

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) $A^{-1} = A^T$, $AA^T = A^TA = I$

$\det(A) \in \{-1, 1\}$, $\det(A^TA) = 1$. preserves: inner product norm, distance, angle, rank, matrix orthogonality

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

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

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

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

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

Probability/ Statistics Basics

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

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

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

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

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: $\log P_\theta(x^{(i)}) = E_{z \sim q_\theta(z|x^{(i)})} [\log P_\theta(x^{(i)})]$

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 $\forall x f''(x) \leq 0$)

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 $t \cdot f(\hat{x}) + (1-t)f^*(x) < f(\hat{x})$. So we get $f(t\hat{x} + (1-t)x^*) \geq f(\hat{x}) > t \cdot f(\hat{x}) + (1-t)f^*(x)$ 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)

$$\frac{\partial Y}{\partial X} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix} \quad \frac{\partial X}{\partial Y} = \begin{bmatrix} \frac{\partial x}{\partial y_1} \dots \frac{\partial x}{\partial y_m} \end{bmatrix} \quad \frac{\partial Y}{\partial Z} = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} \dots \frac{\partial y_1}{\partial z_n} \\ \vdots \\ \frac{\partial y_m}{\partial z_1} \dots \frac{\partial y_m}{\partial z_n} \end{bmatrix}$$

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

1D Derivation Rules

$$\frac{\partial(f \cdot g)}{\partial x} = f \cdot g' + g \cdot f'$$

$$\frac{\partial f}{\partial g} = \frac{f'g - g'f}{g^2}$$

$$\frac{\partial f}{\partial x} = f'(g(x))g'(x)$$

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

$$(f^g)' = f^g \left(f' \frac{g}{f} + g' \ln(f) \right)$$

Multivariate Chain Rule

For function $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$, $x \mapsto y$ and $f: \mathbb{R}^m \rightarrow \mathbb{R}$, $y \mapsto z$ we have: $\frac{\partial z_k}{\partial x_j} = \sum_{i=1}^m \frac{\partial z_k}{\partial y_i} \frac{\partial y_i}{\partial x_j}$

$$\frac{\partial z_k}{\partial y} = \frac{\partial z_k}{\partial y} \frac{\partial y}{\partial x}$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$$\frac{\partial f(x), y(t)}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$$

x, y, t can be scalar, vector or matrix.

Derivatives:

$$\frac{\partial(b^T x)}{\partial x} = \frac{\partial}{\partial x}(x^T b) = b \quad \frac{\partial(x^T x)}{\partial x} = 2x \quad \frac{\partial(x^T Ax)}{\partial x} = (A^T - A)x$$

$$\frac{\partial(b^T Ax)}{\partial x} = A^T b \quad \frac{\partial}{\partial x}(c^T x b) = cb^T \quad \frac{\partial}{\partial x}(c^T x^T b) = bc^T$$

$$\frac{\partial(\|x - b\|_2)}{\partial x} = \frac{x - b}{\|x - b\|_2} \quad \frac{\partial(\|x\|_2^2)}{\partial x} = 2x \quad \frac{\partial(\|x\|_F^2)}{\partial x} = 2x$$

$$\frac{\partial \log(x)}{\partial x} = \frac{1}{x} \quad \frac{\partial}{\partial x} \frac{1}{f(x)} = \frac{-f'(x)}{(f(x))^2} \quad \frac{\partial}{\partial x} |x| = (-1)^{|x|} \mathbb{1}_{x>0}$$

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

Challenges (with images): Viewpoint variations, illumination, deformation, occlusion, background clutter, intra-class variation.

Basic linear classifier: $f(x, w) = w^T x + b$

Treat pixels as a high-dimensional space, let the classifier draw hyperplanes through that space that separate the classes. Inner product between image and weights to compute a similarity score.

Multiclass SVM loss: $L^i = \sum_{j \neq y^{(i)}} \max(0, S_j - S_{y^{(i)}} + 1)$

What if we drop $j \neq y^{(i)}$: scaled a bit, can also use mean or square

Regularization: $L1: R(w) = \sum_k \sum_l |w_{k,l}|$, $L2: R(w) = \sum_k \sum_l \beta w_{k,l}^2 + |w_{k,l}|$, MaxNorm, Dropout,...

Elastic Net $R(w) = \sum_k \sum_l \beta w_{k,l}^2 + \lambda |w_{k,l}|$, MaxNorm, Dropout,...

Image Features: color histograms (hue), HOG, SIFT features, bag of words (for patches)

Traditional ML: hand-crafted features and kernel, then simple classifier (vs lots of features & learnable kernel)

Perceptron learning: Init $w^{(0)} = 0$, then while $\exists y$ that is misclassified do $\{w^{k+1} = w^k + \eta(y - \hat{y})x^{(i)} \text{ where } \hat{y} = (\omega^T x^{(i)}) > 0\}$ 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 $D = \{(x^{(1)}, y^{(1)}, \dots, x^{(n)}, y^{(n)})\}$

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

MLE: $\omega_{MLE} = \arg \max_w P(D|w) = \arg \max_w \prod_{i=1}^n P(y^{(i)}|x^{(i)}, w)$

$= \arg \max_w \prod_{i=1}^n \left[\sigma(w^T x^{(i)})^{y^{(i)}} (1 - \sigma(w^T x^{(i)}))^{1-y^{(i)}} \right]$

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

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

Loss Functions

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

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

Negative log-likelihood: $-\sum_k y_k \log \hat{y}_k$

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 $|g(x) - f(x)| < \epsilon \forall x \in C^\infty(\mathbb{R})$ (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: $X_s = (x - \mu)/\sigma$ 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!

