

ESSENTIALS

Matrix/Vector Basics

Multiplication: $C = AB \Leftrightarrow c_{ik} = \sum_{j=1}^m a_{ij} b_{jk}$
Orthogonal Matrix: (full rank, square matrix with orthonormal columns) $A^T = A^{-1}$, $AA^T = A^TA = I$
 $\det(A) \in \{-1, 1\}$, $\det(AA^T) = 1$. preserves: inner product norm, distance, angle, rank, matrix orthogonality

Inner product: $\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$.

$$\langle x \pm y, x \pm y \rangle = \langle x, x \rangle \pm 2\langle x, y \rangle + \langle y, y \rangle$$

$$\langle x, y+z \rangle = \langle x, z \rangle + \langle x, y \rangle. \quad \langle x, y \rangle = \|x\|_2 \|y\|_2 \cos(\theta)$$

Outer product: $UV^T \Leftrightarrow (UV^T)_{ij} = u_i v_j$

$$\text{Norms: } \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\langle x, x \rangle}, \quad \|M\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n M_{ij}^2}$$

Probability / Statistics Basics

$$P(x) = \sum_{y \in Y} P(x,y) \quad P(x,y) = P(x|y) P(y)$$

$$\forall y \in Y: \sum_{x \in X} P(x|y) = 1$$

$$P(x|y) = \frac{P(x,y)}{P(y)} \quad \text{if } P(y) > 0$$

$$\text{posterior } p(A|B) = \frac{\text{prior } p(A) \text{ likelihood } p(B|A)}{\text{evidence } p(B)}$$

If i.i.d: $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i)$

If X & Y independent: $P(x|y) = P(x) \Leftrightarrow P(y|x) = P(y)$

Expectation: $E[f(X)] = \sum_{x \in X} P(x) f(x) = \int P(x) f(x) dx$

Variance: $\text{Var}[X] = E[(X - \mu)^2] = \sum_{x \in X} (x - \mu)^2 P(x) = E[X^2] - E[X]^2$

Expectation Properties

Linearity: $E[X+Y] = E[X] + E[Y]$, $E[aX] = aE[X]$

For constant: $X = c$ then $E[X] = c$, so we can take the expectation of unrelated stuff, and get the same: $\log P_\theta(x^{(i)}) = E_{z \sim q_\theta}(z|x^{(i)}) [\log P_\theta(z)]$

Double: For any X : $E[E[X]] = E[X]$.

Convex Function

$\forall x_1, x_2 \in X, \forall t \in [0, 1]:$ (also if $\forall x f'(x) \leq 0$)

$f(tx_1 + (1-t)x_2) \leq f(x_1) + (1-t)f(x_2)$

Example: Show that if f is convex, any local optimum is global. Assume there is a local optimum \hat{x} that is not the global optimum x^* , then if we choose t in the ball of the local optimum we know that: $f(t\hat{x} + (1-t)x^*) \geq f(\hat{x})$. Since $f(x^*) < f(\hat{x})$, we have $t \cdot f(\hat{x}) + (1-t) \cdot f^*(x^*) < f(\hat{x})$. So we get $f(t\hat{x} + (1-t)x^*) \geq f(\hat{x}) > t \cdot f(\hat{x}) + (1-t) \cdot f^*(x^*)$ which contradicts the convexity of f .

Jensen Inequality

For convex ϕ : $\phi(E[X]) \leq E[\phi(X)]$ and also

$\forall \alpha_i \geq 0, \sum \alpha_i = 1$ and any $x_i > 0$: $\log(\sum \alpha_i x_i) \geq \sum \alpha_i \log(x_i)$

Matrix calculus (numerator layout)

$$\frac{\partial \mathbf{v}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial v_1}{\partial x_1} & \dots & \frac{\partial v_n}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_1}{\partial x_m} & \dots & \frac{\partial v_n}{\partial x_m} \end{bmatrix} \quad \frac{\partial \mathbf{x}}{\partial \mathbf{y}} = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \dots & \frac{\partial x_n}{\partial y_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_1}{\partial y_m} & \dots & \frac{\partial x_n}{\partial y_m} \end{bmatrix} \quad \frac{\partial \mathbf{v}}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial v_1}{\partial z_1} & \dots & \frac{\partial v_n}{\partial z_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_1}{\partial z_m} & \dots & \frac{\partial v_n}{\partial z_m} \end{bmatrix}$$

$$\frac{\partial \mathbf{A}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial a_{11}}{\partial x_1} & \dots & \frac{\partial a_{1n}}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial a_{m1}}{\partial x_1} & \dots & \frac{\partial a_{mn}}{\partial x_1} \end{bmatrix} \quad \frac{\partial \mathbf{x}}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial x_1}{\partial a_{11}} & \dots & \frac{\partial x_1}{\partial a_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_1}{\partial a_{m1}} & \dots & \frac{\partial x_1}{\partial a_{mn}} \end{bmatrix}$$

1D Derivation Rules

$$\frac{\partial(f \circ g)}{\partial x} = f' \cdot g + g' \cdot f \quad \frac{\partial f}{\partial x} = \frac{f'g - g'f}{g^2} \quad \frac{\partial f(g(x))}{\partial x} = f'(g(x)) \cdot g'(x)$$

$$\frac{\partial x^r}{\partial x} = r \cdot x^{r-1} \quad \frac{\partial 1}{\partial x f(x)} = \frac{f'(x)}{(f(x))^2} \quad (f \circ g)' = f'(f(g(x))) \cdot g'(x)$$

Multivariate Chain Rule

For function $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $x \mapsto y$ and $f: \mathbb{R}^m \rightarrow \mathbb{R}^l$, $y \mapsto z$ we have: $\frac{\partial z_k}{\partial x_j} = \sum_{i=1}^m \frac{\partial z_k}{\partial y_i} \frac{\partial y_i}{\partial x_j}$
 $\frac{\partial z_k}{\partial x} = \frac{\partial z_k}{\partial y} \frac{\partial y}{\partial x}$ $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$
 $\frac{\partial f(x, y, t)}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$ x, y, t can be scalar, vector or matrix.

Derivatives:

$$\frac{\partial}{\partial x}(b^T x) = \frac{\partial}{\partial x}(x^T b) = b \quad \frac{\partial}{\partial x}(x^T x) = 2x \quad \frac{\partial}{\partial x}(x^T Ax) = (A+A^T) \times x \\ \frac{\partial}{\partial x}(b^T Ax) = A^T b \quad \frac{\partial}{\partial x}(x^T b) = cb^T \quad \frac{\partial}{\partial x}(c^T x^T b) = bc^T \\ \frac{\partial}{\partial x}(\|x-b\|_2) = \frac{x-b}{\|x-b\|_2} \quad \frac{\partial}{\partial x}(\|x\|_2^2) = 2x \quad \frac{\partial}{\partial x}(\|x\|_2) = 2x \\ \frac{\partial}{\partial x} \log(x) = \frac{1}{x} \quad \frac{\partial}{\partial x} \frac{1}{f(x)} = -\frac{f'(x)}{(f(x))^2} \quad \frac{\partial}{\partial x} |x| = (-1)^{|x|} \quad |x| > 0$$

Kullback - Leibler Divergence

$$D_{KL}(P||Q) = -\sum_{x \in X} P(x) \log \left(\frac{Q(x)}{P(x)} \right) = \sum_{x \in X} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx$$

Is not symmetric: $D_{KL}(p||q) \neq D_{KL}(q||p)$

Entropy: $H(p) = -\sum_i p_i \log p_i$

Cross entropy: $H(p, q) = E_p[-\log(q)] = H(p) + D_{KL}(p||q)$

Jensen - Shannon Divergence

$$D_{JS}(P||Q) = \frac{1}{2} D_{KL}(P||\frac{1}{2}(P+Q)) + \frac{1}{2} D_{KL}(Q||\frac{1}{2}(P+Q))$$

NEURAL NETWORK BASICS

Perceptron learning: Init $w^{(0)} = 0$, then while $\exists y$ that is misclassified do $\{w^{(k+1)} = w^k + \eta(y - \hat{y})x^k\}$ where $\hat{y} = (\omega^k x^k) > 0\}$ if lin. sep. converges in finite time.

