

1 Probability

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

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

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

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

Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X, Y|Z) = P(X|Z)P(Y|Z)$

Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X|Y, Z) = P(X|Z)$

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

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

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

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

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

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

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

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

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

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

$f_{\alpha}(y) = \frac{1}{\alpha} \delta(y - \frac{x}{\alpha})$

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

$B(p) = p^t + (1-p)^{1-t}$ $\mathcal{N}(p, \sigma) = \exp(pt + \frac{1}{2}\sigma^2 t^2)$

$\text{Bin}(n, p) = (pe + (1-p))^n$ $\text{Gam}(\alpha, \beta) = (\frac{1}{\alpha} \beta^{\frac{1}{\alpha}})^{\alpha}$ for $t < 1/\beta$

$\text{Pois}(\lambda) = e^{\lambda(e^{t-1})}$

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

T. Let X, Y be indep., then the p.d.f. of $Z = X + Y$ is the conv. of the p.d.f. of X and Y : $f_Z(z) = \int_R f_X(t) f_Y(z-t) dt = \int_R f_X(z-t) f_Y(t) dt$

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

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

T. $\mathcal{P}(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}) = \mathcal{N}(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} | \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix})$

$\mathbf{a}_1, \mathbf{u}_1 \in \mathbb{R}^c, \Sigma_{11} \in \mathbb{R}^{c \times c}$ p.s.d., $\Sigma_{12} \in \mathbb{R}^{c \times f}$ p.s.d., $\mathbf{a}_2, \mathbf{u}_2 \in \mathbb{R}^f, \Sigma_{22} \in \mathbb{R}^{f \times c}$ p.s.d.

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

T. (Chebyshev) Let X be a rv with $\mathbb{E}[X] = \mu$ and variance $\text{Var}[X] = \sigma^2 < \infty$. Then for any $\epsilon > 0$, we have $P(|X - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$.

2 Analysis

Log-Trick (Identity): $\nabla_\theta [p_\theta(\mathbf{x})] = p_\theta(\mathbf{x}) \nabla_\theta [\log(p_\theta(\mathbf{x}))]$

T. (Cauchy-Schwarz) $\forall \mathbf{u}, \mathbf{v} \in V: \langle \mathbf{u}, \mathbf{v} \rangle \leq \|\mathbf{u}\| \|\mathbf{v}\|$

$\mathbf{u}, \mathbf{v} \in V: 0 \leq \langle \mathbf{u}, \mathbf{v} \rangle \leq \|\mathbf{u}\| \|\mathbf{v}\|$

Special cases: $\langle \sum_i x_i y_i \rangle \leq (\sum_i x_i^2)(\sum_i y_i^2)$. 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

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

$f'(x) \geq 0$

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 of the loss function are nearly spherical. Gradient Descent converges quickly and directly toward the minimum.
- Ill-conditioned ($\kappa \gg 1$): 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.

T. (Taylor-Lagrange Formula)

$f(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k + \int_{x_0}^x f^{(n+1)}(t) \frac{(x-t)^n}{n!} dt$

T. (Jensen) f convex/concave, $\forall i: \lambda_i \geq 0, \sum_i \lambda_i = 1$

Special case: $\mathbb{E}[f(X)] \leq \mathbb{E}[f(X)]$.

D. (*L*-Lipschitz Continuous Function) Given two metric spaces (X, d_X) and (Y, d_Y) , a function $f: X \rightarrow Y$ is called Lipschitz continuous, if there exists a real constant $L \in \mathbb{R}_0^+$ (Lipschitz constant), such that $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n: \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\| \leq L \cdot \|x_1 - x_2\|$.

Com. If the objective function is L -smooth a step size of $\eta = 1/L$ guarantees convergence.

D. (Lagrangian Formulation) of $\arg \max_{x,y} f(x, y)$ s.t. $g(x, y) = c: \mathcal{L}(x, y) = f(x, y) - \gamma(g(x, y) - c)$

D. (PL Condition) A differentiable function $f(x)$ with global minimum f^* satisfies the μ -Polyak-Lojasiewicz (PL) condition if there exists a constant $\mu > 0$ such that for all x :

$\frac{1}{2} \|\nabla f(x)\|^2 \geq \mu(f(x) - f^*)$

Significance:

$JSD(P, Q) = \frac{1}{2} KL(P, M) + \frac{1}{2} KL(Q, M) \in [0, \log(n)]$ $M = \frac{1}{2}(P + Q)$

C. The JSD is symmetric!

Com. The JSD is a symmetrized and smoothed version of the KL-divergence.

6 Activation Functions

Activation functions are non-linear functions that are applied element-wise to the output of a linear function.