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 $M(T(I)) = T(M(I))$

Early Stopping: Save parameter where the lowest Validation loss occurs then stop if there is no improvement for p (patience) epochs. Then retrain model on entire training set for final test error (either for the same number of epochs (easiest & decent performance) or continue training).

As regularizer: Assuming $\|\nabla C(\theta)\| \leq K$ bounded, learning rate ϵ , initial parameters θ_0 , optimal number of training iterations i^* and param. θ^* learned in i^* steps: $\|\theta_0 - \theta^*\|_2 \leq K\epsilon i^*$.

Pseudocode early stopping

$i, j, i^* \leftarrow 0; V \leftarrow \infty; \theta, \theta^* \leftarrow \theta_0$. Initialize

While $j < p$ do: p : patience (how long to wait)

$\theta \leftarrow \text{train}(\theta, n)$; $i \leftarrow i + n$ train n steps

$V_c \leftarrow \text{validate}(\theta)$

if $V_c < V$ then

$j \leftarrow 0; \theta^* \leftarrow \theta; i^* \leftarrow i; V \leftarrow V_c$

else

$j \leftarrow j + 1$

return θ^*, i^* optimal param. w.r.t. validation set

Dropout: Let $r_i \sim \text{Bernoulli}(p)$, for $p \in [0, 1]$

$\$

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

Why adaptive LR: magnitude of gradients can be very different for different layers. Gradients can get very small in the early layers of very deep networks. If a hidden unit has a big fan-in, small changes on many of its incoming weights can cause the learning to overshoot. So we adapt the learning rate to the parameters by performing larger updates for infrequent and smaller update to frequent parameters.

AdaGrad: $\theta_{t+1,i} = \theta_{t,i} - \eta_i / \sqrt{G_{t,i}} + \epsilon \cdot \nabla_{\theta} J(\theta_i)$ where $G_i \in \mathbb{R}^{d \times d}$ is a diagonal matrix where each diagonal element i,i is the sum of the squares of the gradients w.r.t. θ_i up to time step t . + Smaller update for larger gradients: avoids overshooting - accumulation of the squared gradients in the denominator keeps growing.

RMS prop/Adadelta: $\theta_{t+1} = \theta_t - \eta / \sqrt{E[g^2]_{t-1}} + \epsilon \cdot \nabla_{\theta} J(\theta_t)$, $E[g^2]_t = \gamma E[g^2]_{t-1} + (1-\gamma) g_t^2$. Attempts to reduce AdaGrad's aggressive, monotonically decreasing learning rate. The sum of squared gradients is recursively defined as a decaying average of all past squared gradients.

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

$\hat{m}_t = m_t / (1-B_1)$, $\hat{v}_t = v_t / (1-B_2)$, $\hat{\theta}_{t+n} = \theta_t - \eta / \sqrt{\hat{v}_t} + \epsilon \hat{m}_t$

Learning rate: If too large, can 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: $\eta_{t+1} = \eta_0 / (1+kt)$

Backpropagation:

$$\delta^{(l-1)} = \frac{\partial C}{\partial z^{(l-1)}} = \frac{\partial C}{\partial z^{(l)}} \frac{\partial z^{(l)}}{\partial z^{(l-1)}} = \sum_j \delta_j^{(l)} \frac{\partial z^{(l)}}{\partial z_j^{(l-1)}}$$

$$z^{(l)} = \sigma(\theta^{(l)} z^{(l-1)})$$

$$\frac{\partial C}{\partial \theta^{(l)}} = \sum_j \frac{\partial C}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial \theta^{(l)}} = \sum_j \delta_j^{(l)} \frac{\partial z_j^{(l)}}{\partial \theta^{(l)}}$$

Receptive fields: particular part of the body surface in which a stimulus will trigger firing in sensory neurons. (everything connected to single input unit form receptive field of that unit)

Receptive Fields Example:

Kernel Size

3	$1 \rightarrow (3, 1) \rightarrow (5, 1) \rightarrow (7, 1) \rightarrow (9, 1)$
3	$1 \rightarrow (3, 2) \rightarrow (7, 1) \rightarrow (11, 1) \rightarrow (15, 1)$
3	$1 \rightarrow (3, 2) \rightarrow (7, 2) \rightarrow (15, 1) \rightarrow (23, 1)$

Backprop in MLP

We have $F = F^L \circ F^{L-1} \circ \dots \circ F^1$ and $F(x) = z^L$, $z_i^L = \tanh(w_i^L z^{L-1} + b_i^L)$