Sigmoid function: $\text{sigm}(x) = 1/(1+e^{-x}) = e^x/(e^x+1)$ for $\sigma(w^T x + b)$: 

Logistic regression (MLE): Given $D = \{(x^{(i)}, y^{(i)}) \}_{i=1}^n$

Model: $y^{(i)} \sim \text{Bernoulli}(\sigma(w^T x))$

$$\text{MLE: } W_{\text{MLE}} = \underset{W}{\operatorname{argmax}} P(D|W) = \underset{W}{\operatorname{argmax}} \prod_{i=1}^n P(y^{(i)}|x^{(i)}, W) = \underset{W}{\operatorname{argmax}} \prod_{i=1}^n [\sigma(w^T x^{(i)})]^{y^{(i)}} [1 - \sigma(w^T x^{(i)})]^{1-y^{(i)}}$$

Neg. log. likelihood: $-\log P(D|W) = \sum_i y^{(i)} \log \pi_i + (1-y^{(i)}) \log (1-\pi_i)$

Softmax function: $\text{softmax}(x_i) = e^{x_i} / \sum_i e^{x_i}$

Loss Functions

Squared Loss: $\frac{1}{2}(y - \hat{y})^2$

Cross entropy: $-\log(\hat{y}) - (1-y) \log(1-\hat{y})$

Negative log-likelihood: $-\sum_i y_i \log \hat{y}_i$

Why not accuracy: is a step function so a small change in weights might not change accuracy. Gradient is either zero or undefined.

MLE: Write down prob. dist, decompose into per sample prob., minimize negative log likelihood.

Universal Approximation Theorem

$\exists g(x)$ as NN (with non-linear activation), $g(x) \approx f(x)$ and $|g(x) - f(x)| < \epsilon \quad \forall x \in C^0(\mathbb{R})$ (continuous functions given enough hidden units, one layer is enough)

Training NN: Init weights to small random values, init biases to 0 or small positive values, update the params with GD. (non-lin. activation make loss non-convex)

TRAINING NEURAL NETWORKS

Regularization: Any technique which aims to reduce generalization error (not training error)

Data Standardization: $x_i = (x - \mu)/\sigma$ where $\mu = \sum_{i=1}^n x_i$ and $\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$. Make sure to use μ and σ for the test set!

Batch Normalization

Apply Batch Norm on each layer while computing θ^i, μ over the entire batch.

$$\mu = \frac{1}{n} \sum_{i=1}^n z_i \quad \theta^i = \frac{1}{n-1} \sum_{i=1}^n (z_i - \mu)^2 \quad z_i = \frac{x_i - \mu}{\sqrt{\theta^i + \epsilon}} \quad \tilde{z}_i = \theta^i z_i + \beta$$

with $\epsilon > 0$ and θ, β learnable parameters

- Bias in linear / conv layer becomes redundant w/BN

- Makes weights in deeper layers more robust to changes in shallow layers - Slight regularization effect

- Test time: μ, σ estimated using exponential average of training time values

Data Augmentation: Generally more data leads to better performance, but acquiring training data is tedious (especially labels). Can generate synthetic data with a model-based approach, but that can lead to discrepancies between the real and synthetic data. Can also expand the data set by transforming existing samples to produce new samples (rotate, scale, noise injection (Gaussian blur, Salt & Pepper noise, channelwise dropout)), need to ensure consistency between transformed data sample and label. For classification introduces invariance $M(T(I)) = M(I)$, for regression introduces equivariance $M(T(I)) = T(M(I))$

Semi-supervised Learning: Acquiring labeled data is expensive, but unlabeled data usually plentiful. Want to learn from both. Assumes $P(x)$ related to $P(y|x)$. Example of this:

Multitask learning: for datasets with similar samples and related labels. Divide model into shared and task-specific parameters. Assumes that $P(y_i|x)$ related to $P(y_j|x)$, ...

Activation Functions:

Sigmoid: $\sigma(x) = 1/(1+e^{-x})$, $\sigma'(x) = \sigma(x)(1-\sigma(x))$ as interpretation as the firing rate of a neuron. Finite range \rightarrow stable training

- Saturation towards either end, so large activations get small gradients. Not zero-centered, always positive (cannot change "direction").

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

+ has finite range, zero-centered, has stronger steeper gradients than σ . - activations also saturate.

Rectified Linear Unit: $\text{ReLU}(x) = \max(0, x)$
 $\text{ReLU}'(x) = 1 \text{ if } x > 0 \text{ else } 0$

+ inexpensive operation, sparse activations, greatly accelerates convergence.

- can blow up the activation and destabilize training. "dying" ReLU problem: units with negative activations get no update

Leaky ReLU: with small pos. const. α : $L\text{ReLU}(x) = \begin{cases} x & \text{for } x > 0 \\ \alpha x & \text{for } x \leq 0 \end{cases}$ $L\text{ReLU}'(x) = \begin{cases} 1 & \text{for } x > 0 \\ \alpha & \text{for } x \leq 0 \end{cases}$

+ Attempts to solve "dying" ReLU by introducing small negative slope for neg. activ.

Randomized Leaky ReLU: $\alpha \sim U(a, b)$: $f(x) = \begin{cases} \alpha x & \text{for } x > 0 \\ x & \text{for } x \leq 0 \end{cases}$ $f'(x) = \begin{cases} \alpha & \text{for } x > 0 \\ 1 & \text{for } x \leq 0 \end{cases}$

At test time set $\alpha = (a+b)/2$.

Optimization Algorithms:

Gradient Descent: $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$. Follows the direction of the loss surface downhill.

SGD: $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$, + prevents redundant calculations, enables to jump to new and potentially better local minima. - has risk of overshooting.

SGD Problem: Gradient big in the direction where we do not need to travel and small in the direction where we want to travel \rightarrow oscillations, slow convergence

Mini-batch GD: $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; X^{(1:n)}, Y^{(1:n)})$ + reduces the variance of the parameter updates which can lead to more stable convergence, can make use of highly optimized matrix operations.