D. (Tanh)

$\tanh(\mathbf{x}) = \frac{e^{\mathbf{x}} - e^{-\mathbf{x}}}{e^{\mathbf{x}} + e^{-\mathbf{x}}} \quad \tanh'(\mathbf{x}) = 1 - \tanh^2(\mathbf{x})$

D. (ReLU)

$\text{ReLU}(\mathbf{x}) = \max(\mathbf{x}, 0) \quad \text{ReLU}'(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} \geq 0 \\ 0 & \text{otherwise} \end{cases}$

D. (Sigmoid/Logistic)

$\sigma(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}}} \in (0, 1) \quad \frac{\partial \sigma(\mathbf{x})}{\partial \mathbf{x}} = \sigma(\mathbf{x}) \odot (1 - \sigma(\mathbf{x}))$

$\sigma^{-1}(y) = \ln\left(\frac{y}{1-y}\right) \quad \nabla_{\mathbf{x}} \sigma(\mathbf{x}) = \mathbf{J}_\sigma(\mathbf{x}) = \text{diag}(\sigma(\mathbf{x}) \odot (1 - \sigma(\mathbf{x})))$

$\sigma'(\mathbf{x}) = \frac{1}{4} \tanh'\left(\frac{1}{2}\mathbf{x}\right) = \frac{1}{4}(1 - \tanh^2(\frac{1}{2}\mathbf{x})) \quad \sigma(-\mathbf{x}) = 1 - \sigma(\mathbf{x})$

3 Linear Algebra

Kernels are positive semi-definite matrices.

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

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

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.

Otherwise: indefinite.

D. (Trace) of $A \in \mathbb{R}^{n \times n}$ is $\text{Tr}(A) = \sum_{i=1}^n a_{ii}$.

Properties:

- $\text{Tr}(A) = \sum_i \lambda_i$ (sum of eigenvalues).
- Cyclic property: $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$.
- Linear: $\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$ and $\text{Tr}(cA) = c\text{Tr}(A)$.
- $\text{Tr}(A) = \text{Tr}(A^T)$.

D. (Frobenius Norm $\|\cdot\|_F$) The square root of the sum of the absolute squares of its elements.

$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}$

Properties:

- Relation to Trace: $\|\mathbf{A}\|_F = \sqrt{\text{Tr}(A^T A)}$.
- Invariant under orthogonal rotations: $\|\mathbf{Q}\mathbf{A}\|_F = \|\mathbf{A}\|_F$ for orthogonal \mathbf{Q} .
- Relation to Singular Values: $\|\mathbf{A}\|_F = \sqrt{\sum_i \sigma_i^2}$.

7 Connectionism

McCulloch & Pitts (1943): MP-Neuron. Abstract model of neurons as linear threshold units. Inputs $\mathbf{x} \in \{0, 1\}^n$, synapses $\sigma \in \{-1, 1\}^n$, $f(\mathbf{x}) = \prod_i \sigma_i x_i \geq 0$. No learning (fixed weights).

Perceptron (Rosenblatt 1958): Pattern recognition. $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$. Update rule (if misclassified): $\mathbf{w}_{old} \leftarrow \mathbf{w}_{old} + y_i \mathbf{x}_i$, $b_{new} \leftarrow b_{old} + y_i$. Cannot learn the XOR function.

T. (Novikoff) Guaranteed convergence if data is linearly separable.

T. (Goodfellow, 2013) Maxout networks with two maxout units that are applied to m linear functions are universal function approximators.

4 Derivatives

4.1 Numerator and Denominator Convention

Jacobian Layout Convention We use the denominator-layout. For a vector-valued function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ we define

$\left(\frac{\partial f}{\partial \mathbf{x}}\right)_{ij} := \frac{\partial f_j}{\partial x_i}, \quad \frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{m \times n}$

Hence, gradients of scalar-valued functions are column vectors, and the chain rule takes the form

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

Remark. Sometimes the numerator-layout is used, where the Jacobian is defined as $(J_{ij})_{ij} = \frac{\partial f_i}{\partial x_j} / \frac{\partial x_j}{\partial x_j}$ in $\mathbb{R}^{m \times n}$. The two conventions are related by transposition.

4.2 Scalar-by-Vector

Denominator Convection $\frac{\partial}{\partial \mathbf{x}} [u(\mathbf{x})] = u(\mathbf{x}) \frac{\partial}{\partial \mathbf{x}}$

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

Hopfield Networks: Associative Memory. Hebbian Learning: "Neurons that fire together, wire together". $w_{ij} \propto \sum_i x_i x_j$. Energy Min.: $E(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \mathbf{W} \mathbf{x}$. Note that $\mathbf{h}^T = \mathbf{x}$.

PD (Rumelhart et al. 1986): Parallel Distributed Processing. Introduction of Backpropagation (Generalized Delta Rule). Differentiable activations allow δ propagation to hidden layers.

