we want to find $\boldsymbol{\eta}$ such that $\boldsymbol{x} + \boldsymbol{\eta}$ is classified as 2 such that $\|\boldsymbol{\eta}\|_2$ is minimized \rightarrow Convex program:

$$\text{minimize} \quad \frac{1}{2} \|\boldsymbol{\eta}\|_2^2$$

$$\text{subject to} \quad (\boldsymbol{w}_1 - \boldsymbol{w}_2)^\top (\boldsymbol{x} + \boldsymbol{\eta}) + b_1 - b_2 \leq 0.$$

Set gradient of Lagrangian to zero $\rightarrow \boldsymbol{\eta}^* = \lambda(\boldsymbol{w}_2 - \boldsymbol{w}_1)$. Then, find λ that satisfies constraint $\rightarrow \lambda \geq f_1(\boldsymbol{x}) - f_2(\boldsymbol{x}) / \|\boldsymbol{w}_1 - \boldsymbol{w}_2\|_2^2$.

Thus:

$$\boldsymbol{\eta}^* = \frac{f_1(\boldsymbol{x}) - f_2(\boldsymbol{x})}{\|\boldsymbol{w}_2 - \boldsymbol{w}_1\|_2^2} (\boldsymbol{w}_2 - \boldsymbol{w}_1).$$

This can be generalized to any source i and target j . In the general case, we can linearize the model and iteratively solve the above convex program.

Robust training

Robust training systematically makes models robust to adversarial attacks by extending the loss function to neighborhoods of training points:

$$\ell(\boldsymbol{x}) \mapsto \max_{\boldsymbol{\eta}: \|\boldsymbol{\eta}\|_p \leq \epsilon} \ell(\boldsymbol{x} + \boldsymbol{\eta}).$$

This can be solved with projected gradient ascent. For $p = 2$, we have

$$\boldsymbol{\eta}_{t+1} = \epsilon \Pi(\boldsymbol{\eta}_t + a \boldsymbol{\nabla}_{\boldsymbol{x}} \ell(\boldsymbol{x} + \boldsymbol{\eta}_t)), \quad \Pi(\boldsymbol{z}) \doteq \frac{\boldsymbol{z}}{\|\boldsymbol{z}\|_2}.$$

Fast gradient sign method performs one iteration with $p = \infty$ resulting in $\boldsymbol{\eta} = \epsilon \cdot \operatorname{sgn}(\boldsymbol{\nabla}_{\boldsymbol{x}} \ell(\boldsymbol{x}))$