Momentum: $V_t = \gamma V_{t-1} + \eta \nabla_{\theta} J(\theta)$, $\theta = \theta - V_t$. Instead of using the gradient to change the position of a "weight particle", use it to change the velocity. (helps with ravines i.e. areas where the surface curves more steeply in one dimension than the other dimension)

Nesterov accelerated gradient (NAG): $\theta = \theta - V_t$, $V_t = \gamma V_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma V_{t-1})$. Looks ahead where the momentum would take it, then calculate there

AdaGrad

Compute Gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}, \theta), y^{(i)})$

Accumulate S_t: $r \leftarrow r + g \otimes g$

Compute Update: $\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{r}} \odot \theta$

Apply Update: $\theta \leftarrow \theta + \Delta \theta$

+ greater progress in more gently sloped parts of θ -space
- accumulation of r results in too early decrease in LR

RMS Prop

Compute Gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}, \theta), y^{(i)})$

Accumulate S_t: $r \leftarrow r + P_t + (1-P_t) g \otimes g$

Compute Update: $\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{r}} \odot \theta$

Apply Update: $\theta \leftarrow \theta + \Delta \theta$

- P_t exponentially saturates at 1

+ performs better in non-convex settings

Adaptive Moment Estimation (Adam):

$m_t = B_1 m_{t-1} + (1-B_1) g_t$ estimate of the first moment

$v_t = B_2 v_{t-1} + (1-B_2) g_t^2$ estimate of the second moment

$\hat{m}_t = \bar{m}_t / (1-B_1^t)$, $\hat{v}_t = \bar{v}_t / (1-B_2^t)$, $\theta_{t+1} = \theta_t - \eta_t \sqrt{\hat{v}_t} + \epsilon \hat{m}_t$

Learning rate: If too large, can cause very high lr

overshoot, if too low can be slow and get trapped in local minima.

Step decay: decay learning rate by half after certain epochs.

Exponential decay: $\eta_{t+1} = \eta_0 e^{-kt}$

1/t decay: $\eta_{t+1} = \eta_0 / (1+kt)$

Fast Geometric Ensembles

- fewer surfaces of constant loss to explore/min

① 80% normal training ② 20% use cyclic LR

③ save checkpoints (fixed)

④ ensemble for inference

Stochastic Weight Averaging

① 80% normal training ② 20% of time re-init with w_{swa} & adopt cyclic LR

③ when LR lowest update w_{swa}

$w_{swa} \leftarrow \frac{w_{swa} \cdot n_{models} + w}{n_{models} + 1}$

④ Use w_{swa} for inference (Batch Norm

has to be re-initialized for w_{swa})

HMAX Model: simple (S) cells tuned to specific features ($y = \exp(-\frac{1}{2\sigma^2} \sum_{j=1}^n (w_j \cdot x_j)^2)$), complex (C) cells combine output of various S cells to increase invariance and receptive field ($y = \max_{j=1, \dots, C_h} x_j$)

CONVOLUTIONAL NEURAL NETWORKS

Linear Transform: $T(u+v) = \alpha T(u) + \beta T(v)$

Invariant to f: $T(f(u)) = T(u)$

Equivariant to f: $T(f(u)) = f(T(u))$

↳ any linear, shift-equivariant transform T can be written as a convolution.

Correlation: $I'(i,j) = \sum_{k=-K}^K \sum_{l=-L}^L K(m,n) I(i+m, j+n)$

Convolution: $I'(i,j) = \sum_{m=-M}^M \sum_{n=-N}^N K(-m,-n) I(i+m, j+n)$

Called shift-invariant if K does not depend on (i,j). Same for $K(i,j) = K(-i,-j)$ kernels.

and kernels: matrix-vector mult. ($I \cdot K = \sum_{i=1}^I \sum_{j=1}^J K_{ij} I_{ij}$)

Differentiation: is a linear, shift invariant

operation, can write as convolution: $\frac{\partial f}{\partial x} = \sum_{x \in \mathcal{X}} f(x) - f(x-\Delta x)$

CNN Backprop: $z^{(0)} = w^{(1)} * z^{(1-1)} + b^{(1)} = \sum_m \sum_n w_{m,n}^{(1)} z^{(1-1)}_{m-n} + b^{(1)}$

Backward pass w.r.t. inputs

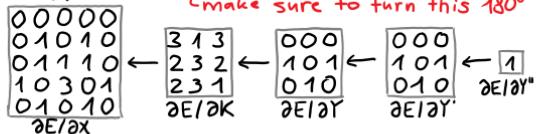
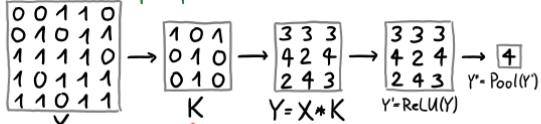
$$\delta^{(1-1)}_{i,j} = \frac{\partial C_{i,j}}{\partial z^{(1-1)}} = \sum_{i' \neq i} \sum_{j' \neq j} \frac{\partial C_{i,j}}{\partial z^{(1-1)}} \frac{\partial z^{(1-1)}}{\partial z^{(1-1)}} = \sum_{i' \neq i} \sum_{j' \neq j} \delta^{(1)}_{i',j'} \cdot w_{i'-i, j'-j}^{(1)} = \text{where } i'=i-m, j'=j-n. \text{ What remains is } w_{m,n}, \text{ but as } i=m-n, m=n-i, j=j-n, j=n-j$$

$$= \sum_{i' \neq i} \sum_{j' \neq j} \delta^{(1)}_{i',j'} \cdot w_{i'-i, j'-j}^{(1)} = \delta^{(1)} * W_{i,-i, j,-j}^{(1)} = \delta^{(1)} * \text{ROT}_{180}(W^{(1)})$$

Weight Update:

$$\frac{\partial C}{\partial w_{m,n}} = \sum_{i' \neq i} \sum_{j' \neq j} \frac{\partial C}{\partial z^{(1-1)}} \frac{\partial z^{(1-1)}}{\partial w_{m,n}} = \sum_{i' \neq i} \sum_{j' \neq j} \delta^{(1)}_{i',j'} \frac{\partial}{\partial w_{m,n}} \left(\sum_{i''=m}^M \sum_{j''=n}^N w_{i'',j''} z^{(1-1)}_{i'',j''} + b^{(1)} \right) = \sum_{i' \neq i} \sum_{j' \neq j} \delta^{(1)}_{i',j'} = \delta^{(1)} * Z_{-m,-n} = \delta^{(1)} * \text{ROT}_{180}(Z^{(1)})$$

CNN Backprop Exercise



$$\frac{\partial E}{\partial X} \otimes K = \frac{\delta E}{\delta Y} \times \text{rot}_1(K) \cdot \frac{\delta E}{\delta Y} \otimes K = \text{rot}_1(X \otimes \frac{\delta E}{\delta Y})$$

$$\frac{\delta E}{\delta X} \otimes K = \frac{\delta E}{\delta Y} \otimes K \quad \text{with 2x zero padding}$$

Pooling Layer: makes the representation smaller and more manageable. $Z^{(l)} = \max_i \{z^{(l-1)}_i\}$

$$\frac{\partial z^{(l)}}{\partial z^{(l-1)}} = \begin{cases} 1 & \text{if } i = \arg\max_i \{z^{(l-1)}_i\} \\ 0 & \text{otherwise.} \end{cases}$$

Deep Networks: many layers: large receptive field despite smaller filters (few parameters)

A deeper network should always perform at least as good as a shallow network. If not \rightarrow failure to optimize

Can be improved with residual connections, reintroducing gradients.

1x1 Convolutions: reduces depth component which leads to fewer parameters, fewer computations. Also increases model capacity.

ResNet: Double convolution with ReLU form Single block: $C_{t+1} = \text{ReLU}(W_h h_t + b_h)$

↳ stack many on top of each other
↳ problem: convolutions at full res are expensive, so often uses down and then upsampling.

Unpooling: nearest neighbor (duplicate), bed of nails (only top left, rest 0), max unpooling (put where max was, else 0)

Transpose Convolution: (also de- & up-convolution, fractionally Strided conv., backward strided conv.)