We define $\delta^L := \frac{\partial C}{\partial z^L}$ which is:

$$\delta^L = \frac{\partial C}{\partial z^L} = \frac{\partial}{\partial z^L} (R \circ F^{L-1} \circ \dots \circ F^1(z^L))$$

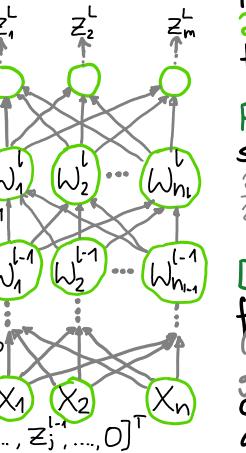
$$= \delta^L \cdot \frac{\partial z^L}{\partial z^{L-1}} \dots \frac{\partial z^L}{\partial z^1} = \delta^L \cdot \frac{\partial z^L}{\partial z^1}$$

We can also calculate:

$$\frac{\partial C}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} (R \circ F^{L-1} \circ \dots \circ F^1(z^L)) = \frac{\partial C}{\partial z^L} \frac{\partial z^L}{\partial w_{ij}}$$

$$\text{Let } z^L = \tanh(w^L z^{L-1} + b^L) = \tanh(a) \quad a = \tanh(w^L z^{L-1} + b^L)$$

$$\frac{\partial z^L}{\partial w_{ij}} = \frac{\partial \tanh(a)}{\partial a} \frac{\partial a}{\partial w_{ij}} = \text{diag}(1 - \tanh^2(a)) \cdot [0, \dots, z_j^{L-1}, \dots, 0]^T$$



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

CONVOLUTIONAL NEURAL NETWORKS

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

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

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

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

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

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

Called shift-invariant if K does not depend on (i,j)

Same for $K(i,j) = K(-i,-j)$ kernels.

and kernels: matrix-vector multipl. $(I * K) = \begin{bmatrix} k_1 & 0 & \dots & 0 \\ k_2 & k_1 & \dots & 0 \\ k_3 & k_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ k_m & k_{m-1} & \dots & k_1 \end{bmatrix} \cdot \begin{bmatrix} I_1 \\ I_2 \\ \vdots \\ I_n \end{bmatrix}$

Differentiation: is a $\frac{\partial f}{\partial x} = \lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon y) - f(x)}{\epsilon} \approx f(x_m, y) - f(x_n, y)$

linear, shift invariant

operation, can write as convolution: $\boxed{-1 \mid 1}$

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

Backward pass w.r.t. inputs

$$\delta^{(l-1)} = \frac{\partial C}{\partial z^{(l-1)}} = \sum_i \sum_j \frac{\partial C}{\partial z^{(l-1)}} \frac{\partial z^{(l-1)}}{\partial z^{(l-1)}}$$

$$= \sum_i \sum_j \delta^{(l)}_{i,j} \frac{\partial}{\partial z^{(l-1)}} \left(\sum_m \sum_n w_{m,n}^{(l)} z_{i-m, j-n}^{(l-1)} + b^{(l)} \right)$$

↳ everything vanishes except where $i=i'-m, j=j'-n$. What remains is $w_{m,n}$, but as $i=i'-m \sim m-i'; j=j'-n \sim n-j'$

$= \sum_i \sum_j \delta^{(l)}_{i,j} \cdot w_{i'-i, j'-j}^{(l)} = \delta^{(l)} * w_{i'-i, j'-j}^{(l)} = \delta^{(l)} * \text{ROT}_{180}(w^{(l)})$

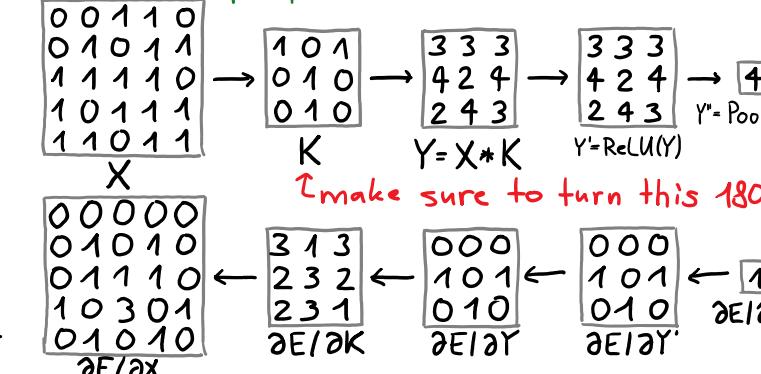
Weight Update:

$$\frac{\partial C}{\partial w_{m,n}} = \sum_i \sum_j \frac{\partial C}{\partial z^{(l-1)}} \frac{\partial z^{(l-1)}}{\partial w_{m,n}}$$

$$= \sum_i \sum_j \delta^{(l)}_{i,j} \frac{\partial}{\partial w_{m,n}} \left(\sum_m \sum_n w_{m,n}^{(l)} z_{i-m, j-n}^{(l-1)} + b^{(l)} \right)$$

$$= \sum_i \sum_j \delta^{(l)}_{i,j} z_{i-m, j-n}^{(l-1)} = \delta^{(l)} * z_{i-m, j-n}^{(l-1)} = \delta^{(l)} * \text{ROT}_{180}(z^{(l-1)})$$

CNN Backprop Exercise



$\partial E / \partial Y'$: Just duplicate value where max was

$\partial E / \partial Y$: Turn to 1 if greater than 0 else 0

$\partial E / \partial K$: Look which elements of X were involved to make that value in Y , then turn that patch 180°, multiply with the corresponding in $\partial E / \partial Y$ and add all of these together.

