

1 Probability

Sum Rule $P(X = x_i) = \sum_{j=1}^J p(X = x_i, Y = y_i)$

Product rule $P(X, Y) = P(Y|X)P(X)$

Independence $P(X, Y) = P(X)P(Y)$

$$\text{Bayes' Rule } P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \frac{\sum_{i=1}^K P(X|Y_i)P(Y_i)}{\sum_{i=1}^K P(X|Y_i)}$$

$$\text{Cond. Ind. } X \perp Y|Z \implies P(\bar{X}, Y|Z) = P(X|Z)P(Y|Z)$$

$$\text{Cond. Ind. } X \perp Y|Z \implies P(X|Y, Z) = P(X|Z)$$

$$\mathbb{E}[X] = \int_X t \cdot f_X(t) dt =: \mu_X$$

$$\text{Cov}(X, Y) = \mathbb{E}_{x,y}[(X - \mathbb{E}_x[X])(Y - \mathbb{E}_y[Y])]$$

$$\text{Cov}(X) := \text{Cov}(X, X) = \text{Var}[X]$$

$$X, Y \text{ independent} \implies \text{Cov}(X, Y) = 0$$

$$\mathbf{XX}^T \geq 0 \text{ (symmetric positive semidefinite)}$$

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{Var}[\mathbf{AX}] = \mathbf{A} \text{Var}[\mathbf{X}] \mathbf{A}^T = \text{Var}[aX + b] = a^2 \text{Var}[X]$$

$$\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$$

$$\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{i \neq j} a_i a_j \text{Cov}(X_i, X_j)$$

$$\frac{\partial}{\partial t} P(X \leq t) = \frac{\partial}{\partial t} F_X(t) = f_X(t) \text{ (derivative of c.d.f. is p.d.f.)}$$

$$f_{\alpha Y}(z) = \frac{1}{\alpha} f_Y\left(\frac{z}{\alpha}\right)$$

$$\mathbf{T.} \text{The moment generating function (MGF) } \psi_X(t) = \mathbb{E}[e^{tX}] \text{ characterizes the distr. of a rv}$$

$$\mathbf{T.} \text{If } X_1, \dots, X_n \text{ are indep. rvs with MGFs } M_{X_i}(t) = \mathbb{E}[e^{tX_i}], \text{ then the MGF of } Y = \sum_{i=1}^n a_i X_i \text{ is } M_Y(t) = \prod_{i=1}^n M_{X_i}(a_i t).$$

$$\mathbf{T.} \text{Let } X, Y \text{ be indep., then the p.d.f. of } Z = X + Y \text{ is the conv. of the p.d.f. of } X \text{ and } Y: f_Z(z) = \int_{\mathbb{R}} f_X(x) f_Y(z-x) dt = \int_{\mathbb{R}} f_X(x) f_Y(z-x) dt$$

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \boldsymbol{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

$$\mathbf{T.} P\left(\begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_L \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_L \end{bmatrix} \middle| \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_L \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \vdots & \vdots \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^d, \Sigma_{11} \in \mathbb{R}^{d \times d} \text{ p.s.d.}, \Sigma_{12} \in \mathbb{R}^{d \times d} \text{ p.s.d.}$$

$$\mathbf{a}_2, \mathbf{a}_2 \in \mathbb{R}^d, \Sigma_{22} \in \mathbb{R}^{d \times d} \text{ p.s.d.}, \Sigma_{21} \in \mathbb{R}^{d \times d} \text{ p.s.d.}$$

$$P(\mathbf{a}_2 | \mathbf{a}_1 = \mathbf{z}) = \mathcal{N}(\mathbf{a}_2 | \mathbf{u}_2 + \Sigma_{21} \Sigma_{11}^{-1} (\mathbf{z} - \mathbf{u}_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

3 Linear Algebra

Kernels are positive semi-definite matrices.

D. (Positive Semi-Definite Matrix) A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is PSD if for all non-zero vectors $\mathbf{x} \in \mathbb{R}^n$: $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ Properties:

- All eigenvalues $\lambda_i \geq 0$.
- The trace $\text{Tr}(\mathbf{A}) \geq 0$ and determinant $\det(\mathbf{A}) \geq 0$.
- Cholesky Decomposition exists: $\mathbf{A} = LL^T$.

T. (Sylvester Criterion) A $d \times d$ matrix is positive semi-definite if and only if all the upper left $k \times k$ for $k = 1, \dots, d$ have a positive determinant.

Negative definite: $\det < 0$ for all odd-sized minors, and $\det > 0$ for all even-sized minors

Special case: $(\sum_i x_i)^2 \leq (\sum_i x_i^2)(\sum_i y_i^2)$.

Special case: $\mathbb{E}[XY]^2 \leq \mathbb{E}[X^2]\mathbb{E}[Y^2]$.

D. (Convex Set) A set $S \subseteq \mathbb{R}^d$ is called convex if $\forall \mathbf{x}, \mathbf{x}' \in S, \forall \lambda \in [0, 1]: \lambda \mathbf{x} + (1-\lambda)\mathbf{x}' \in S$.

Com. Any point on the line between two points is within the set. \mathbb{R}^d is convex.