Output = copies of filter weighted by input, summed where it overlaps. $\text{Conv}: X \cdot K = MX \quad \text{Conv}^T: X \cdot K = M^T X$

Adversarial Defense: Adversarial training (brute force train with lots of adversarial samples to increase robustness), defensive distillation (train with "soft" class probabilities instead of hard labels. Soft-probabilities stem from 2nd model trained on same task with hard labels. Smooths network with respect to attack)

Convolutional layers

kernel size: k, stride padding p, stride s, dilation d, dilation d_x, d_y

Forward: $(H+2p-d(k-1)-1) \cdot S^d + 1$

Backward: $(W+2p-d_s(d_k-1)-1) \cdot S^d + 1$

RECURRENT NEURAL NETWORKS

Dynamical System with inputs: $h^t = f(h^{t-1}, x^t, \theta)$

With f: transition function (same for all timesteps)

Instead of $g^t(x^{t-1}, \dots, x^t)$, has always same input size

RNN Types: One-to-one (vanilla RNN, character-level language model), one-to-many (image captioning), many-to-one (sentiment classification), many-to-many (delayed) (machine translation), many-to-many (immediate) (video classification on a frame level)

Vanilla RNN: $y^t = W_h h^t + b^t, h^t = \tanh(W_u x^t + W_h h^{t-1}) + f(h^{t-1}, x^t, \theta)$

$L^t = \|y^t - y^t\|^2, L = \sum L^t$

Problems with long-term dependencies: (Intuition)

$h^t = W^t h^{t-1} = (W^t)^T h^{t-1}$ (note W to the power of t)

If λ eigenvalue of W: $W = Q \Lambda Q^T$ then

$h^t = ((Q \Lambda Q^T)^T)^t h^0$. If Q is orthogonal: $h^t = Q^T \Lambda^t Q h^0$

↳ if λ gets small \rightarrow vanishes, if λ gets large \rightarrow explodes

Backprop through time:

$$\frac{\partial L}{\partial \theta} = \frac{\partial L^t}{\partial \theta} \frac{\partial \theta}{\partial W}, \frac{\partial \theta}{\partial W} = \sum_{k=t}^T \frac{\partial L^k}{\partial \theta} \frac{\partial \theta}{\partial h^k} \frac{\partial h^k}{\partial W}$$

Note: $\frac{\partial h^t}{\partial h^{t-1}}$ denotes the immediate derivative (keep h^{t-1} fixed)

Vanishing/Exploding Gradient: We look at $\frac{\partial h^t}{\partial h^{t-1}} = \prod_{i=1}^t \frac{\partial h^i}{\partial h^{i-1}} = \prod_{i=1}^t W_{hi} \text{diag}[f'(h^{i-1})]$ of the gradient

Assume: λ_1 is the largest singular value of W_{hi} and $|\text{diag}[f'(h^{i-1})]| < \gamma$ where $\gamma \in \mathbb{R}$ (fine for tanh)

Case 1: It is sufficient $\lambda_1 < \gamma$ for the vanishing gradient problem.

$$V_i: \left| \frac{\partial h^t}{\partial h^{t-1}} \right| \leq \|W_{hi}\| \cdot \left| \text{diag}[f'(h^{i-1})] \right| < \frac{1}{\gamma} \gamma = 1$$

Let $\eta \in \mathbb{R}$ s.t. $V_i: \left| \frac{\partial h^t}{\partial h^{t-1}} \right| \leq \eta < 1$. We can see that

$$\left| \frac{\partial h^t}{\partial h^{t-1}} \right| \leq \eta \rightarrow \text{this goes to zero for } t \rightarrow \infty$$

Case 2: $\lambda_1 > \gamma \rightarrow$ exploding

LSTM

cell state

$$C_{t+1} = \text{tanh}(W_{hc} h_t + W_{cc} c_t + b_c)$$

hidden state

$$h_{t+1} = \text{tanh}(W_{hx} x_t + W_{hh} h_t + b_h)$$

GRU

$$z_{t+1} = \sigma(W_{xz} \cdot [h_t, x_t])$$

$$$\tilde{h}_{t+1} = \tanh(W_{xh} \cdot x_t + W_{zh} \cdot z_{t+1} \cdot h_t)$$$

$$h_{t+1} = (1-\alpha) \tilde{h}_{t+1} + \alpha \cdot h_t$$

↳ Gradient clipping: A way to avoid exploding gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

↳ gradients. Limit the length of the gradient vector

ELBO Derivation: start with data log-likelihood:
 $\log P_{\theta}(x^{(i)}) = \mathbb{E}_{z \sim p_{\theta}(z|x^{(i)})} [\log(p_{\theta}(x^{(i)}|z))] = [p_{\theta}(x^{(i)}) \text{ does not depend on } z]$
 $E_z [\log \frac{p_{\theta}(z|x^{(i)})}{P_{\theta}(z) \cdot p_{\theta}(z|x^{(i)})}] = E_z [\log p_{\theta}(x^{(i)}|z)] - E_z [\log P_{\theta}(z)] + E_z [\log \frac{p_{\theta}(z|x^{(i)})}{p_{\theta}(z)}]$
 $= E_z [\log p_{\theta}(x^{(i)}|z)] - D_{KL}(p_{\theta}(z|x^{(i)}) \| p_{\theta}(z)) + D_{KL}(p_{\theta}(z|x^{(i)}) \| p_{\theta}(z|x^{(i)})) = 0$

Monte Carlo Gradient Estimator

Often encountered: $\nabla_{\theta} \mathbb{E}_{p_{\theta}(z|x)}[f(z)]$ with some differentiable function f . We also have a transformation $z = f(\epsilon, \eta)$.

$\nabla_{\theta} \mathbb{E}_{p_{\theta}(z|x)}[f(z)] = \nabla_{\theta} \int p_{\theta}(z|x) f(z) dz = \nabla_{\theta} \int p(\epsilon) f(z) d\epsilon$
 Note: this works because of the transformation and the differential area under the integral remains the same ($D_{KL}(z) d\epsilon = p(z) d\epsilon$)

$= \nabla_{\theta} \int p(\epsilon) f(g(\epsilon, \eta)) d\epsilon = \nabla_{\theta} \mathbb{E}_{p_{\theta}(z|x)}[f(g(\epsilon, \eta))] = \mathbb{E}_{p_{\theta}(z|x)}[\nabla_{\theta} f(g(\epsilon, \eta))]$

We can then estimate this using: $\frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p_{\theta}(z|x)}[\nabla_{\theta} f(g(\epsilon^{(s)}, \eta))]$

Or in practice draw samples from the noise distribution and average the gradients.