S. Regression

D. (Linear Regression): $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$

Proper initialization is crucial to prevent vanishing or exploding gradients at the start of training.

D. (Feedforward Network) is a set of computational units arranged in a DAG (layer-wise processing)

F. $F = F^L \circ \dots \circ F^1$ where each layer $L \in \{1, \dots, L\}$ is a composition of the following functions

$F^l: \mathbb{R}^{m_{l-1}} \rightarrow \mathbb{R}^{m_l} \quad F^l = \sigma^l \circ \overline{F}^l$ where $\overline{F}^l: \mathbb{R}^{m_{l-1}} \rightarrow \mathbb{R}^{m_l}$ is the linear function in layer l

$\overline{F}^l(\mathbf{h}^{l-1}) = \mathbf{W}^l \mathbf{h}^{l-1} + \mathbf{b}_l \quad \mathbf{W}^l \in \mathbb{R}^{m_{l-1} \times m_l}$ and $\sigma^l: \mathbb{R}^{m_{l-1}} \rightarrow \mathbb{R}^{m_l}$ element-wise non-linearity at layer l .

Note that $\mathbf{h}^0 = \mathbf{x}$.

T. Feedforward Networks are invariant under permutations of units. The units within a hidden layer are interchangeable (along with their corresponding weights).

T. The units of feedforward networks are invariant along directions orthogonal to the weight vector $\mathbf{F}[\mathbf{w}, \mathbf{b}](\mathbf{x}) = [\mathbf{F}[\mathbf{w}, \mathbf{b}](\mathbf{x} + \delta\mathbf{x}) - \mathbf{F}[\mathbf{w}, \mathbf{b}](\mathbf{x})] / \delta\mathbf{x}$.

D. (Maxout) is just the max non-linearity applied to k groups of linear functions. So the input $[1 : d]$ (of the previous layer) is partitioned into k sets A_1, \dots, A_k , and then we define the activations $G_j(\mathbf{x})$ for $j \in \{1, \dots, k\}$

$G_j(\mathbf{x}) = \max_{i \in A_j} [\mathbf{w}_i^T \mathbf{x} + b_i] \quad (i \in \{1, \dots, d\})$

D. (Scale and Shift): The network learns parameters γ and β to restore representation power (allows the network to undo the normalization if needed):

$y = \gamma \hat{x} + \beta$

Benefits:

- Effectively increases the effective batch size and reduces the variance of the estimate.
- Allows the use of larger learning rates (longer steps) while still converging to the optimal solution asymptotically.
- Often achieves the optimal convergence rate of $O(1/t)$ for convex problems.

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:

- Normalization: Calculate batch mean μ_B and variance σ_B^2 to normalize input \mathbf{x} :

$\hat{x} = \frac{x - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$

- Scale and Shift: The network learns parameters γ and β to restore representation power (allows the network to undo the normalization if needed):

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.

2.2 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 doing "arrays" (discrete-time samples).

D. (Translation- (or Shift-) Invariant) A transform T is translation (or shift) invariant, if for any f and scalar τ ,

$f_\tau := f(t - \tau)$ (Def. shift operator f_τ)

$(Tf_\tau)(t) = (Tf)(t - \tau)$ (commuting of operators holds)

So, an operator T is shift-invariant iff it commutes with the shift operator s_Δ .

$\forall f: (T(s_\Delta f)) = (s_\Delta(Tf))$.

So, the commutative diagram for this is:

$$\begin{array}{ccc} f & \xrightarrow{s_\Delta} & \Delta f = f_\Delta \\ \downarrow & & \downarrow T \\ Tf & \xrightarrow{s_\Delta} & \Delta(Tf) = T(f_\Delta) \end{array}$$

Convolutions are translation and shift invariant.

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

10.3 Gradient Descent

D. (Gradient Descent (GD)) Iteratively moves parameters θ in the direction of the negative gradient of the loss function $J(\theta)$.

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

Com. In Stochastic GD (SGD), the gradient is approximated using a single sample (or mini-batch) to introduce noise and escape local minima. Although, the gradient is unbiased it adds variance, this can help to escape local minima and saddle points.

Com. Gradient Flow can be seen as the numerical integration of the continuous-time ordinary differential equation (ODE) $\dot{x} = -\nabla f(x)$.

D. (Polak-Ribiere (Averaged SGD)) Instead of using the final parameter

12.6 Efficient Comp. of Convolutional Activities
A naive way to compute the convolution of a signal of length n and a kernel of length m gives an effort of $\mathcal{O}(m \cdot n)$. A faster way is to transform both with the FFT and then just do element-wise multiplication (effort: $\mathcal{O}(n \log n)$). However, this is rarely done in CNNs as the filters usually are small ($m \ll n, m \approx \log(n)$).