D. (Convex Function) A function $f: S \rightarrow \mathbb{R}$ defined on a convex set $S \subseteq \mathbb{R}^d$ is called convex if

$$\forall \mathbf{x}, \mathbf{x}' \in S, \lambda \in [0, 1]: f(\lambda \mathbf{x} + (1-\lambda)\mathbf{x}') \leq \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{x}')$$

$$\text{Com. A function is strictly convex if the line segment between any two points on the graph of the function lies strictly above the graph. This guarantees that there is a unique global minimum.}$$

T. (Properties of Convex Functions)

$f(y) \geq f(x) + \nabla f(x)^T (y - x)$

$f'(y) \geq 0$

Local minima are global minima, strictly convex functions have a unique global minimum

If f, g are convex then $\alpha f + \beta g$ is convex for $\alpha, \beta \geq 0$

If f, g are convex then $\max(f, g)$ is convex

If f is convex and g is convex and non-decreasing then $g \circ f$ is convex

D. (Strongly Convex Function) A function f is μ -strongly convex if it curves up at least as much as a quadratic function with curvature $\mu > 0$. For all x, y :

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

Relation to Optimization:

Guarantees: Ensures a unique global minimum exists.

Convergence: Gradient Descent on strongly convex (and Lipschitz smooth) functions guarantees a linear convergence rate ($O(c^k)$ for some $c < 1$).

Condition Number: The convergence speed depends on the condition number $\kappa = L/\mu$. If κ is large (poor conditioning), convergence slows down.

D. (Condition Number) The condition number $\kappa(A)$ measures the sensitivity of a function's output to small perturbations in the input. For a symmetric positive semi-definite matrix (like the Hessian H of a loss function), it is the ratio of the largest to the smallest eigenvalue:

$$\kappa(H) = \frac{|\lambda_{\max}|}{|\lambda_{\min}|} \geq 1$$

Implications for Optimization:

Well-conditioned ($\kappa \approx 1$): The contours form narrow, elongated ellipses (steep valleys). Gradient Descent tends to oscillate ("zigzag") across the narrow valley rather than moving down the slope, leading to very slow convergence.

Com. Momentum and Adaptive Learning Rate methods (like Adam) are specifically designed to mitigate the issues caused by high condition numbers.

3.3.3 Vector-by-Vector

A, C, D, a, b, e not a function of \mathbf{x} ,

$\mathbf{f} = \mathbf{f}(\mathbf{x}), \mathbf{g} = \mathbf{g}(\mathbf{x}), \mathbf{h} = \mathbf{h}(\mathbf{x}), u = u(\mathbf{x}), v = v(\mathbf{x})$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{u}(\mathbf{x}) \mathbf{v}(\mathbf{x})] = \mathbf{u}(\mathbf{x}) \frac{\partial \mathbf{v}(\mathbf{x})}{\partial \mathbf{x}} + \mathbf{v}(\mathbf{x}) \frac{\partial \mathbf{u}(\mathbf{x})}{\partial \mathbf{x}}$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})] = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{g}(\mathbf{x}) + \mathbf{f}(\mathbf{x})^T \frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}}$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})^T \mathbf{A} \mathbf{g}(\mathbf{x})] = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}}^T \mathbf{A} \mathbf{g}(\mathbf{x}) + \mathbf{f}(\mathbf{x})^T \frac{\partial \mathbf{A} \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}}$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial \mathbf{x}}^T \mathbf{x} + \mathbf{a}^T \frac{\partial \mathbf{x}}{\partial \mathbf{x}}$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b}$

$\frac{\partial}{\partial \mathbf{x}} [\mathbf{a}^T \mathbf{A} \mathbf{x}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$

3.4.3 Scalar-by-Matrix

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.4 Matrix-by-Matrix

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.5 Vector-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.6 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.7 Vector-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.8 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.9 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.10 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.11 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.12 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.13 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.14 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.15 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.16 Matrix-by-Matrix (Generalized Gradient)

$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n: \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|$

3.4.17 Matrix-by-Matrix (Generalized Gradient)

D. (Learning Rate (Robbins-Monro Conditions)) For Stochastic Gradient Descent (SGD) to guarantee convergence to a local minimum (in non-convex cases) or global minimum (in convex cases), the step size schedule α_t must satisfy two conditions:

- 1. **Explore Forever:** The steps must sum to infinity to ensure the algorithm can reach the optimum from any starting point, no matter how far.

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

2. Decay Fast Enough: The squared steps must sum to a finite value to ensure the variance (noise) of the updates tends to zero, preventing the parameters from oscillating forever around the minimum.

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Example: A schedule of $\alpha_t = \frac{1}{\sqrt{t}}$ satisfies both, whereas $\alpha_t = \frac{1}{t}$ satisfies the first but not the second.

D. (Momentum) Accelerates SGD by navigating along the relevant direction and softening oscillations in irrelevant directions. It maintains a velocity vector v (exponential moving average of past gradients).

$$\theta^{t+1} = \theta^t - \eta \nabla J(\theta^t) + \beta(\theta^t - \theta^{t-1})$$

where $\beta \in [0, 1]$ is the momentum term (friction).

D. (Nesterov Accelerated Gradient (NAG)) A "look-ahead" version of GD. It computes the gradient at the approximate future position of the parameters rather than the current position.

$$\begin{aligned} \theta^{t+1} &= \theta^t + \beta(\theta^t - \theta^{t-1}) \\ \theta^{t+1} &= \theta^{t+1} - \eta \nabla f(\theta^{t+1}) \end{aligned}$$

D. (Adaptive Learning Rate Methods) These methods adjust the learning rate for each parameter individually, scaling them based on the history of gradient magnitudes.