$\partial E / \partial X$: Duplicate the 180° turned K at all the locations of $\partial E / \partial Y$. (like upconvolution)

Pooling Layer: makes the representation smaller and more manageable. $\bar{z}^{(l)} = \max \{z_i^{(l-1)}\}$
 $\frac{\partial \bar{z}^{(l)}}{\partial z^{(l-1)}} = \begin{cases} 1 & \text{if } i = \arg\max \{z_i^{(l-1)}\} \\ 0 & \text{otherwise.} \end{cases}$ $\bar{z}^{(l-1)} = \{\bar{z}^{(l)}\}; \text{ where } i^* = \arg\max \{z_i^{(l-1)}\}$

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 → 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: block → stack many on top of each other

Semantic Segmentation: pixel-wise CV task
Using patches or all at once (fully convolutional)
↳ 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.

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. Smoothes network with respect to attack)

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 $g(t; x^t, \dots, x^1)$, 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: $y^t = W_h h^t$, $h^t = \tanh(W_h h^{t-1} + W_x x^t) = f(h^{t-1}, x^t; \theta)$

$L^t = \|y^t - y^t\|_2^2$, $L = \sum L^t$.

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 \exists eigendecomposition 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 → vanishes, if Λ gets large → explodes

Backprop through time:

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T \frac{\partial L^t}{\partial w} \quad \frac{\partial L^t}{\partial w} = \sum_{k=1}^T \frac{\partial L^t}{\partial y^t} \frac{\partial y^t}{\partial h^t} \frac{\partial h^t}{\partial w}$$

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

RNN Backprop Exercise

Given an RNN with formulas $a_t = Wh_{t-1} + Ux_t + b$, $h_t = \tanh(a_t)$, $\hat{y}_k = \frac{e^{o_k}}{\sum_i e^{o_i}}$, $L_t = -\sum_k y_{t,k} \log \hat{y}_{t,k}$, $L = \sum_t L_t$

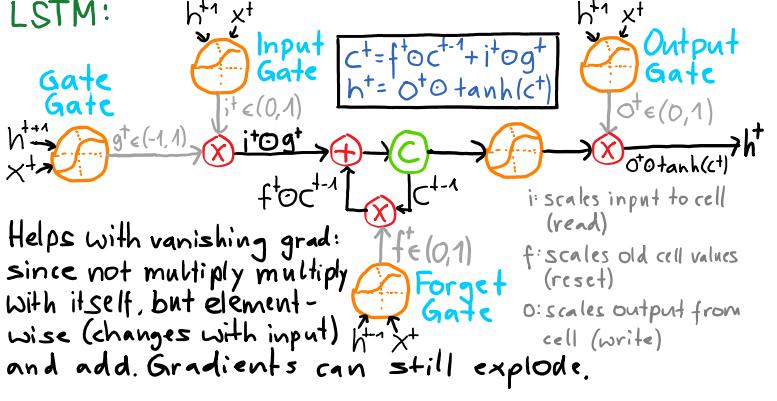
We want to compute some gradients:

$$\frac{\partial L}{\partial o_t} = -\sum_k y_{t,k} \frac{\partial \log \hat{y}_{t,k}}{\partial o_t} = -\sum_k y_{t,k} \frac{\partial \log \hat{y}_{t,k}}{\$$

Vanishing / Exploding Gradient: We look at $\frac{\partial h^t}{\partial h^k} = \prod_{i=k+1}^t \frac{\partial h^i}{\partial h^{i-1}} = \prod_{i=k+1}^t W_{hh}^\top \text{diag}[f'(h^{i-1})]$ a small part of the gradient. Assume: λ_1 is the largest singular value of W_{hh} and $\|\text{diag}[f'(h^{i-1})]\| < \gamma$ where $\gamma \in \mathbb{R}$ (fine for tanh). Case 1: It is sufficient $\lambda_1 < \frac{1}{\gamma}$ for the vanishing gradient problem.

$$\forall i, \left\| \frac{\partial h^i}{\partial h^{i-1}} \right\| \leq \|W_{hh}^\top\| \cdot \|\text{diag}[f'(h^{i-1})]\| < \frac{1}{\gamma} \cdot \gamma = 1$$

Let $\eta \in \mathbb{R}$ s.t. $\forall i, \left\| \frac{\partial h^i}{\partial h^{i-1}} \right\| \leq \eta < 1$. We can see that $\left\| \frac{\partial h^t}{\partial h^k} \right\| \leq \eta^{(t-k)}$ \Rightarrow this goes to zero for $t \rightarrow \infty$. Case 2: $\lambda_1 > \frac{1}{\gamma}$ \Rightarrow exploding