DKL Non-Negativity: $-D_{KL} = - \int_x p(x) \log \frac{p(x)}{q_{\theta}(x)} dx = \int_x p(x) \log \frac{p(x)}{q_{\theta}(x)} dx \leq \log \int_x p(x) \frac{p(x)}{q_{\theta}(x)} dx = \log \int_x q(x) dx = \log(1) = 0$
 Trick: $\mathbb{E}[p(x)] \leq \mathbb{E}[q(x)]$ iff q is concave (\log is)

Reparameterization Trick: (to enable VAE backprop)
 Given Gaussian with μ and σ^2 : $z \sim N(\mu, \sigma^2)$

Assume \exists underlying random variable $\epsilon \sim N(0, 1)$
 $z = \mu + \sigma \epsilon = f(x, \epsilon, \theta) \Rightarrow$ can take derivative w.r.t. M, σ (tells us how changes in μ, σ affects output for fixed ϵ)

Advantage of using N : when modelling $p_{\theta}(z)$ and $q_{\theta}(z|x^{(i)})$ with normal distributions with diagonal covariance matrix, it allows drawing easily from the prior & can calculate DKL analytically
Equation for analytical DKL with N :
 For any $p(z) = N(\mu_p, \Sigma_p)$ and $q(z) = N(\mu_q, \Sigma_q)$
 $\int p(z) \log(q(z)) dz = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \frac{\sigma_{p,j}^2}{\sigma_{q,j}^2} + \frac{1}{2} \sum_{j=1}^J \frac{\mu_{p,j}^2 - \mu_{q,j}^2}{\sigma_{q,j}^2}$

VAE Backprop: $\mathbb{E}_z [\log p_{\theta}(x^{(i)}|z)] - D_{KL}(q_{\theta}(z|x^{(i)}) \| p_{\theta}(z))$
 $= \mathbb{E}_z [\log p_{\theta}(x^{(i)}|z)] + \mathbb{E}_z [\log \frac{p_{\theta}(z)}{q_{\theta}(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_{\theta}(x^{(i)}|z)}{q_{\theta}(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_{\theta}(x^{(i)})}{q_{\theta}(x^{(i)})}]$
 $D_{KL} = \mathbb{E}_z [\log p_{\theta}(x^{(i)})] - \mathbb{E}_z [\log q_{\theta}(x^{(i)})] = \mathbb{E}_z [\log p_{\theta}(x^{(i)})] - \mathbb{E}_z [\log \frac{p_{\theta}(x^{(i)})}{q_{\theta}(x^{(i)})}]$
 $\approx \frac{1}{K} \sum_{k=1}^K \mathbb{E}_z [\log(p_{\theta}(x^{(i)}) / q_{\theta}(x^{(i)}))]$

Limitations of VAEs: Tendency to generate blurry images (believed to be due to injected noise and weak inference models), recently better with inverse autoregressive flow
 β -VAE: learns disentangled representation without supervision. For $\beta=1$ same as VAE
 $\min: L_{\beta} = \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)] + \beta D_{KL}(q_{\theta}(x) \| p_{\theta}(x))$

AUTOREGRESSIVE MODELS

Autoregression: for timeseries, $x_t = b_0 + b_1 x_{t-1} + b_2 x_{t-2}$ taking observations from previous timesteps as input. (e.g. sequence models)

Tabular approach: $p(x) = \prod_{i=1}^T p(x_i|x_{i-1}, \dots, x_1)$ #params 2^{T-1}
 can also specify $p(x_i|x_{i-1}, \dots, x_1) = \text{Ber}(f(x_{i-1}, \dots, x_1))$

Fully visible belief networks: $x_i = p(x_i = 1 | x_{i-1}, \dots, x_1)$ modeled via logistic regression: $f(x_{i-1}, \dots, x_1) = \sigma(a_0 + \dots + a_{i-1} x_{i-1})$, has $O(n^2)$ parameters

Neural Autoregressive Density Estimator (NADE): autoencoder-like network to learn $p(x_i = 1 | x_{i-1}, \dots, x_1)$
 $K = XW_i \in \mathbb{R}^{T \times D}$ $V = XW_i \in \mathbb{R}^{T \times D}$ $Q = X^T W_i \in \mathbb{R}^{D \times T}$
 $h_i = \sigma(b + W_i \cdot \text{softmax}(X^T h_{i-1}))$ for binary data (first i cols)
 (we can use that $(b + W_i \cdot \text{softmax}(X^T h_{i-1})) - (b + W_i \cdot \text{softmax}(X^T h_{i-1}) - W_i \cdot x_{i-1}) = W_i \cdot x_{i-1}$)
 Train by max. $\frac{1}{T} \sum_{i=1}^T \log(p(x_i)) = \frac{1}{T} \sum_{i=1}^T \log(p(x_i|x_{i-1}, \dots, x_1))$
 Advantages: efficient computations in O(TD)

could make use of second order optimizers, easily extendable to other types of observations (reals, multinomials). Random input order is fine.
 Training: take data from training, Inference: sampling
Extensions: RNADE: conditionals modelled by mixture of Gaussians, DeepNADE: a DNN trained to assign a cond. distr. to any given var. given any subset of others (order-less), ConvNADE

Masked Autoencoder Distribution Estimator (MADE)
 Constrain autoencoder s.t. output can be used as conditional $P(x_i | x_{i-1})$. To fulfill autoregressive property: no computational path output unit x_d and any of the input units x_1, \dots, x_0 must exist. Can also be seen as masks $m_i^{(i)}$: if unit i of layer i is connected to unit j of layer $i-1$. Train with NLL of binary x . Computing $p(x)$ is just forward pass, Sampling requires D forward passes. In practice very large hidden layers necessary.

Generative Models for Natural Images: Use an explicit density model & use chain rule to decompose likelihood of image x into prod. of 1D dist. $p(x) = \prod_{i=1}^T p(x_i | x_{i-1}) \rightarrow$ need to define ordering of pixels, complex distributions \rightarrow paramet. via NN. Then max. LL.

PixelRNN: Generate pixel starting from the corner and dependency to prev. pixels modelled using LSTM
 Issue: sequential gen's slow due to explicit pixel dependencies

PixelCNN: Still generate starting from corner, but model the dependency as CNN: (masked convolutions)
 $p(x_i | x_{i-1}) = p(x_i | x_{i-1}, x_{i-2}, \dots, x_1) \cdot p(x_i | x_{i-1}, x_{i-2}, \dots, x_1) \cdot \dots$
 stacks can be used to speed it up. Add information between vertical and horizontal stack (preserve pixel dependencies)

$h_{k+1} = \tanh(W_{k,k} * h_k) \odot (W_{k,k+1} * h_k)$. Training is faster than PixelRNN (can parallelize convolutions) but still slow (sequential generation).

Advantages: explicit likelihood, likelihood of training data is good eval. metric, good samples.

Dilated convolution: Receptive field size can increase exponentially with dilated layers.

WaveNet: Temporal convolutional networks uses dilated convolution. Cannot use strided conv, due to need to preserve resolution

Variational RNN (VRNN): RNNs are autoregressive models increase expressive power of RNNs by incorporating stochastic latent variables into hidden state of RNN. Include one VAE per timestep. Including a dynamic prior explicitly models temporal dependencies between time steps. Can include conditional (e.g. style and content of handwriting).

Stochastic Temporal Convolutional Networks (STCN)
 Combines computational advantages of TCNs with expressiveness of stochastic RNNs: integrate a hierarchical VAE with a TCN.