D. (RMSPprop) Designed to resolve the diminishing learning rates of Adagrad. It uses a decaying average of squared gradients.

$$E[g^2]_t = \beta E[g^2]_{t-1} + (1-\beta)(\nabla J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{E[g^2]_t + \epsilon}} \nabla J(\theta_t)$$

D. (Adam (Adaptive Moment Estimation)) Combines Momentum (first moment m_t) and RMSProp (second moment v_t). It also includes bias correction terms \hat{m}_t, \hat{v}_t to account for initialization at zero.

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1-\beta_1) \nabla J(\theta_t) \\ v_t &= \beta_2 v_{t-1} + (1-\beta_2) (\nabla J(\theta_t))^2 \\ \theta_{t+1} &= \theta_t - \frac{\alpha}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t \end{aligned}$$

- 11.4 - Backpropagation Graphs

The approach a backpropagation graph for a loss L from a FF network graph is as follows:

1. (RED) Starting at the loss L backward: for each node n and for each of its outputs o which has a path to the loss, add the node $\frac{\partial L}{\partial n}$ (at the height of node n). Connect the output o and that node n to that newly created node.
2. (BROWN) Starting at the loss backward: for each node n create a node $\frac{\partial L}{\partial n}$ (if it doesn't exist already) and connect each previously created partial node $(\frac{\partial L}{\partial n})$ to it, and the previously created (in this step) $\frac{\partial L}{\partial o}$'s too.

12 Convolutional Neural Networks

- 12.1 - Convolutional Layers

D. (Transform) A transform T is a mapping from one function space \mathcal{F} to another function space \mathcal{F}' . So $T: \mathcal{F} \rightarrow \mathcal{F}'$.

D. (Linear Transform) A transform T is linear, if for all functions f, g and scalars α, β , $T(\alpha f + \beta g) = \alpha(Tf) + \beta(Tg)$.

D. (Integral Transform) An integral transform is any transform T of the following form

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

Com. The fourier transform is an example of an integral transform.

T. Any integral transform is a linear transform.

D. (Convolution) Given two functions $f, h: \mathbb{R} \rightarrow \mathbb{R}$, their convolution is defined as

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

Com. Whether the convolution exists depends on the properties of f and h (the integral might diverge). However, a typical use is f = signal, and h = fast decaying kernel function.

T. (Convolution Theorem) Any linear, translation-invariant transformation T can be written as a convolution with a suitable h .

T. (Convolutions are commutative and associative)

T. (Convolutions are shift-invariant), we define $f_\Delta(t) := f(t + \Delta)$. Then

$$(f * h)_\Delta(u) = (f * h)_\Delta(u)$$

D. (Fourier Transform) The fourier transform of a function f is defined as

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

and its inverse as

$$(\mathcal{F}^{-1}f)(u) := \int_{-\infty}^{\infty} f(t) e^{2\pi i u t} dt$$

Com. Convolutional operators can be efficiently computed with point wise multiplication using the Fourier transform.

$$\mathcal{F}(f * h) = \mathcal{F}f \cdot \mathcal{F}h$$

and then transformed back using the inverse Fourier transform.

$$\mathcal{F}^{-1}(\mathcal{F}(f * h)) = \mathcal{F}^{-1}(\mathcal{F}f \cdot \mathcal{F}h) = f * h$$

- 12.2 - Discrete Time Convolutions

D. (Discrete Convolution)

For $f, h: \mathbb{Z} \rightarrow \mathbb{R}$, we can define the discrete convolution via

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

Com. Note that the use of rectangular brackets suggests that we're using "arrays" (discrete-time samples).

Com. Typically we use a h with finite support (window size).

D. (Multidimensional Discrete Convolution)

For $f, h: \mathbb{R}^d \rightarrow \mathbb{R}$ we have

$$\begin{aligned} (f * h)[u_1, \dots, u_d] &= \sum_{t_1=-\infty}^{\infty} \dots \sum_{t_d=-\infty}^{\infty} f(t_1, \dots, t_d)h(u_1 - t_1, \dots, u_d - t_d) \\ &= \sum_{t_1=-\infty}^{\infty} \dots \sum_{t_d=-\infty}^{\infty} f(u_1 - t_1, \dots, u_d - t_d)h(t_1, \dots, t_d) \end{aligned}$$

D. (Discrete Cross-Correlation)

Let $f, h: \mathbb{Z} \rightarrow \mathbb{R}$, then

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

$$= (h * f)[u] = (f * \bar{h})[u] \quad \text{where } \bar{h}(t) = h(-t).$$

aka "sliding inner product", non-commutative, kernel "flipped over" ($u + t$ instead of $u - t$). If kernel symmetric: cross-correlation = convolution.

- 12.3 - Convolution via Matrices

Represent the input signal, the kernel and the output as vectors. Copy the kernel as columns into the matrix offsetting it by one more very time (gives a band matrix (special case of Toeplitz matrix)). Then the convolution is just a matrix-vector product.