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

GRU: Newer, simpler LSTM alternative

$$r_t = O(W_r x_t + W_h h_{t-1} + b_r), \quad z_t = O(W_z x_t + W_h h_{t-1} + b_z), \\ h_t = \tanh(W_x x_t + W_h (r_t \odot h_{t-1}) + b_h), \quad h_t = z_t \odot h_{t-1} + (1 - z_t) \cdot h_t$$

Image Captioning: Take some CNN (e.g. AlexNet), drop the fully connected part to label classes (not all FC), append a RNN. Get token probabilities, sample and feed in again (until reaches some END token)

GENERATIVE MODELS

Generative Modelling: Non-probabilistic (discriminative models) vs probabilistic (generative that are given random vector or discriminative that are also given input (conditionally generative))

(Taxonomy of) Generative models



Autoencoders: Finding meaningful DoF that describe the high-dimensional signal with fewer dimensions. encoder f projects input space X into latent space Z . decoder g projects Z back into input space X . $\hat{\theta}_f, \hat{\theta}_g = \arg \min_{\theta_f, \theta_g} \sum_{n=1}^N \|x_n - g(f(x_n))\|^2$. If f and g are linear, PCA is optimal solution.

Generative Capabilities: by inducing simple density model over the latent space Z , g can be used to generate new data: $f(x) \sim N(\mu, \sigma^2 I)$. Problem: gaps in the latent space \rightarrow not meaningful samples.

VARIATIONAL AUTOENCODER

Variational Autoencoders: probabilistic version that allows for sampling from the learned model to generate data with natural variation.

Assume training data D is generated from underlying latent distribution Z . x : desired output (e.g. image). Idea: Sample from prior $p_\theta(z)$ z : latent factors that control generation process of x (e.g. scale). \hookrightarrow conditional complex (generate image) \Rightarrow param. with neural net. \hookrightarrow choose prior to be a simple distribution (e.g. Gaussian).

VAE Training: Want to capture this process via estimation of the parameters $p_\theta(x) = \int_Z p_\theta(x|z)p_\theta(z)dz$. Problem: integral over continuous RV. Posterior also intractable. \hookrightarrow solution: define additional encoder network $q_\phi(z|x)$ to approx. $p_\theta(x)$ using Bayes' rule.

$$\mathcal{L}(x^{(i)}, \theta, \phi) = \mathbb{E}_z [\log p_\theta(x^{(i)}|z)]$$

(maximize reconstruction likelihood) \rightarrow encourages samples

$$-D_{KL}(q_\phi(z|x^{(i)}) \| p_\theta(z^{(i)}))$$

(make approximate posterior similar to prior)

↳ needed, otherwise model learns a zero variance (like autoencoder)

↳ encourage encoder to project latent space representations even to ground the center of the latent space. At runtime, only use decoder network, sample z from prior

ELBO Derivation: start with data log-likelihood:

$$\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log(p_\theta(x^{(i)}))] = (\text{Po}(x^{(i)}) \text{ does not depend on } z) \\ \mathbb{E}_z \left[\log \frac{p_\theta(x^{(i)}|z) \cdot p_\theta(z)}{p_\theta(z|x^{(i)}) \cdot q_\phi(z|x^{(i)})} \right] = \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}] + \mathbb{E}_z [\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}] \\ = \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - D_{KL}(q_\phi(z|x^{(i)}) \| p_\theta(z|x^{(i)})) \geq 0$$

Monte Carlo Gradient Estimator

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

$$\nabla_\phi \mathbb{E}_{p_\phi(z)} [f(z)] = \nabla_\phi \int p_\phi(z) f(z) dz = \nabla_\phi \int p(\epsilon) f(g(\epsilon, \varphi)) d\epsilon$$

Note: this works because of the transformation and the differential area under the integral remains the same ($Dp(z) dz = p(\epsilon) d\epsilon$)

$$= \nabla_\phi \int p(\epsilon) f(g(\epsilon, \varphi)) d\epsilon = \nabla_\phi \mathbb{E}_{p(\epsilon)} [f(g(\epsilon, \varphi))] = \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(g(\epsilon, \varphi))]$$

We can then estimate this using: $\frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(g(\epsilon^{(s)}, \varphi))]$

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

DKL Non-Negativity: $-D_{KL} = - \int_X P(x) \log \frac{P(x)}{q(x)} dx = \int_X P(x) \log \frac{q(x)}{P(x)} dx \leq \log \int_X P(x) \frac{q(x)}{P(x)} dx = \log \int_X q(x) dx = \log(1) = 0$
Trick: $E(\varphi(x)) \leq \varphi E(x)$ iff φ is concave (\log is)

Reparameterization Trick: (to enable VAE backprop)

Given Gaussian with μ and σ^2 $\sim N(\mu, \sigma^2 I)$

Assume \exists underlying random variable $\epsilon \sim N(0, 1)$

$$Z = \mu + \sigma \epsilon = f(x, \epsilon, \theta) \Rightarrow \text{can take derivative w.r.t. } \mu, \sigma$$

(tells us how changes in μ, σ affects output for fixed ϵ)

Advantage of using N : when modelling $p_\theta(z)$ and $q_\phi(z|x^{(i)})$ with normal distributions with diagonal covariance matrix, it allows drawing easily from the prior & can calculate DKL analytically.