Self Attention and Transformers
 Transforms input Embeddings to keys, values, queries
 $K = XW_k \in \mathbb{R}^{T \times D}$ $V = XW_v \in \mathbb{R}^{T \times D}$ $Q = X^T W_q \in \mathbb{R}^{D \times T}$
 Attention weights: $\alpha = \text{softmax}(QK^T / \sqrt{D} + M) \in \mathbb{R}^{T \times T}$
 where $M = \begin{pmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$ masks out future influences
 Prediction: $x_t = \alpha V$

NORMALIZING FLOWS

Have analytic model density $p_{\theta}(x)$ and latent variable space

$$z \xrightarrow{f_{\theta}} x \quad x = f_{\theta}(z) \quad z = f_{\theta}^{-1}(x)$$

where f_{θ} is deterministic & invertible, preserves

dimensionality, Jacobian is efficient

Change of variables

Maps one distribution to another

For a transformation $f'(x) = x$

$$p_x(x) = p_z(f^{-1}(x)) \mid \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right)$$

$$= p_z(f^{-1}(x)) \mid \det \left(\frac{\partial f(x)}{\partial x} \right)$$

Coupling Layer

$$\begin{array}{c} x^a \xrightarrow{h} y^a \\ x^b \xrightarrow{h} y^b \end{array}$$

forward: $(y^a, y^b) = (h(x^a, h(x^b)))$
 backward: $(x^a, x^b) = (h^{-1}(y^a, h(y^b)))$
 $\det(h(x^a, x^b)) = \begin{vmatrix} h^a & h^b \\ x^a & x^b \end{vmatrix}$ h can be arbitrary
 h is any elementwise function

We get complex transforms via composition

$$x = f(z) = f_0 \circ f_1 \circ \dots \circ f_T$$
 (training) (each f_i is invertible)

$$\text{now } \log p(x) = \sum_{i=0}^T (\log p_{\theta}(f_i(x))) + \log |\det(\frac{\partial f_i(x)}{\partial x})|$$

$$\text{Inference: } p_{\theta}(x) = p_z(f^{-1}(x)) \mid \det(\frac{\partial f^{-1}(x)}{\partial x})$$

\Rightarrow Also conveniently via ODE: $x = f_{\theta}(x(t), t)$

Multi-Scale Architecture

Steps: squeeze, split, step of flow

By reducing the spatial resolution

flow step: θ Activation norm

like BatchNorm: $q = \text{scale} \cdot x + \text{shift}$

② invertible 1x1 conv: h is a HWx1 tensor, $W = PLU$ ($U = \text{diag}(s)$) rand. rotation

$$\log |\det(h_{\theta}(h))| = H \cdot W \cdot \log |\det(U)| = \sum \log(s)$$

③ (conditional) Coupling layer

Applications in Computer Vision

Super Resolution: Normalizing flow learns a dist. of HR

U-Net: (low resolution image as conditioning)

Disentanglement: Conditions on learned attributes

C-Flow: Use two flows, condition one on the other.

Each flow handles one representation, conditioning gives transfer

Training is still slow

GENERATIVE ADVERSERIAL NETWORKS

GAN Intuition: Want to sample high dimensional data points from true data distribution (intractable), so draw from simple distribution and use a neural network to transform into output

Use generator $G: \mathbb{R}^n \rightarrow X$ and discriminator $D: X \rightarrow [0, 1]$

GAN Objective: Consider only the discriminator
 $L(D) = -\frac{1}{N} \sum_{i=1}^N (y_i \cdot \log(D(x^{(i)})) + (1-y_i) \cdot \log(1-D(x^{(i)})))$
 for GANs we know 50% are $y=1$ and 50% are $y=0$. Our goal is to train D & G jointly:
 $= \frac{1}{2} (\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_{\theta}} [\log(1-D(G(z)))] + \mathbb{E}_{x \sim p_{\theta}} [\log(1-D(x))])$
 \Rightarrow a good G maximizes loss function of D. Value function: $V(G, D) = \mathbb{E}_{x \sim p_{\text{data}}} [\log(D(x))] + \mathbb{E}_{z \sim p_{\theta}} [\log(1-D(G(z)))]$
 $\Rightarrow G$ will fool any D: $G^* = \arg \max_{G} V(G, D)$

Optimum of $V(G, D)$ iff $P_d = P_g$: Let $D_g^* = \arg \max_{D} V(G, D)$ then the objective becomes $G^* = \arg \min_{G} V(G, D_g^*)$.

If $P_d = P_g$ then $D_g^* = \frac{P_d}{P_d + P_g}$. Proof: $V(G, D) = \mathbb{E}_{x \sim p_{\text{data}}} [\log(D(x))] + \int_{x \sim p_{\theta}} \log(1-D(G(z))) dz$

Optimal D: $\forall a, b \in \{0, 1\}$ $f(y) = a \log(y) + b \log(1-y)$
 $f'(x) = 0 \Leftrightarrow \frac{a}{y} - \frac{b}{1-y} = 0 \Rightarrow y = \frac{a}{a+b}$. $f''(x) = -\frac{a+b}{y(1-y)^2} < 0 \Rightarrow D_g^* = \frac{P_d}{P_d + P_g}$

Obj. min JS: $V(D_g^*, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log \frac{P_d}{P_d + P_g}] + \mathbb{E}_{z \sim p_{\theta}} [\log \frac{P_g}{P_d + P_g}]$
 Global min is achieved if $P_d = P_g$ and at $\min V(\cdot) = -\log 4$. Proof: $\mathbb{E}_{x \sim p_{\text{data}}} [\log \frac{P_d}{P_d + P_g}] + \mathbb{E}_{z \sim p_{\theta}} [\log \frac{P_g}{P_d + P_g}] = -\log(2) + \mathbb{E}[\log(\frac{P_d}{P_d + P_g})] - \log(2) + \mathbb{E}[\log(\frac{P_g}{P_d + P_g})] = -\log(4) + D_{KL}(P_d || \frac{P_d + P_g}{2}) + D_{KL}(P_g || \frac{P_d + P_g}{2})$ & iff $P_d = P_g$ $= -\log(4) + 2 \cdot D_{JS}(P_d || P_g)$. Note: last term $\gg 0$ then = 0

Jensen-Shannon Divergence Symmetric DKE
 $D_{JS}(p||q) = \frac{1}{2} D_{KL}(p || \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q || \frac{p+q}{2})$

Obj: $\min_{\theta} \max_{\theta_D} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_{\theta}(x)] + \mathbb{E}_{z \sim p_{\theta}} [\log(1-D_{\theta}(G(z)))]$

GAN Training: $\min_{\theta_D} \max_{\theta_G} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_{\theta_D}(x)] + \mathbb{E}_{z \sim p_{\theta_G}} [\log(1-D_{\theta_D}(G_{\theta_G}(z)))]$

Alternate: 1) $\max_{\theta_G} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_{\theta_D}(x)] + \mathbb{E}_{z \sim p_{\theta_G}} [\log(1-D_{\theta_D}(G_{\theta_G}(z)))]$ (both GD)
 2) $\min_{\theta_D} \mathbb{E}_{z \sim p_{\theta_G}} [\log(1-D_{\theta_D}(G_{\theta_G}(z)))]$

GAN Training Problems: In this problematic region G produces bad samples and gradient is very flat. Can lead to overtraining of the discriminator (G cannot catch up). Solution: In 2. gradient ascend: $\max_{\theta_G} \mathbb{E}_{z \sim p_{\theta_G}} [\log(D_{\theta_D}(G_{\theta_G}(z)))]$. Generally: need to find Nash-equilibrium in two-player game. This leads to mode collapse (no sample diversity). happens if there is a saddle in dual energy landscape (the generator learns only variations of a single sample).

Issues with JS divergence: correlates badly with sample quality (do not know when to stop training).