- 12.4 - Border Handling

There are different options to do this

- **D. (Padding of p)** Means we extend the image (or each dimension) by p on both sides (so $+2p$) and just fill in a constant there (e.g., zero).
- **D. (Same Padding)** Padding with zeros = same padding ("same" constant, i.e., 0, and we'll get a tensor of the "same" dimensions)
- **D. (Valid Padding)** Only retain values from windows which are fully-contained within the support of the signal f (see 2D example below) = valid padding

- 12.5 - Backpropagation for Convolutions

D. (Receptive Field I_i^l of x_i^l)

The receptive field I_i^l of node x_i^l is defined as $I_i^l := \{j \mid W_{ij}^l \neq 0\}$ where W^l is the Toeplitz matrix of the convolution at layer l .

Com. Hence, the receptive field of a node x_i^l are just nodes which are connected to it and have a non-zero weight.

Com. One may extend the definition of the receptive field over several layers. The further we go back in layer, the bigger the receptive field becomes due to the nested convolutions. The receptive field may be even the entire image after a few layers. Hence, the convolutions have to be small.

We have $\forall j \neq I_i^l: \frac{\partial x_i^l}{\partial x_j^l} = 0$,

Due to weight-sharing, the kernel weight h_j^l is re-used for every unit in the target layer at layer l , so when computing the derivative $\frac{\partial R}{\partial h_j^l}$ we just build an additive combination of all the derivatives (note that some of them might be zero).

$$\frac{\partial R}{\partial h_j^l} = \sum_{i=1}^{m_l} \frac{\partial R}{\partial x_i^l} \frac{\partial x_i^l}{\partial h_j^l}$$

Backpropagations of Convolutions as Convolutions

$y^{(l)}$ output of l -th layer $y^{(l-1)}$ output of $(l-1)$ -th layer / input to l -th layer w convolution filter $\frac{\partial R}{\partial y^{(l)}}$ known $y^{(l+1)} = y^{(l)} * w$

$$\begin{aligned} \frac{\partial R}{\partial w_i} &= \sum_k \frac{\partial R}{\partial y_k^{(l)}} \frac{\partial y_k^{(l)}}{\partial w_i} = \sum_k \frac{\partial R}{\partial y_k^{(l)}} \frac{\partial}{\partial w_i} [y^{(l)} * w]_k \\ &= \sum_k \frac{\partial R}{\partial w_i} \frac{\partial}{\partial w_i} \left[\sum_{o=p}^p y_{k-o}^{(l-1)} w_o \right] = \sum_k \frac{\partial R}{\partial y_k^{(l)}} y_{k-p}^{(l-1)} \\ &= \sum_k \frac{\partial R}{\partial y_k^{(l)}} y_{k-p}^{(l-1)} = \sum_k \frac{\partial R}{\partial y_k^{(l)}} \text{rot180}(y^{(l-1)})_{k-p} \end{aligned}$$

Com. Whether the convolution exists depends on the properties of f and h (the integral might diverge). However, a typical use is f = signal, and h = fast decaying kernel function.

T. (Convolution Theorem) Any linear, translation-invariant transformation T can be written as a convolution with a suitable h .

T. (Convolutions are commutative and associative)

T. (Convolutions are shift-invariant), we define $f_\Delta(t) := f(t + \Delta)$. Then

$$(f * h)_\Delta(u) = (f * h)_\Delta(u)$$

D. (Fourier Transform) The fourier transform of a function f is defined as

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

and its inverse as

$$(\mathcal{F}^{-1}f)(u) := \int_{-\infty}^{\infty} f(t) e^{2\pi i u t} dt$$

Com. Convolutional operators can be efficiently computed with point wise multiplication using the Fourier transform.

D. (Rotation180) $\forall i: \text{rot180}(x)_i = x_{-(i)}$.

- 12.6 - Pooling

There are min, max, avg, and softmax pooling. Max pooling is the most frequently used one.

D. (Max-Pooling)

- 1D: $x_{i,\max} = \max \{x_{i+k} \mid 0 \leq k < r\}$
- 2D: $x_{i,j,\max} = \max \{x_{i+k,j+l} \mid 0 \leq k, l < r\}$

- 12.7 - Sub-Sampling (aka "Strides")

Often, it is desirable to reduce the size of the feature maps. This can be achieved by skipping some of the input values in the convolution. The stride is the number of steps the kernel takes in each direction.

- 12.8 - Channels

Ex. Here we have

- an input signal that is 2D with 3 channels ($7 \times 7 \times 3$) (image x channels)
- and we want to learn two filters W_0 and W_1 , which each process the 3 channels, and sum the results of the convolutions across each channel leading to a tensor of size $3 \times 3 \times 2$ (convolution result x num convolutions)

13 Recurrent Neural Networks (RNNs)

- 13.1 - Simple Recurrent Networks

D. (Concept) Unlike CNNs (fixed filter widths), RNNs model temporal/sequence data of variable length. They maintain a hidden state z_t acting as a "memory" of the history.

Formulation: Given input sequence x_1, \dots, x_T :

$$\begin{aligned} z_t &= \phi(Uz_{t-1} + Vx_t) \\ \hat{y}_t &= \psi(Wz_t) \end{aligned}$$

where U, V, W are shared weight matrices and ϕ, ψ are non-linearities.

Unrolling: An RNN is equivalent to a deep feedforward network with T layers and *shared weights*.