12.7 Typical Convolutional Layer Stages

A typical setup of a convolutional layer is as follows:

1. Convolution stage: affine transform
2. Detector stage: nonlinearity (e.g., ReLU)
3. Pooling stage: locally combine activities in some way (max, avg, ...)
4. Locality of the item that activated the neurons isn't too important, further we profit from dimensionality reduction. Alternative: deconvolution with stride. Another thing that turns out to be so is that most of the kernels that are learned resemble a low-pass filter. Hence, when we sub-sample the images most of the information is still contained.

12.8 Pooling

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

D. (Max-Pooling)

• 1D: $x_{ij}^{\text{max}} = \max\{x_{i+k} \mid 0 \leq k < r\}$

• 2D: $x_{ij}^{\text{max}} = \max\{x_{i+k, j+l} \mid 0 \leq k, l < r\}$

12.9 Sub-Sampling (aka "Strides")

Often, it is desirable to reduce the size of the feature maps. That's why sub-sampling was introduced.

D. (Sub-Sampling)

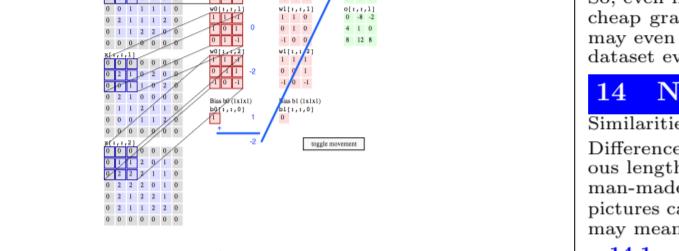
Hereby the temporal/spatial resolution is reduced.

12.10 Channels

Ex: Here we have

• an input signal that is 2D with 3 channels ($7 \times 7 \times 3$) (image \times 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 \times num convolutions)



Usually we convolve over all of the channels together, such that each convolution has the information of all channels at its disposition and the order of the channels hence doesn't matter.

12.11 CNNs in Computer Vision

So the typical use of convolution that we have in vision is: a sequence of convolutions.

1. that reduce the spatial dimensions (sub-sampling)

The deeper we go in the network, we transform the spatial information into a semantic representation. Usually, most of the parameters lie in the fully connected layers

12.12 Famous CNN Architectures

12.12.1 LeNet, 1989

MNIST, 2 Convolutional Layers + 2 Fully-connected layers

12.12.2 LeNet5

MNIST, 3 Convolutional Layers (with max-pool subsampling) + 1 Fully connected layer

12.12.3 AlexNet

ImageNet: similar to LeNet5, just deeper and using GPU (performance breakthrough)

12.12.4 Inception Module

Now, a problem that arose with this ever deeper and deeper networks was that the filters at every layer were getting longer and longer and lots of their coefficients were becoming zero (so no information flowing through). So, Arora et al. came up with the idea of an inception module.

What this inception module does is just taking all the channels for one element in the space, and reduces their dimensionality. Such that we don't get too deep channels, and also compress the information (learning the low-dimensional manifold).

This is what gave rise to the inception module:

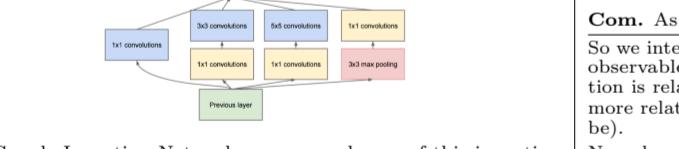
D. (Dimension Reduction)

m channels of a $1 \times 1 \times k$ convolution $m \leq k$:

$$x^T i j = \sigma(W x_{i,j}), \quad W \in \mathbb{R}^{m \times k}.$$

So it uses a 1×1 filter over the k input channels (which is actually not a convolution), aka "network within a network".

12.12.5 Google Inception Network



The Google Inception Network uses many layers of this inception module along with some other tricks

- dimensionality reduction through the inception modules
- convolution at various sizes, as different filter sizes turned out to be useful
- a max-pooling of the previous layer, and a dimensionality reduction of the result
- 1×1 convs for dimension reduction before convolving with larger kernels
- then these informations are passed to the next layer
- gradient shortcuts: connect softmax layer at intermediate stages to have the gradient flow until the beginnings of the network
- decomposition of convolution kernels for computational performance
- all-in-all the dimensionality reductions improved the efficiency

12.13 Networks Similar to CNNs

D. (Locally Connected Network) A locally connected network has the same connections that a CNN would have, however, the parameters are not shared. So the output nodes do not connect to all nodes, just to a set of input nodes that are considered "near" (locally connected).