Equation for analytical DKL with N :

$$\text{For any } p(z) = N(\mu_p, \Sigma_p^2 I) \text{ and } q(z) = N(\mu_q, \Sigma_q^2 I) \\ \int p(z) \log(q(z)) dz = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J \log \Sigma_{q,j}^2 - \frac{1}{2} \sum_{j=1}^J \frac{\Sigma_{p,j}^2 + (\mu_{p,j} - \mu_{q,j})^2}{\Sigma_{q,j}^2}$$

VAE Backprop: $\mathbb{E}_z [\log p_\theta(x^{(i)}|z)] - D_{KL}(q_\phi(z|x^{(i)}) \| p_\theta(z))$

$$= \mathbb{E}_z [\log p_\theta(x^{(i)}|z)] + \mathbb{E}_z [\log \frac{p_\theta(z)}{q_\phi(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_\theta(x^{(i)}, z)}{q_\phi(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_\theta(x^{(i)}, z)}{q_\phi(z|x^{(i)})}]$$

$$\nabla_\theta \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log \frac{p_\theta(x^{(i)}, z)}{q_\phi(z|x^{(i)})}] = \nabla_\theta \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\log \frac{p_\theta(x^{(i)})}{q_\phi(x^{(i)}|z)}] = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})} [\nabla_\theta \log \frac{p_\theta(x^{(i)})}{q_\phi(x^{(i)}|z)}]$$

$$\approx \frac{1}{K} \sum_{k=1}^K \nabla_\theta \log \frac{p_\theta(x^{(i)})}{q_\phi(x^{(i)}|z_k)}$$

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

Hierarchical Latent Variable Models Exercise

To increase expressiveness we make a hierarchy by stacking. Circles are stochastic variables, diamonds are deterministic variables. Equivalent to VAE for $L=1$.

$$q_\phi(z_1, \dots, z_L|x) = q_\phi(z_L|x) \cdot \dots \cdot q_\phi(z_1|x)$$

the h are all deterministic operations. (Inference)

$$p_\theta(z_1, \dots, z_L) = p_\theta(z_L) p(z_{L-1}|z_L) \dots p(z_1|z_2)$$

ELBO: $\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - D_{KL}(q_\phi(z_1, \dots, z_L|x) \| p_\theta(z_1, \dots, z_L))$

$$= \mathbb{E}_{q_\phi(z_1, \dots, z_L|x)} [\log p_\theta(x|z_1, \dots, z_L)] - D_{KL}(q_\phi(z_1, \dots, z_L|x) \| p_\theta(z_1, \dots, z_L))$$

The last transformation is possible, as x only depends on z_1 . We now look at the DKL term:

$$= \int_{z_1, \dots, z_L} q_\phi(z_1, \dots, z_L|x) \cdot \log \frac{p_\theta(z_1, \dots, z_L)}{q_\phi(z_1, \dots, z_L|x)} dz_1 \dots dz_L$$

$$= \int_{z_1, \dots, z_L} \left(\prod_{i=1}^L q_\phi(z_i|x) \right) \cdot \left(\sum_{i=1}^L \log \frac{p_\theta(z_i|z_{i+1})}{q_\phi(z_i|z_{i+1})} \right) dz_1 \dots dz_L = p_\theta(z_L|z_{L+1}) = p_\theta(z)$$

$$= \sum_{i=1}^L \int_{z_1, \dots, z_L} \prod_{j=1}^{i-1} q_\phi(z_j|x) \cdot q_\phi(z_i|z_{i+1}) \log \frac{p_\theta(z_i|z_{i+1})}{q_\phi(z_i|z_{i+1})} dz_1 \dots dz_{i+1}$$

$$= \int_{z_L} q_\phi(z_L|x) \log \frac{p_\theta(z_L)}{q_\phi(z_L|x)} dz_L + \sum_{i=1}^{L-1} \int_{z_1, \dots, z_L} \int_{z_{i+1}} q_\phi(z_i, z_{i+1}|x) \log \frac{p_\theta(z_i|z_{i+1})}{q_\phi(z_i|z_{i+1})} dz_i dz_{i+1}$$

$$= -D_{KL}(q_\phi(z_L|x) \| p_\theta(z_L)) + \sum_{i=1}^{L-1} \int_{z_1, \dots, z_L} q_\phi(z_{i+1}|x) \int_{z_i} q_\phi(z_i|x) \log \frac{p_\theta(z_i|z_{i+1})}{q_\phi(z_i|z_{i+1})} dz_i dz_{i+1}$$

$$= -D_{KL}(q_\phi(z_L|x) \| p_\theta(z_L)) - \sum_{i=1}^{L-1} \int_{z_1, \dots, z_L} q_\phi(z_{i+1}|x) D_{KL}(q_\phi(z_i|x) \| p_\theta(z_i|z_{i+1})) dz_{i+1}$$

$$= -D_{KL}(q_\phi(z_L|x) \| p_\theta(z_L)) - \sum_{i=1}^{L-1} \int_{z_1, \dots, z_L} q_\phi(z_{i+1}|x) \left[D_{KL}(q_\phi(z_i|x) \| p_\theta(z_i|z_{i+1})) \right]$$