Backpropagation Through Time (BPTT): Gradients are propagated backward through the unrolled graph. Since weights are shared, the total gradient is the sum of gradients at each time step:

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T \frac{\partial L}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial z_t} \frac{\partial z_t}{\partial w}$$

Teacher Forcing: During training, feed the ground truth y_{t-1}^* as input to step t rather than the model's own prediction \hat{y}_{t-1} . This speeds up convergence but causes "exposure bias" (train-test discrepancy).

Seq2Seq (Encoder-Decoder):

- **Encoder:** Compresses input sequence into a fixed context vector z_T .
- **Decoder:** Generates output sequence from z_T .
- **Attention:** Allows Decoder to "look back" at specific Encoder states z_T via learned weights, solving the bottleneck of a fixed-size context vector.

14 Optimization

- 14.1 - Objectives as Expectations

$\nabla_{\theta} \mathcal{R}(D) = \mathbb{E}_{S_N \sim P_D} [\nabla_{\theta} \mathcal{R}(S_N)] = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \mathcal{R}(\theta; \{x_i, y_i\}) \right]$

D. (Bidirectional RNNs): Process sequence in both directions to capture past and future context.

$$\hat{y}_t = \psi(Wz_t^{\rightarrow} + \tilde{W}z_t^{\leftarrow})$$

D. (Deep RNNs): Stacking multiple RNN layers to increase representational power. The output of layer l becomes the input to layer $l+1$.

D. (Deep RNNs): Compresses input sequence into a fixed context vector z_T .

14.2 Optimization Challenges in NNs: Curvatures

14.2.1 - Convergence Rates

Under certain conditions SGD converges to the optimum:

- if we have a convex, or strongly convex objective,
- and if we have Lipschitz continuous gradients,
- and a deaying learning rate, s.t.

$$\sum_{t=1}^{\infty} \eta_t = \infty, \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty$$

we get far enough our steps get always smaller

typically $\eta_t = Ct^{-\alpha}, \frac{1}{2} < \alpha < 1$ (c.f. harmonic series.)

- or we use iterate (Polyak) averaging (once we start jumping around, we average the solutions over time).

Then, we can get the following convergence rates:

- strongly-convex case: can achieve a $\mathcal{O}(1/t)$ suboptimality rate (only polynomial convergence)
- non-strongly convex case: $\mathcal{O}(1/\sqrt{t})$ suboptimality rate (even worse than polynomial convergence)

So, even if the convergence rates are not super nice, thanks to the cheap gradient computation (only one example at the time), we may even converge faster than computing the gradient on the full dataset everytime.

15 Geometric Deep Learning

- 15.1 - Sets & Point Clouds (Deep Sets)

Standard NNs assume fixed input order. Sets require **Permutation Invariance**.

D. (Dilated Convolutions) To handle long sequences without losing resolution (pooling), *dilation* introduces gaps between kernel elements.

D. (Gated Recurrent Unit (GRU)): Simplifies LSTM by merging Cell/Hidden states and Forget/Input gates.

$$z^t = (1 - \Gamma_u) \odot z^{t-1} + \Gamma_u \odot \tilde{z}^t$$

where $\Gamma_u = \sigma(G[x_t, z_{t-1}])$ is the update gate.

GRU

D. (Invariance vs. Equivariance): Let π be a permutation of indices $\{1, \dots, M\}$.

- **Invariant:** Output remains unchanged (e.g., classification).
- **Equivariant:** Output permutes exactly as input does (e.g., segmentation).

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

D. (Langevin Dynamics): Since exact inference is intractable, we use sampling. Langevin dynamics injects noise into Gradient Descent to explore the posterior distribution (sampling from $p(\theta|S)$) rather than collapsing to a minimum.

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

This mimics a physical system with friction and thermal noise.

- 16.2 - Bayesian DNNs

D. (Bayesian Paradigm): Instead of finding a point estimate θ^* , we compute a posterior distribution $p(\theta|S)$ to capture uncertainty.

$$p(\theta|S) \propto p(S|\theta)p(\theta)$$

Predictions are made by marginalizing over the posterior: $p(y|x) = \int p(y|x, \theta)p(\theta|S)d\theta$.

D. (Langevin Dynamics): Mechanism: By the Central Limit Theorem, the pre-activations (sums of many independent random variables) become Gaussian.

Com. Result: The network output $f(x)$ is a draw from a GP with mean $\mu(x) = 0$ and a specific kernel $k(x, x')$.

D. (Deep Sets Theorem): Any invariant function f can be decomposed into an element-wise encoder ϕ and an invariant aggregator ρ (e.g., sum, max).

$$f(X) = \rho \left(\sum_{m=1}^M \phi(x_m) \right)$$

D. (PointNet): Architecture for 3D point clouds. Uses a T-Net to predict affine transformations (canonicalization) for rotation invariance.

Input $\xrightarrow{\text{MLP}} \text{Features} \xrightarrow{\text{Max Pool}} \text{Global Feature} \xrightarrow{\text{MLP}} \text{Output}$

- 16.3 - Gaussian Processes (GPs)

D. (Infinite Width Equivalence (Neal's Theorem)): A single-hidden-layer neural network with infinite width ($m \rightarrow \infty$) and i.i.d. priors on weights converges to a Gaussian Process (GP).