12.14 Comparison of #Parameters (CNNs, FC, LC)

Ex: Input image $m \times n \times c$ (c = number of channels)

K convolution kernels: $p \times q$ (valid padding and stride 1)

output dimensions: $(m - p + 1) \times (n - q + 1) \times K$

#parameters CNN: $K(pq + 1)$

#parameters of fully-conn.: N^2 with same number of outputs as CNN: $mnc(m - p + 1)(n - q + 1) + 1K$

#parameters of locally-conn.: NN with same connections as CNN: $pqc(m - p + 1)(n - q + 1) + 1K$

Ex: Assume we have an $m \times n$ image (with one channel).

And we convolve it with a filter $(2p + 1) \times (2q + 1)$

Then the convolved image has dimensions (assuming stride 1)

• valid padding (only where it's defined): $(m - 2p) \times (n - 2q)$

• same padding (extend image with constant): $m \times n$ where the extended image has size $(m + 2p) \times (n + 2q)$.

13 Optimization

13.1 Objectives as Expectations

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

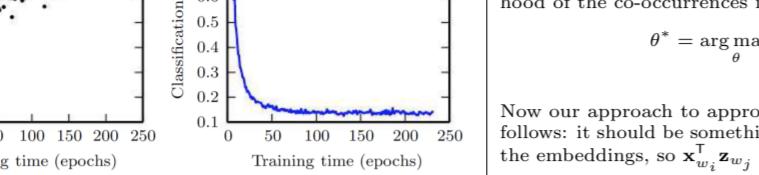
T: We have the following chain of inclusions for functions over a closed and bounded (i.e., compact) subset of the real line.

Continuously differentiable \subseteq Lipschitz continuous \subseteq (Uniformly) continuous
If we use Nesterov acceleration (in the general case), then we get a polynomial convergence rate of $O(t^{-2})$.

Com: The trick used in the Nesterov approach is momentum.

13.2 Optimization Challenges in NNs: Curvatures

A typical setup of a convolutional layer is as follows:



13.2.1 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 learning rate, then s.t.

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

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 the idea is to use two different latent vectors \mathbf{x}_w and \mathbf{z}_w per word (to allow for the asymmetry for the conditional probability)

$$\theta = (\mathbf{x}_w, \mathbf{z}_w)_{w \in V}$$

• \mathbf{x}_w is used to predict w 's conditional probability, and

• \mathbf{z}_w is used to use w as an evidence in the cond. prob.

Note that \mathbf{C} is actually symmetric (we're computing the joint probabilities), but the probabilities that we're computing are asymmetric (of conditional probability).

Problem: Note however that it's too expensive to compute the partition function as we have to do a full sum over V (can be $\sim 10^5$ up to $\sim 10^7$). And we'd have to do this every time we pass a new batch-sample (w_{t+1}, w_t) through the network.

Bright Idea of skip-grams: instead of computing the partition function, then the part of determining $P_\theta(w_i | w_j)$ into a classification problem (logistic regression). So, we create a classifier that determines the co-occurring likelihood of the words on a scale from 0 to 1.

For this reason we'll introduce the following function

$$D_{w_i, w_j} = \begin{cases} 1, & \text{if } (i, j) \in C_R, \\ 0, & \text{if } (i, j) \notin C_R. \end{cases}$$

and we'll squash the dot-product to something between 0 and 1 to make it serve as a probability

$$P_\theta(D_{w_i, w_j}) \cong P_\theta(D_{w_i, w_j} = 1 | \mathbf{x}_{w_i}, \mathbf{z}_{w_j}) = \sigma(\mathbf{x}_{w_i}^\top \mathbf{z}_{w_j})$$

and the opposite event is given by:

$$P_\theta(D_{w_i, w_j} = 0 | \mathbf{x}_{w_i}, \mathbf{z}_{w_j}) = 1 - \sigma(\mathbf{x}_{w_i}^\top \mathbf{z}_{w_j}) = \sigma(-\mathbf{x}_{w_i}^\top \mathbf{z}_{w_j})$$

Further, instead of just maximizing the likelihood over the dataset of co-occurrences D , we'll also maximize the log likelihood over a dataset of non-co-occurrences \bar{D} (negative samples). So, we'll have the following maximization problem

$$\theta^* = \arg \max_{\theta} \sum_{(i, j) \in D} \log(P_\theta(D_{w_i, w_j} = 1 | \mathbf{x}_{w_i}, \mathbf{z}_{w_j})) + \sum_{(i, j) \in \bar{D}} \log(1 - P_\theta(D_{w_i, w_j} = 1 | \mathbf{x}_{w_i}, \mathbf{z}_{w_j}))$$

Given an observation sequence $\mathbf{x}^1, \dots, \mathbf{x}^T$. We want to identify the hidden activities \mathbf{h}^t of the state of a dynamical system. The discrete time evolution of the hidden state sequence is expressed as a HMM with a non-linearity

$$\theta = (U, \mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c})$$

This ensures that the gradient norm is never greater than γ_{\max} . However, when having vanishing gradients over time, then this means that the RNN is forgetting the past after a few timesteps, and usually then the RNN is not performing very well. This is harder to fix, and we'll see later how this is solved with LSTM.

14.1.1 Bi-Linear Models

The first thing that we could do is to use an information theoretic quantity: the so-called mutual information. The mutual information is described in information theory as how much information one random variable has about another random variable. If two variables are independent, then the mutual information will be zero. So, if we put two words nearby, it's because they have to be related somehow in the meaning of the sentence. Hence, we expect them to have a larger mutual information.

14.1.2 Dynamic CNNs

Kalchbrenner et al. suggested Dynamic CNNs in 2014 (as an alternative to ConvNet). They are exactly the same as ConvNets except for one thing: before doing the max-pooling over time (or for a fixed size representation), they do a dynamic max-pooling (dynamic since it depends on the input size) over the sentence and another convolution.

14.2 Recurrent Networks (RRNs)

Disadvantage of CNNs: need to pick right convolution size, too small: no context, too large: a lot of data needed, anyways: loss of memory at some point.

Advantage of RNNs: capture better the time component, lossy memorization of past in hidden state.

Given an observation sequence $\mathbf{x}^1, \dots, \mathbf{x}^T$. We want to identify the hidden activities \mathbf{h}^t of the state of a dynamical system. The discrete time evolution of the hidden state sequence is expressed as a HMM with a non-linearity

$$\theta = (U, \mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c})$$

This ensures that the gradient norm is never greater than γ_{\max} . However, when having vanishing gradients over time, then this means that the RNN is forgetting the past after a few timesteps, and usually then the RNN is not performing very well. This is harder to fix, and we'll see later how this is solved with LSTM.

14.3.1 Backprop Over Time

For multi-output loss

$$\frac{\partial L}{\partial \theta} = \sum_{t=1}^T \frac{\partial L_t}{\partial \theta} \frac{\partial y_t}{\partial \theta} = \sum_{t=1}^T \frac{\partial h_t}{\partial \theta} \frac{\partial y_t}{\partial h_t}$$

There are two scenarios for producing outputs

- 1. Only one output at the end:

$$h^T \mapsto \mathbf{H}(h^T; \theta) = \mathbf{y}$$

And then we just pass this \mathbf{y} to the loss \mathcal{R} .

- 2. Output a prediction at every timestep: $\mathbf{y}^1, \dots, \mathbf{y}^T$. And then use an additive loss function

$$\mathcal{R}(\mathbf{y}^1, \dots, \mathbf{y}^T) = \sum_{t=1}^T \mathcal{R}(\mathbf{H}(h^t; \theta))$$

So, additionally, we're define a reverse order sequence

$$\math$$

The probability that we want to determine is
 $P(\mathbf{y}^1, \dots, \mathbf{y}^T | \mathbf{x}_1, \dots, \mathbf{x}^T, F(\mathbf{x}^T))$.

The issue that we have here is that T_x and T_y have variable lengths, and the difference between the two lengths is not always the same. So it's very hard to match one sequence to another. Now, sequence learning will compute a function
 $F(\mathbf{x}^1, \dots, \mathbf{x}^{T_x}) = \text{"thought vector"}$

which will be a vector which will have all the information that we need from the input sequence to compute the output sequence. This F is the so-called "thought vector" (Hinton). So F will be computed via an LSTM.

To produce the output sequence we'll use another LSTM that takes as input the thought vector F plus the output that we'll be producing (output feedback).

- How to make the RNN Encoder/Decoder Work?

The following things were discovered by Sutskever, Vinals & Le in 2011:

- Use Deep LSTMs (multiple layers, e.g., 4)
- Use different RNNs for encoding and decoding
- Apply beam search for decoding
- Reverse the order of the source sequence
- Ensemble-ing

For a machine translation task this gave state-of-the-art results on WMT benchmarks. However, traditional approaches use *sentence alignment models*. We still don't know what is the equivalent in a neural architecture.

- 15.3.1 — Seq2Seq with Attention

The issue with the encoder-decoder architecture is that if we're translating a very long sequence, it might have the issue that suddenly we have to store the entire sequence in a single vector. But when we as humans translate we translate small parts into small parts. In order to understand this better let's have a look at a concrete example. Let's say that we want to translate the following sentence from English to French:

- bi-directionality (it's good to know future and past context)
- selected words/selected based on attention
- sizes of sentences might not be the same
- outputted words might have slightly different order
- Note that we don't have dependencies that are out of order we can use the CTC approach.

- 15.4 — Recurrent Networks

Good to process tree-structure, e.g., from a parser (more depth efficient $O(\log(n))$). Gives a single output at the root.