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

GAN Objective: Consider only the discriminator

$$\mathcal{L}(D) = -\frac{1}{N} \sum_{i=1}^N (y_i \log(D(x^{(i)})) + (1-y_i) \log(1-D(G(z^{(i)})))$$

for GANs we know 50% are $y=1$ and 50% are $y=0$. Our goal is to train D & G jointly:

$$= \frac{1}{2} \left(\mathbb{E}_{x \sim p_d} [\log D(x)] + \mathbb{E}_{z \sim p_g} [\log (1 - D(G(z)))] \right)$$

$$= \frac{1}{2} \left(\mathbb{E}_{x \sim p_d} [\log D(x)] + \mathbb{E}_{z \sim p_g} [\log (1 - D(x))] \right)$$

\Rightarrow a good G maximizes loss function of D .

Value function: $V(G, D) = \mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{z \sim p_g} [\log(1 - D(G(z)))]$

$\Rightarrow G^*$ will fool any D : $G^* = \arg \max_D V(G, D)$

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: $V(G, D) = \int_X P_d(x) \log(D(x)) + \int_X P_g(x) \log(1 - D(x)) dx$

Optimal D : $\nabla_\theta \{0, 1\} f(y) = a \cdot \log(y) + b \cdot \log(1 - y)$

$$f'(x) = 0 \Leftrightarrow \frac{a}{y} - \frac{b}{1-y} = 0 \Rightarrow y = \frac{a}{a+b}, f'(\frac{a}{a+b}) = -\frac{a}{y^2} - \frac{b}{(1-y)^2} < 0 \Rightarrow D_G^* = \frac{P_d}{P_d + P_g}$$

$$\nabla_\theta \{0, 1\} f(y) = a \cdot \log(y) + b \cdot \log(1 - y)$$

$$f'(x) = 0 \Leftrightarrow \frac{a}{y} - \frac{b}{1-y} = 0 \Rightarrow y = \frac{a}{a+b}, f'(\frac{a}{a+b}) = -\frac{a}{y^2} - \frac{b}{(1$$

can 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:** RNNADE: conditionals modelled by mixture of Gaussians, DeepNADE: a DNN trained to assign a cond. distr. to any given var. given any subset of others (orderless), ConvNADE

Masked Autoencoder Distribution Estimator (MADE): constrain autoencoder s.t. output can be used as conditional $p(x_i | x_{<i})$. To fulfill autoregressive property: no computational path output unit x_d and any of the input units x_1, \dots, x_d must exist. Can also be seen as masks $m_i^{(l)}$: 1 if unit i of layer l is connected to unit j of layer $l-1$. 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. $p(x) = \prod_{i=1}^n p(x_i | x_{<i}) \rightarrow$ 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 slow due to explicit pixel dependencies

PixelCNN: Still generate starting from corner, but model the dependency as CNN:  (masked convolutions) $p(x_i | x_{<i}) = p(x_{i,R} | x_{<i}) \cdot p_{i,G}(x_{i,R}, x_{<i}) \cdot p_{i,B}(x_{i,G}, x_{i,R}, x_{<i})$ Stacks can be used to speed it up. Add information between vertical and horizontal stack (preserve pixel dependencies) $h_{k+1} = \tanh(W_{k,f} * h_k) \odot (W_{k,g} * h_k)$. 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 timesteps. 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.

EYE GAZE ESTIMATION

Eye Gaze Estimation: **Challenges**: head poses, gaze directions, illumination conditions, personal appearances. **Applications**: explicit eye input (eye typing, gaming), user modelling (what are they

doing/focusing on), attentive user interface (autopause), passive eye monitoring (ads).

Typical Gaze Estimation Pipeline: Grab frame, segment face, detect landmarks, segment eye, 3D gaze direction or 2D point of regard/gaze algorithm

Geometry-based Methods: hand-picked features (e.g. iris center-ey-vector) often with 3D head pose feature, not robust to head movements, eyeball model fitting calibration from few samples, robust to head movement.

LeNet/GazeNet: uses image, DNN. Can also include both eyes, face & face grid.

Synthetic data: need domain adaption (network that takes unlabeled real images & refines them).