Com. Mechanism: By the Central Limit Theorem, the pre-activations (sums of many independent random variables) become Gaussian.

Com. Result: The network output $f(x)$ is a draw from a GP with mean $\mu(x) = 0$ and a specific kernel $k(x, x')$.

D. (Deep GPs): This equivalence extends to deep networks. The kernel is defined recursively:

$$K_l(x, x') = E[\phi(f_{l-1}(x))\phi(f_{l-1}(x'))]$$

where the expectation is taken over the GP of the previous layer f_{l-1} .

- 16.4 - Statistical Learning Theory

D. (Generalization Error): The gap between performance on training data (empirical risk) and unseen data (expected risk).

$$\text{Gen}(f) = R[f] - R_{\text{emp}}[f]$$

Classical theory (VC-dimension) predicts overfitting for huge models, but DNNs exhibit Double Descent: test error decreases, rises (at interpolation threshold), and then decreases again as width grows.

D. (PAC-Bayesian Bounds): Provides generalization bounds for stochastic classifiers (posterior Q) based on their distance from a prior P (KL-divergence).

$$EQ[R(f)] \leq R_{\text{emp}}(Q) + 2sKL(Q||P) + \ln(2s)$$

Com. This links generalization to the "flatness" of minima: flat minima allow for a posterior Q with high variance (large entropy) that still fits the data, minimizing the KL term.

17 Chain-Rule and Jacobians for Tensors

D. (k-Dimensional Tensor) $T \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_k}$

D. (Tensor Multiplication)

$$\begin{aligned} \mathbf{T} &= \mathbf{P} \times_{\mathbb{R}^{(a+c)}} \mathbf{Q} \times_{\mathbb{R}^{(a+b)}} \mathbf{R} \\ \mathbf{T} &= \mathbf{P} \times_{\mathbb{R}^{(a \times c)}} \mathbf{Q} \times_{\mathbb{R}^{(a \times b)}} \mathbf{R} \times_{\mathbb{R}^{(b \times c)}} \mathbf{S} \end{aligned}$$

where each entry of \mathbf{T} is computed as follows: $T_{i_1, \dots, i_a, k_1, \dots, k_c} := \sum_{j_1, \dots, j_b} P_{i_1, \dots, i_a, j_1, \dots, j_b} Q_{j_1, \dots, j_b, k_1, \dots, k_c}$

Note that this is just the sum of the multiplications of two numbers which are in corresponding locations in \mathbf{P} and \mathbf{Q} . Essentially, it's the

18.1 — Variational Autoencoders (VAEs)

Relation to Autoencoders

Recall, that with autoencoders, we had defined a concatenation of two differentiable (non-linear) mappings $\mathbf{x} \xrightarrow{E} \mathbf{z} \xrightarrow{D} \hat{\mathbf{x}}$ (an encoder E and a decoder D) and trained it with the following loss $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ (approximating identity function) in order to learn some compressed representation \mathbf{z} of \mathbf{x} which just contains the essence of \mathbf{x} according to some meaningful feature-dimensions.

Further, we could then use E to bootstrap a classification model, by first applying E , and then a classification network C and fine-tuning them jointly on the cross-entropy loss.

So the lower-dimensional features \mathbf{z} capture the factors of variation in the data. And we can reconstruct an \mathbf{x} from its compressed representation \mathbf{z} .

Now the question is can we us a similar kind of setup to use new images?

VAEs define an *intractable* density function $p_{\text{model}}(\mathbf{x})$ with a latent \mathbf{z} . Having this latent variable allows us to build a network similar to an autoencoder. A tractable lower bound for the intractable density function is then derived and optimized.

$$p_{\text{model}}(\mathbf{x}) := p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z}) d\mathbf{z}$$

"Decoder" sample from conditional (complex NN)
"latent representation" sample from prior (simple, e.g., MV - Gaussian)

With VAEs we assume that our data is generated from some underlying unobserved representation \mathbf{z} . We first sample \mathbf{z} and then generate some \mathbf{x} from the conditional distribution. Now, we just have to learn the parameters θ that maximize the likelihood of the training data. Unfortunately, we cannot optimize the likelihood of our model for the data directly (as it's intractable).

$$\int_{\mathbf{z}} p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x} | \mathbf{z}) d\mathbf{z} \quad \text{intractable!}, \text{ tractable}$$

Note that different factorizations of the distributions would also be intractable

$$p_{\theta}(\mathbf{z} | \mathbf{x}) = \frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{x})}$$

Now, the solution to this is that in addition to the encoder network modeling $p_{\theta}(\mathbf{x} | \mathbf{z})$, we define an additional encoder network $q_{\phi}(\mathbf{z} | \mathbf{x})$ that approximates $p_{\theta}(\mathbf{z} | \mathbf{x})$.

$$p_{\phi}(\mathbf{z} | \mathbf{x}) \approx p_{\theta}(\mathbf{z} | \mathbf{x}) = \frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{x})}$$

This will allow us to derive a lower bound on the data likelihood that is tractable, which we can optimize.

A common way to define these distributions is as follows: We assume the latent variable to follow a multivariate gaussian (assumed to be a reasonable prior for latent attributes).