$F: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$

$\mathbf{h}^n = F(\mathbf{h}^{n-1}, \mathbf{h}^{n-1})$

16 Unsupervised Learning

- 16.1 — Density Estimation

D. (Density Estimation) is used to learn the distribution of the data. Classically, we use a *parametric family of densities* $\{p_\theta | \theta \in \Theta\}$ to describe the set of densities that we may model. Usually, the parameters are estimated with MLE (expectation w.r.t. the empirical distribution)

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{emp}}} [\log(p_\theta(\mathbf{x}))].$$

However, real data is rarely gaussian, laplacian, ... e.g., images. So the fact that in general we cannot solve for p_θ for a parametric function makes this task quite complicated.

So when using a *prescribed model* p_θ we have to:

- ensure that p_θ defines a proper density:

$$\int p_\theta(\mathbf{x}) d\mathbf{x} = 1.$$

and to be able to evaluate the density p_θ at *various sample points* \mathbf{x} :

- this may be trivial for models such as exponential families (simple formulas)
- but impractical for complex models (Markov networks, DNNs)

Now, the question is what strategies can we use for more complex models.

A typical example for an non-parametric and unnormalized model is kernel-density estimation.

D. (Kernel Density Estimator) Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be a sample, and k a kernel with bandwidth $h > 0$ then the estimator is defined as:

$$\bar{p}_\theta(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k_h(\mathbf{x} - \mathbf{x}_i) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right).$$

The problem with this is that the rate of convergence is $\log(\log(n))$

- this is extremely painfully slow. This is just a guarantee in general when we know nothing about our density.

An alternative is to use *unnormalized models* (non-parametric: the number of parameters depends on dataset size). These then represent improper density functions:

$$\bar{p}_\theta(\mathbf{x}) = \frac{c_\theta}{\text{represented}} \cdot \frac{p_\theta(\mathbf{x})}{\text{unknown normalized}}.$$

Finding the normalization constant c_θ might be really complicated, so we can only evaluate relative probabilities. Further, here we cannot use the log-likelihood, because scaling up p_θ leads to an unbounded likelihood.

So the question still is: is there an alternative *estimation method* for the log-likelihood?

What we do in practice is we do not look for the exact p_θ , but we look for properties of p_θ . In many cases these properties depend on our prior knowledge of p_θ . We need to understand what the problem is in order to put the prior knowledge into the model that we want to do. This was already important in supervised learning (e.g., CNNs with several layers for images), and is even more important in unsupervised learning. We have to do the same thing there without knowing what our final goal is.

Finally, Hyvärinen came up with the following idea in 2005. He asked himself whether there's an *operator* that we can apply to \bar{p}_θ that does not depend on normalization. - The answer was yes! Instead of estimating p_θ , we estimate $\log p_\theta$.

D. (Score Matching (Hyvärinen 2005))

$$\psi_\theta := \nabla_{\mathbf{x}} \log \bar{p}_\theta, \quad \psi = \nabla_{\mathbf{x}} \log p$$

Minimize the criterion

$$J(\theta) = \mathbb{E}[\|\psi_\theta - \psi\|^2]$$

or equivalently (by eliminating ψ by integration by parts)

$$J(\theta) = \mathbb{E} \left[\sum_i \partial_i \psi_{\theta,i} - \frac{1}{2} \psi_{\theta,i}^2 \right].$$

This expectation can be approximated by sampling.

The main problem with this is that it assumes that the two normalization constants are the same!

- 16.2 — Factor Analysis

Latent Variable Analysis provides a generic way of defining probabilistic, i.e., *generative models* - the so-called *latent variable models*. They usually work as follows

1. Define a *latent variable* \mathbf{z} , with a distribution $p(\mathbf{z})$
2. Define *conditional models* for the observables \mathbf{x} conditioned on the latent variable: $p(\mathbf{x} | \mathbf{z})$
3. Construct the *observed data model* by integrating/summing out the latent variables

$$p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) \mu(d\mathbf{z}) = \left(\int p(\mathbf{z}) \right) \left(\sum_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}) \right), \quad \mu = \text{Lebesgue}$$

Ex. (Gaussian Mixture Models GMMs)

$\mathbf{z} \in \{1, \dots, K\}$, $p(\mathbf{z})$ = mixing proportions

$p(\mathbf{x} | \mathbf{z})$: conditional densities (Gaussians for GMMs)

The idea of latent variable models is very similar to the one of autoencoders.

The idea is to have some

- $\mathbf{x} \in \mathbb{R}^d$
- and we want to embed it into \mathbb{R}^k ($k \ll d$)
- so we'll use $\mathbf{z} \in \mathbb{R}^k$ (latent-space)
- and look at the conditional probabilities $p(\mathbf{x} | \mathbf{z})$ for some \mathbf{x}

Depending on whether \mathbf{z} is continuous (e.g., as with PCA) or discrete random variable (e.g., GMMs) we'll be using the Lebesgue integral or counting to integrate/sum it out.