Sample quality (do not know when to stop training). D tends to be optimal, JSD saturates \rightarrow bad gradients
 Solution: use Wasserstein distance

GAN Pseudocode

for #training iterations do
 for k steps do

sample minibatch of m noise samples $\{z \sim p_{\theta}(z)\}$
 sample minibatch of examples $\{x\}$ from data gen. dist.
 update discriminator by ascending on its stochastic grad:
 $\nabla_{\theta_D} \frac{1}{m} \sum_{i=1}^m [\log(D(x^{(i)}) + \log(1-D(G(z^{(i)})))]$

sample minibatch of m noise samples $\{z \sim p_{\theta}(z)\}$
 update the generator (ascent): $\nabla_{\theta_G} \frac{1}{m} \sum_{i=1}^m \log(D(G(z^{(i)}))$

GAN: Comparison to VAE: No variational bound is needed, only samples. Generally sharper images. GANs are asymptotically consistent (not proven for VAE)

DC-GAN: G is an up-sampling network with de-convolutions, D is a CNN. Replace pooling by strided convolutions in D and strided deconvolutions in G. Use batch norm for both. Remove fully connected hidden layers, use ReLU in G except tanh for output, use LeakyReLU in D for all layers. Can do "vector math": smiling φ - neutral φ + neutral θ^{φ} = smiling θ^{φ}
 Image translation

Pix2Pix: $\mathcal{L}(G, D) = \mathcal{L}_{\text{GAN}}(G, D) + \lambda \mathcal{L}_{\text{L1}}(G)$ where $\mathcal{L}_{\text{GAN}}(G, D) = \mathbb{E}_{x,y} [\log D(x, y)] + \mathbb{E}_{x,z} [1 - \log D(x, G(x, z))]$ and $\mathcal{L}_{\text{L1}}(G) = \mathbb{E}_{x,y,z} [\|y - G(x, z)\|_1]$

Limitations: requires pairs of images as training data +
circular image rotation
Cycle GAN: $L_{GAN}(G, D_x, D_y) = L_{GAN}(G, D_y, X, Y) + L_{GAN}(F, D_x, Y, X) + \lambda L_{Cycle}(G, F)$ where
 $L_{cycle} = \mathbb{E}_{x \sim \text{rand}}[\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim \text{rand}}[\|G(F(y)) - y\|_1]$

Progressive Growing of GANs: Grow G, D resolution by adding layers during training

Style GAN: ① Learn an intermediate latent space ② Inject latent space (style) and noise in every layer of G

Diffusion Models: Gradually add Gaussian Noise and learn to reverse it:

BODY MODELLING

Deep Features: • Directly regress on (x,y)
Coordinate nodes + refinement • Estimate Headness for each joint (correlation with Conv/Net)

Thin-Slicing Networks: ① predict Joint Headmap ② Use flow warp to track joint location estimates over time ③ Minimize spring Energy in tracked body graph over thin time slice

Linear Blend Skinning: w_i : Blend skinning weights, $t_i = \sum_i w_i c_i(G, J) t_i$, c_i : Rigid Bone transformation

J : Joint locations \ominus Unrealistic Deformations

SMPL Body Model: B : Body Shape params, $\epsilon_i = \sum_i w_i c_i(\theta, J(B))(t_i + s_i(B) + p_i(\theta))$ (off of base mesh and scaled template)

$s_i(B)$: shape correctives: β learned from

$p_i(\theta)$: pose correctives: β template poses

④ Joints $J(B)$ depend on shape B ④ pose / shape correctives

Optimization Based Fitting: Directly optimize θ^* : $\arg\min \| \text{SMPL} - 20 \text{joints} \|$ ④ slow convergence $\theta^{t+1} = \theta^t - \lambda \left(\frac{\delta \text{loss}_{\text{regress}}}{\delta \theta} + \frac{\delta \text{loss}_{\text{shape}}}{\delta \theta} \right)$ ④ sensitive to init of G, B

Learned Gradient Descent: Fnn: Neural Net

$\theta^{t+1} = \theta^t + F_{\text{grad}} \left(\frac{\delta \text{loss}_{\text{regress}}}{\delta \theta}, \theta^t, x \right) \frac{\delta \text{loss}_{\text{shape}}}{\delta \theta}$: Actual gradients

θ^t : current state x ; Target pose; F_{grad} is learned by simulating 2D poses with Ground truth

Challenges: ① Self-Occlusions ② Ambig-
uous Depth Information ③ Articulated Motion
④ Non-rigid Deformations

Template- & Regression Based Capture: Template Mesh provides Data for challenging Poses.

NEURAL IMPLICIT SURFACES

Coordinate MCPs: Takes fixed size positional Encoding as input, RGB output
3D Representations: Voxels \ominus Cubic in memory
Points \ominus Does not model connectivity
Meshes \ominus Difficult to work with

Implicit Differentiation: Take Derivative on both sides of Equation, then substitute

Neural Implicit Representations: Occupancy Networks $f_o: \mathbb{R}^2 \times \mathcal{X} \rightarrow [0, 1]$ occupancy prob.

Signed Distance Fields: $f_d: \mathbb{R}^2 \times \mathcal{X} \rightarrow \mathbb{R}$ signed distance

Learning 3D Implicit Shapes From Waterfall meshes: Randomly sample 3D points in space

$g(O, V) = \sum_i^k \text{SCE}(f_d(p_i, z_i), O_i) \leftarrow \text{Cross entropy Loss}$

From Point Clouds: A Sparse supervision only on the surface. Ekard POE: $\|\nabla f(x)\|_1 = 1; f(x) = 0 \forall x \in \mathcal{S}$

\rightarrow solution gives distance from \mathcal{S}
 $L(\theta) = \sum_i \|f_\theta(x_i)\|^2 + \lambda \mathbb{E}_x (\|f_\theta(x)\| - 1)^2$

④ Converges to reproducing plane w/prob 1 for model

Differentiable rendering Forward Pass

① Find surface root \hat{p} along ray w ② Evaluate texture field $t_o(\hat{p})$ **Second Method**

$x_i = x_i - f(x_i) \frac{x_i - x_o}{\|x_i - x_o\|}$ x_i, x_o need to be on opposite sides of the root

Backward Pass: $R(I, I') = \sum_u \|I_u - I'_u\|$

$\frac{\delta L}{\delta \theta} = \sum_u \frac{\delta L}{\delta I_u} \frac{\delta I_u}{\delta \theta} \in \mathbb{R}^d$ depends on intersection point $\hat{p} = t_o(\hat{p})$ and texture value

$\frac{\delta I_u}{\delta \theta} = \frac{\delta I_u(\theta)}{\delta \theta} + \frac{\delta t_o(\theta)}{\delta \theta} \cdot \frac{\delta I_u(t_o(\theta))}{\delta t_o(\theta)} = \frac{\delta I_u(\theta)}{\delta \theta} + \omega \frac{\delta t_o(\theta)}{\delta \theta} \cdot \frac{\delta I_u(t_o(\theta))}{\delta t_o(\theta)}$

Neural Radiance Fields: $F_o: (x, y, z, t) \rightarrow (R, G, B, \sigma)$

④ Can model transparency by predicting density

④ Conditioning on view (θ, t) can model view dep. effects

④ Only add view dir (θ, t) to keep geometry consistent across views

Volume Rendering: $\pi = 1 - e^{-\int_{t_i}^{t_f} \sigma(s) ds}$ $\sigma = \text{sum}_{i=1}^n \sigma_i$

Transmittance: $T_i = \prod_{j=i}^{t_f} (1 - \sigma_j)$ Color $c = \sum_{i=1}^n T_i \sigma_i$ output

Training: $R(\theta) = \sum_i \| \text{render}(F_o) - I_i \|^2 \leftarrow \text{RGB Loss}$

Faster with Hierarchical Sampling: $C = \sum_i T_i \sigma_i$

④ Can model transparency & skin structure

④ Generally worse geometry than implicit surfaces

Positional Encoding (Fourier Features)

MCPs prefer smooth functions \rightarrow Bad frequency Details

Pass coordinates through random Fourier Features

$v \rightarrow \boxed{\boxed{\boxed{\boxed{\quad}}}}$ $\rightarrow y \Rightarrow \boxed{\boxed{\boxed{\boxed{\quad}}}} \rightarrow \boxed{\boxed{\boxed{\boxed{\quad}}}} \rightarrow y$ Can model a higher frequency function

CLASSICAL REINFORCEMENT LEARNING

Reinforcement Learning: Uses: games, robot control, attention based models, logistics & operations

Goal: achieve high cumulative reward over time.

Value function of policy π , V_π , expresses expected

cumulative reward for each state when following π .

$V_\pi(s) = \mathbb{E}_{G \sim \pi}[G_t | s_t = s]$ $= \mathbb{E}_{G \sim \pi}[r_t + \gamma \mathbb{E}_{G \sim \pi}[s_{t+1}, r_{t+1}, \dots, G_{T-1}, r_{T-1}, G_T | s_t = s]]$

$= \sum_a \pi(a, s) \sum_s \sum_r P(s', r | s, a) [r + \gamma V_\pi(s')]$ where $G_t = \sum_{k=t}^T \gamma^k R_{k+1}$

MDP: states S , actions A , reward function r , transition function P , initial state s_0 , discount factor γ

Value Iteration: Given a value function V_π , we can derive a greedy policy π^* :

$\pi^*(s) = \arg\max_{a \in A} \{r(s, a) + \gamma V_\pi(P(s, a))\}$. One can show $V_\pi(s) \geq V^*(s)$. The optimal policy π^* is greedy wr.t. V^* . V^* fulfills the Bellman Optimality Equation:

$V^*(s) = \max_{a \in A} \{r(s, a) + \gamma V^*(P(s, a))\}$. Can use step cost instead of γ . Update backwards each

time doing the action that gives the highest reward until nothing changes anymore.

Theorem: Value iteration

while $\Delta > 0$ converges. At convergence found the optimal V^* . Need to find the optimil (stationary & greedy) policy.

$\Delta \leftarrow 0$ for all $s \in S$

$V(s) \leftarrow V(s)$

$V(s) \leftarrow \max_{a \in A} \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |V - V(s)|)$

$\pi^* \leftarrow \text{greedy policy wr.t. } V$

Advantages of DP: Exact methods, policy/value iteration guaranteed to converge in finite number of iterations, value iteration typically more efficient easy implementation

Disadvantages of DP: need to know transition probability matrix, need to iterate over the whole state space (very expensive), requires memory proportional to size of state space, generally computationally expensive

Problem with large state space: we will never visit all states, but only a small subset. We cannot estimate the values of these states.

Monte Carlo Methods: value function is $E[\cdot]$:

$V_\pi(s) = \mathbb{E}_{G \sim \pi}[r(s, a_1)] = \sum_{t=0}^T r(s_t, a_t)$, so we can just run policies with the given policy and compute samples of: $\sum_{t=0}^T r(s_t, a_t)$

follow policy, update values, update policy, repeat

Advantages: Unbiased estimate, do not need to know system dynamics. **Disadvantages**: High variance, exploration/exploitation dilemma, need termination state (slow for long episodes)

Temporal Difference Learning: For each step with action a from s to s' that we take with our policy, we compute the difference to our current estimate and update our value function: $\Delta V(s) = r(s, a) + \gamma V(s') - V(s)$, $V(s) \leftarrow V(s) + \alpha \Delta V(s)$, $a > 0$

SARSA: On policy: Computes the Q-Value according to a policy and then the agent follows that policy.

Q-Learning: Off policy: Computes the Q-Value according to a greedy policy, but the agent follows a different exploration policy:

$\Delta Q(s, a) = R_{t+1} + \gamma \max_a \{Q(s', a)\} - Q(s, a)$, $Q(s, a) \leftarrow Q(s, a) + \alpha \Delta Q(s, a)$ off policy since π for updating Q different from π' that we use for data collection.

Advantages: Less variance than Monte Carlo sampling due to bootstrapping. More sample efficient, do not need to know the transition probability matrix. **Disadvantages**: biased due to bootstrapping, exploration/exploitation dilemma can behave poorly in stochastic environments. requires exploring full state-action space (does not work for continuous action space)

TAXONOMY OF RL METHODS



DEEP REINFORCEMENT LEARNING

Policy gradients: $\pi(a_t | s_t) = \mathcal{N}(\mu_t, \sigma^2 | s_t)$. We can collect rollout trajectories: $\pi(\tau) = \pi(s_1, a_1, \dots, s_t, a_t) = \prod_{i=1}^t \pi(s_i | a_i)$ then we update the policy: good traj. more & bad traj. less likely.

Policy gradient: goal: $\theta^* = \arg\max J(\theta)$ by updating $\theta \leftarrow \theta + \nabla_\theta J(\theta)$ with gradient ascent.

$J(\theta) = \mathbb{E}_{\tau \sim \pi(\theta)} [\sum_t r^t(s_t, a_t)] = \mathbb{E}_{\tau \sim \pi(\theta)} [\log \pi_\theta(\tau)] = \int_{\tau \sim \pi(\theta)} \log \pi_\theta(\tau) d\tau = \int_{\tau \sim \pi(\theta)} \frac{\partial \log \pi_\theta(\tau)}{\partial \theta} d\tau = \frac{\partial \log \pi_\theta(\tau)}{\partial \theta} = \frac{\partial \log \pi_\theta(s_t, a_t)}{\partial \theta} = \log p(s_t, a_t, \dots, s_1, a_1) = \log p(s_1, a_1, \dots, s_t, a_t) \prod_{i=1}^t \log \pi_\theta(s_i, a_i)$

Valid even if reward is discontinuous or unknown Sample space of path is a discrete set.

Gradient tries to increase probability with pos. r and decrease probability with small/neg. r.

Reinforce: sample trajectories τ by rolling out the policy, i.e. sampling a random action from the current policy until the end of the episode. Monte Carlo type of RL method that evaluates full trajectories (episodes) to update policy

$\pi_\theta(\theta) \approx \frac{1}{N} \sum_i^N \sum_{t=0}^T \log \pi_\theta(a_t | s_t)$, $\theta \leftarrow \theta + \eta \nabla_\theta J(\theta)$ (unbiased estimate)

Variance problem: this produces noisy gradients: If all cumulative rewards are positive, we make all trajectories more likely. If all traj. rewards are 0, we might decrease probability of good trajectories.

Baseline: To reduce variance, we can introduce a baseline $b(s_t)$ (produces okay trajectories, actually any function that does not depend on a_t , e.g. average reward or estimate of state-value function) $r(t) = \sum_{t=0}^T \gamma^t r(s_t, a_t) - b(s_t)$. Still unbiased!

Actor-critic: Use bootstrapping for reward estimation: $\nabla_\theta J(\theta) = \frac{1}{N} \sum_i^N \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) (r(s_t, a_t) + \gamma V(s_{t+1}) - V(s_t))$ Both π and V are represented by neural nets.

|

|