Prior: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Enc. Netw. E: models $q_{\phi}(\mathbf{z} | \mathbf{x})$ with params ϕ and maps $\mathbf{x} \xrightarrow{E} (\mu_{\mathbf{z}|\mathbf{x}}, \Sigma_{\mathbf{z}|\mathbf{x}})$

Dec. Netw. D: models $p_{\theta}(\mathbf{x} | \mathbf{z})$ with params θ and maps $\mathbf{z} \xrightarrow{D} (\mu_{\mathbf{x}|\mathbf{z}}, \Sigma_{\mathbf{x}|\mathbf{z}})$

Note that we use both *diagonal* covariance matrices for $\Sigma_{\mathbf{z}|\mathbf{x}}$ and $\Sigma_{\mathbf{x}|\mathbf{z}}$. So the output of both networks are just two vectors (one for mean, other for diagonal).

Com. Encoder and decoder networks are also called "recognition/inference" and "generation" networks.

Now, equipped with our encoder and decoder networks, we can rewrite the data (log) likelihood as follows (note that we omit the product for all points - you'd just have to put a sum over all the instances in front of everything)

$$\begin{aligned} \log(p_{\theta}(\mathbf{x})) &= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z} | \mathbf{x})} [\log(p_{\theta}(\mathbf{x}))] - \log(p_{\theta}(\mathbf{x})) \text{ does not} \\ &= \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{x})} \right) \right] \quad (\text{Bayes Rule}) \\ &= \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_{\theta}(\mathbf{x} | \mathbf{z}) p_{\theta}(\mathbf{z})}{p_{\theta}(\mathbf{z} | \mathbf{x}) q_{\phi}(\mathbf{z} | \mathbf{x})} \right) \right] \quad (\text{Multiplying by 1}) \\ &= \mathbb{E}_{\mathbf{z}} [\log(p_{\theta}(\mathbf{x} | \mathbf{z}))] - \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{q_{\phi}(\mathbf{z} | \mathbf{x})}{p_{\theta}(\mathbf{z})} \right) \right] \\ &= \underbrace{\mathbb{E}_{\mathbf{z}} [\log(p_{\theta}(\mathbf{x} | \mathbf{z}))]}_{(1)} - \underbrace{KL(q_{\phi}(\mathbf{z} | \mathbf{x}), p_{\theta}(\mathbf{z}))}_{(2)} + \underbrace{KL(q_{\phi}(\mathbf{z} | \mathbf{x}), p_{\theta}(\mathbf{z}))}_{\text{Redundant}} \end{aligned}$$

- (1) Decoder network gives us $p_{\theta}(\mathbf{x} | \mathbf{z})$ and we can compute estimates of this term through sampling. (Sampling differentiable through *reparametrization trick*!)
- This term ensures that we reconstruct the data well.
- This term (between Gaussians for encoder and \mathbf{z} prior) has a nice closed-form solution.
- This term ensures that the approximate posterior distribution is close to prior.
- (3) $p_{\theta}(\mathbf{z} | \mathbf{x})$ is intractable (as seen earlier). But we know that the KL-divergence is ≥ 0 .

Now what we have is a *tractable lower bound* \mathcal{L} (so-called *variational lower bound*, or *evidence lower bound* "ELBO") for the likelihood

$$\mathcal{L}(\mathbf{x}, \theta, \phi) \leq \log(p_{\theta}(\mathbf{x}))$$

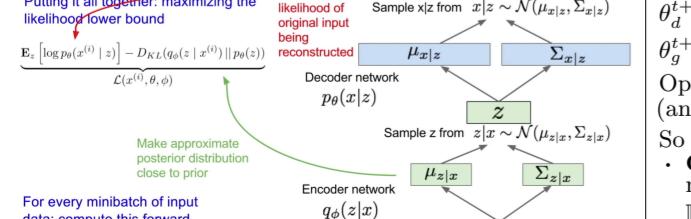
and we can take its gradient to optimize:

$$(\theta^*, \phi^*) = \arg \max_{(\theta, \phi)} \sum_{i=1}^n \mathcal{L}(\mathbf{x}^{(i)}, \theta, \phi).$$

Reparametrization Trick

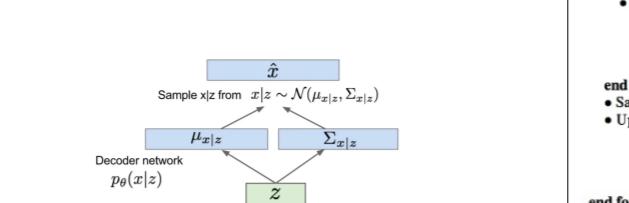
Forward Pass and Backpropagation

Variational Autoencoders



So, first we do the whole backpropagation. And then we just compute the updates to the parameters θ and ϕ via backpropagation.

Generating Data Here we just sample from the prior, and pass it through the decoder network to get the posterior distribution's parameters, then we sample from that one.



for number of training iterations do
for k step do
• Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_{\theta}(\mathbf{z})$.
• Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
• Update the discriminator by ascending its stochastic gradient:
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D_{\theta_d}(\mathbf{x}^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)})))]$$

end for
• Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_{\theta}(\mathbf{z})$.
• Update the generator by ascending its stochastic gradient (improved objective):
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(\mathbf{z}^{(i)})))$$

end for

18.2 — Deep Latent Gaussian Models

Generative Adversarial Networks (GANs)

GANs do not try to model a density function but directly aim to build a function to generated data (implicit generative method). The whole optimization motivated by a game-theoretic approach in a 2-player game. GANs sample from a simple random noise distribution and try to learn a transformation (via a NN) to a data distribution.