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 .

$$V_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s] = \sum_a \pi(a|s) \sum_s \sum_r P(s, r | s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s]]$$

$$= \sum_a \pi(a|s) \sum_s \sum_r P(s, r | s, a) [r + \gamma V_\pi(s)] \text{ where } G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

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' :

$\pi'(s) = \arg\max_{a \in A} (r(s, a) + \gamma V_\pi(s))$. One can show $V_{\pi'}(s) \geq V_\pi(s)$. The optimal policy π^* is greedy w.r.t. $V_{\pi'}$. π^* fulfills the Bellman Optimality Equation: $V^*(s) = \max_{a \in A} \{r(s, a) + \gamma V^*(s)\}$. Can use step cost instead of γ . Update backwards each time doing the action that gives the highest reward until nothing changes anymore.

$i \leftarrow 0, V_0 < 0$ Theorem: value iteration converges. At convergence found the optimal V^* . Need to find the optimal (stationary & greedy) policy.

$\text{while } \Delta > \theta$ for all $s \in S$ $V \leftarrow V(s)$

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

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

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

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 $\mathbb{E}[\cdot]$: $V_\pi(s) = \mathbb{E}_{a \sim \pi(\cdot | s), s_{t+1} \sim P(\cdot | s, a)} \sum_{i=0}^T \gamma^i r(s_{t+i}, a_{t+i})$, so we can just run policies with the given policy and compute samples of: $\sum_{i=0}^T \gamma^i r(s_{t+i}, a_{t+i})$

(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: $\Delta V(s) = r(s, a) + \gamma V(s') - V(s)$, $V(s) \leftarrow V(s) + \alpha \Delta V(s)$, $\alpha > 0$

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

Policy Optimization Methods based on TD learning

Grad-free methods Policy grad. Policy iter. Value iter.

-Param. Perturbation Evolutionary Strat.

-Actor-critic methods Q-learning

DEEP REINFORCEMENT LEARNING

Policy gradients: $\pi(a|s) = \mathcal{N}(\mu_t, \sigma_t^2 | s_t)$. We can collect rollout trajectories: $p(\tau) = p(s_0, a_0, \dots, s_T, a_T)$

$= p(s_0) \prod_{t=0}^T \pi(a_t | s_t) p(s_{t+1} | a_t, s_t)$ then we update the policy: good traj. more & bad traj. less likely.

Difference Policy optimization vs DP: directly optimize desired quantity (not indirect & exploit problem structure and self-consistency), more compatible with modern ML machinery, more versatile & flexible, more compatible with auxiliary objectives, less compatible with off-policy and exploration, less sample efficient.

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

$$J(\theta) = \mathbb{E}_{\tau \sim p(\tau)} [\sum_t \gamma^t r(s_t, a_t)] = \mathbb{E}_{\tau \sim p(\tau)} [\mathbb{E}_{\tau \sim p(\tau)} [\sum_t \gamma^t r(s_t, a_t)]] = \int p(\tau) \mathbb{E}_{\tau \sim p(\tau)} [\sum_t \gamma^t r(s_t, a_t)] d\tau$$

$$= \int p(\tau) \nabla_\theta \log p(\tau) r(\tau) d\tau \quad (\text{since } \frac{d \log(p(\tau))}{d\theta} = \frac{p'(\tau)}{p(\tau)}) = \mathbb{E}_{\tau \sim p(\tau)} [\nabla_\theta \log p(\tau) r(\tau)]$$

$$\log p(\tau) = \log p(s_0, a_0, \dots, s_T, a_T) = \log [p(s_0) \prod_{t=0}^T \pi(a_t | s_t) p(s_{t+1} | a_t, s_t)]$$

$$\Rightarrow \nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim p(\tau)} [\sum_{t=0}^T \log \pi(a_t | s_t) (\sum_{i=0}^T \gamma^i r(s_i, a_i))]$$

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 Trick

$\nabla_\theta p_\theta(x) = p_\theta(x) \nabla_\theta \log p_\theta(x)$ which allows us:

$$\nabla_\theta \mathbb{E}_{x \sim p_\theta(x)} [f(x)] = \nabla_\theta \int f(x) p_\theta(x) dx = \int f(x) \nabla_\theta p_\theta(x) dx$$

$$= \int f(x) p_\theta(x) \nabla_\theta \log p_\theta(x) dx = \mathbb{E}_{x \sim p_\theta(x)} [f(x) \nabla_\theta \log p_\theta(x)]$$

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 $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i (\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t)) (\sum_{t=0}^T \gamma^t r(s_t, a_t))$ ($\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$) (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) $r(\tau) = \sum_{t=0}^T \gamma^t r(s_t, a_t) - b(s_t)$. Still unbiased!

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

TRPO: tries to solve stepsize problem (slow learning vs. disastrous effects) by formulating it as a constraint optimization problem. (second order optimization, computationally heavy).

PPO: Relax hard TRPO constraint to penalize with tuning parameter (allows some bad decisions but generally within trust region). Even better: clipped objective. Simpler, computationally easier than TRPO.

DDPG: When selecting action we add noise to it. Also uses experience replay & target networks (off-policy)

Issues with DRL: Exploration vs Exploitation (get stuck in local minima), Reward hacking (just get power-ups), Sample efficiency (need excessive amounts of data, infeasible in real-world scenarios), domain knowledge (small tricks lead to different results e.g. comparison of different implementations)