A typical approach to for latent variable models is *linear factor analysis*

- Refresher on MGFs and Gaussians

D. (Moment Generating Function (MGF)) The MGF M_X of a random vector $\mathbf{X} \in \mathbb{R}^n$ is defined as

$$M_X: \mathbb{R}^n \rightarrow \mathbb{R} \\ t \mapsto \mathbb{E}_{\mathbf{x}} [e^{t^T \mathbf{x}}].$$

The reason M_X is called *moment generating function* is because it represents the *moments* of \mathbf{x} in the following way. Let $k_1, \dots, k_n \in \mathbb{N}$, then

$$\mathbb{E}_{\mathbf{x}} \left[x_1^{k_1} x_2^{k_2} \cdots x_n^{k_n} \right] = \frac{\partial^k}{\partial t_1^{k_1} \partial t_2^{k_2} \cdots \partial t_n^{k_n}} M_X|_{t=0}.$$

T. (Uniqueness Theorem) If M_X and M_Y exist for the RVs \mathbf{X} and \mathbf{Y} and $M_X = M_Y$ then $\forall t: P(\mathbf{X} = t) = P(\mathbf{Y} = t)$ (distributions are the same).

Now every distribution has its unique kind of MGF form. Hence, MGFs can be very useful to deal with sums of *i.i.d.* random variables:

T. If \mathbf{X}, \mathbf{Y} are i.i.d. then $M_{\mathbf{X}+\mathbf{Y}} = M_{\mathbf{X}} \cdot M_{\mathbf{Y}}$.

- 16.3 — Latent Variable Models

- 16.3.1 — DeFinetti's Theorem

There is another way of looking at latent variable models which is by the DeFinetti's theorem from the 1930s. This is one of the foundations of Bayesian probability (although there is nothing Bayesian in this theorem).

T. (DeFinetti's Theorem) For *exchangeable* data (order of dataset doesn't matter and they come from the same distribution), we can decompose the data by a *latent variable model*

$$P(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \prod_{i=1}^N p_\theta(\mathbf{x}_i | \mathbf{z}) p_\theta(\mathbf{z}) d\mathbf{z}.$$

We expect that those hidden variables are: interpretable and actionable and even show causal relations.

Later we'll put our Bayesian priors into the distributions $P(\mathbf{z})$ and we then hope that the latent structure will tell us something about the data that we didn't know before.

- 16.3.2 — Latent Variable Models

Classically we define complex models via the *marginalization* of a *latent variable model*

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

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 $\mathbf{x} \xrightarrow{\phi} \mathbf{z} | \mathbf{x}$

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

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).

Conn. 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)

$$\log(p_\theta(\mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} [\log(p_\theta(\mathbf{x}))]$$

(log($p_\theta(\mathbf{x})$) does not depend

$$= \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z})}{p_\theta(\mathbf{x})} \right) \right] \quad (\text{Bayes Rule})$$

$$= \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z}) p_\theta(\mathbf{z})}{p_\theta(\mathbf{x})} q_\phi(\mathbf{z} | \mathbf{x}) \right) \right] \quad (\text{Multiply by 1})$$

$$= \mathbb{E}_{\mathbf{z}} \left[\log(p_\theta(\mathbf{x} | \mathbf{z})) - \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z})}{p_\theta(\mathbf{z})} \right) \right] \right] + \mathbb{E}_{\mathbf{z}} \left[\log \left(\frac{p_\theta(\mathbf{x} | \mathbf{z})}{p_\theta(\mathbf{z})} \right) \right]$$

$$= \mathbb{E}_{\mathbf{z}} \left[\log(p_\theta(\mathbf{x} | \mathbf{z})) \right] - KL(q_\phi(\mathbf{z} | \mathbf{x}), p_\theta(\mathbf{z})) + KL(q_\phi(\mathbf{z} | \mathbf{x}), \mathbf{z})$$

$$\stackrel{(1)}{=} \mathbb{E}_{\mathbf{z}} \left[\log(p_\theta(\mathbf{x} | \mathbf{z})) \right] + \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} [\log(1 - D_{\theta_d}(G_{\theta_d}(\mathbf{z})))]$$

$$\stackrel{(2)}{=} \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} [\log(1 - D_{\theta_d}(G_{\theta_d}(\mathbf{z})))]$$

$$+ \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} [\log(1 - D_{\theta_d}(G_{\theta_d}(\mathbf{z})))]$$

$$\stackrel{(3)}{=} \mathbb{E}_{\mathbf{z} \sim p_\theta(\mathbf{z} | \mathbf{x})} [\log(1 - D_{\theta_d}(G_{\theta_d}(\mathbf{z})))]$$

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:

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 on \mathbb{R}^d

Then compute both losses and perform gradient updates.

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