D. (Discriminator D) must be a differentiable function parametrized by θ_d

$$\mathbf{x} \xrightarrow{D} P(\mathbf{x} \text{ comes from true data distr.}) \in [0, 1]$$

D. (Generator G) must be a differentiable function parametrized by θ_g

$$\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \xrightarrow{G} \mathbf{x} \quad (\text{one falsified data sample})$$

(one may also use another prior)

\mathbf{z} is sampled from the prior distr. over latent vars (source of randomness)

D tries to make $D(G(\mathbf{z}))$ near 0 (for fake data)

D tries to make $D(\mathbf{x})$ near 1 (\mathbf{x} sampled from true data)

G tries to make $D(G(\mathbf{z}))$ near 1

Com. In some sense G implicitly tries to make $D(\mathbf{x})$ near 0 (since it uses the negative loss of D) see Minimax game VS Non-Saturating game.

Minimax Game: Both D and G try to minimize and maximize the same value function:

$$\begin{aligned} V(G, D) &= \min_G \max_D \mathcal{R}^{(D)} \\ &= \min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(D_{\theta_d}(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})))] \end{aligned}$$

$$\begin{aligned} \mathcal{R}^{(D)} &= -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(D(\mathbf{x}))] - \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\log(1 - D(G(\mathbf{z})))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{generator}}} [\log(1 - D(\mathbf{x}))] \end{aligned}$$

$$\mathcal{R}^{(G)} = -\mathcal{R}^{(D)}$$

- The loss $\mathcal{R}^{(D)}$ is simply the cross-entropy between D 's predictions and the correct labels in the binary classification task (real/fake)
- The equilibrium of this game is saddle point of the discriminator loss
- If we look for this equilibrium the whole procedure resembles minimizing the Jensen-Shannon divergence between the true data distribution and the generator distribution.
- So G minimizes the log-probability of D being correct

What is the solution $D(\mathbf{x})$ in terms of p_{data} and $p_{\text{generator}}$ at the equilibrium?

In the equilibrium it must hold that the gradient of the discriminator is zero, because the discriminator otherwise would improve (thus change) itself.

$$\begin{aligned} \frac{\partial \mathcal{R}^{(D)}}{\partial D(\mathbf{x})} &= -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{D(\mathbf{x})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{generator}}} \left[\frac{1}{1 - D(\mathbf{x})} \right] \stackrel{!}{=} 0 \\ &\Rightarrow \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[\frac{1}{D(\mathbf{x})} \right] = \mathbb{E}_{\mathbf{x} \sim p_{\text{generator}}} \left[\frac{1}{1 - D(\mathbf{x})} \right] \\ &\Leftrightarrow \int_{\mathbf{x}} \frac{p_{\text{data}}(\mathbf{x})}{D(\mathbf{x})} d\mathbf{x} = \int_{\mathbf{x}} p_{\text{generator}}(\mathbf{x}) \frac{1}{1 - D(\mathbf{x})} d\mathbf{x} \\ &\text{Recall } \mathbf{x} = \text{operator}(\mathbf{z}). \text{ Using the inverse of the operator, we can get rid of it and the solution constraints on the optimal } D(\mathbf{x}) \text{ still remain the same.} \\ &\Leftrightarrow p_{\text{data}}(\mathbf{x}) \frac{1}{D(\mathbf{x})} = p_{\text{generator}}(\mathbf{x}) \frac{1}{1 - D(\mathbf{x})} \end{aligned}$$

Then we get the following stationarity condition: The optimal $D(\mathbf{x})$ for any $p_{\text{data}}(\mathbf{x})$ and $p_{\text{generator}}(\mathbf{x})$ is always

$$D(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{generator}}(\mathbf{x})}$$

Note that by stationarity condition we mean that this must hold in the Nash equilibrium (where both players stop adapting themselves).

Estimating this ratio using supervised learning is the key approximation mechanism used by GANs.

What assumptions are needed to obtain this solution?

We need to assume that both densities are nonzero everywhere. If we don't make this assumption then there's this issue that the discriminator's input space might never be sampled during its training process. Then these points would have an undefined behaviour since they're never trained.

Training Procedure: Use SGD-like algorithm of choice (ADAM) on two minibatches simultaneously. At each iteration, we choose:

- a minibatch of true data samples
- a minibatch of noise vectors to produce minibatch generated samples

Then compute both losses and perform gradient updates.

$$\theta_d^{t+1} \leftarrow \theta_d^t - \eta_t \nabla_{\theta_d} \mathcal{R}^{(D)}(\theta_d)$$

$$\theta_g^{t+1} \leftarrow \theta_g^t - \eta_t \nabla_{\theta_g} \mathcal{R}^{(G)}(\theta_g)$$

Optional: run $k \geq 1$ update steps of D for every iteration (and only 1 update step for G).

So we alternate between

$$\max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(D_{\theta_d}(\mathbf{x}))]$$

$$\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})))]$$

$$\max_{\theta_g} \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})))]$$

Note that this training algorithm uses a heuristically motivated loss (that is a bit different) for the generator to have better gradients when the discriminator is good: