

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. $X \perp Y | Z \Rightarrow P(X, Y | Z) = P(X | Z)P(Y | Z)$
Cond. Ind. $X \perp Y | Z \Rightarrow P(X | Y, Z) = P(X | Z)$
 $E[X] = \int_X f_X(t) dt = \mu_X$
Cov(X, Y) = $E_{x,y}[X - E_x[X](Y - E_y[Y])]$
Cov(X) := Cov(X, X) = Var(X)
 X, Y independent \Rightarrow Cov(X, Y) = 0
 $XX^T \geq 0$ (symmetric positive semi-definite)

Var [X] = $E[X^2] - E[X]^2$
Var [$A\mathbf{x}$] = A Var [\mathbf{x}] A^T Var [$A\mathbf{x} + b$] = a^2 Var [\mathbf{x}]
Var [$\sum_{i=1}^n a_i X_i$] = $\sum_{i=1}^n a_i^2$ Var [X_i] + $2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
Var [$\sum_{i=1}^n a_i X_i$] = $\sum_{i=1}^n a_i^2$ Var [X_i] + $\sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
 $\frac{\partial}{\partial t} P(X \leq t) = P(X \leq t) - P(X | Z)$ (derivative of c.d.f. is p.d.f.)
 $f_{\alpha Y}(z) = \frac{1}{\alpha} f_Y(\frac{z}{\alpha})$

T. The moment generating function (MGF) $\psi_X(t) = E[e^{tX}]$ characterizes the distr. of a rv
Be(p): $p^t + (1-p)^{1-t}$ $\mathcal{N}(\mu, \sigma^2): \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$
Bin(n, p): $(pe + (1-p))^n$ $\text{Gam}(\alpha, \beta): (\frac{\alpha}{\alpha+\beta t})^\alpha$ for $t < 1/\beta$
Pois(λ): $e^{\lambda(e^{t-1})}$

T. If X_1, \dots, X_n are indep. rvs with MGFs $M_{X_i}(t) = E[e^{tX_i}]$, 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 : $f_Z(z) = \int_X f_X(x) f_Y(z-x) dt = \int_X f_X(x) f_Y(z-x) dx$

$N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$
 $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top$

T. $P(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}) = N\left(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \middle| \begin{pmatrix} \mathbf{H}_1 & \mathbf{H}_2 \\ \mathbf{S}_1 & \mathbf{S}_2 \end{pmatrix}\right)$
 $\mathbf{a}_1, \mathbf{u}_1 \in \mathbb{R}^e$, $\mathbf{S}_{12} \in \mathbb{R}^{e \times e}$ p.s.d., $\mathbf{S}_1 \in \mathbb{R}^{e \times f}$ p.s.d.
 $\mathbf{a}_2, \mathbf{u}_2 \in \mathbb{R}^f$, $\mathbf{S}_{22} \in \mathbb{R}^{f \times f}$ p.s.d., $\mathbf{S}_2 \in \mathbb{R}^{f \times e}$ p.s.d.

T. Chebyshev: Let X be a rv with $E[X] = \mu$ 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) $\|\mathbf{u}, \mathbf{v} \in V: \langle \mathbf{u}, \mathbf{v} \rangle \leq \|\mathbf{u}\| \|\mathbf{v}\|$.

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

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

- $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$
- $f'(x) \geq 0$
- 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

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, $\forall i: \lambda_i \geq 0, \sum_i \lambda_i = 1 \Rightarrow \sum_i \lambda_i f(x_i) \leq f(\sum_i \lambda_i x_i)$
Special cases $f(E[X]) \leq E[f(X)]$.

D. (Lagrangian Formulation) of $\arg \max_{x,y} f(x, y)$ s.t. $g(x, y) = c: L(x, y, \gamma) = f(x, y) - \gamma(g(x, y) - c)$
 $x^{t+1} = x^t - \text{Hess}(f)(x^t)^{-1} \nabla_{\mathbf{x}} f(x^t)$.

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 $\mathbf{x} \in \mathbb{R}^n: \mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$ 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. Otherwise: indefinite.

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

- $\text{Tr}(\mathbf{A}) = \sum_i \lambda_i$ (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 $\text{Tr}(cA) = c\text{Tr}(A)$.
- $\text{Tr}(\mathbf{A}^T) = \text{Tr}(\mathbf{A})$.

D. (Frobenius Norm) $\|A\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}$ The square root of the sum of the absolute squares of its elements.

Properties:

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

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

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

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

5 Regression

D. (Linear Regression): $f(x) = w^\top x + b$
Denominator $\frac{\partial}{\partial \mathbf{x}}[\mathbf{u}(\mathbf{x})] = \frac{\partial u}{\partial \mathbf{x}} + v(\mathbf{x}) \frac{\partial v}{\partial \mathbf{x}}$
Analytically solvable: $w^* = (X^\top X)^{-1} X^\top y$

D. (Ridge Regression): $f(x) = w^\top x + b$ with ℓ_2 regularization
 $\ell(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|w\|^2$
Analytically solvable: $w^* = (X^\top X + \lambda I)^{-1} X^\top y$

D. (Lasso Regression): $f(x) = w^\top x + b$ with ℓ_1 regularization
 $\ell(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|w\|_1$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{b}^\top \mathbf{A} \mathbf{x}] = \mathbf{A}^\top \mathbf{b}$

D. (Logistic Regression): $f(x) = \sigma(w^\top x + b)$
It minimizes the cross-entropy loss (negative log-likelihood), which acts as a surrogate loss for the 0/1 classification error.
 $\ell(w) = -\frac{1}{n} \sum_{i=1}^n \| \mathbf{x} - \mathbf{f}_\theta(\mathbf{x}) \|^2 = \frac{2}{\partial \mathbf{x}}[\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x})$

4.3 Vector-by-Vector

A. C, D, a, b not a function of \mathbf{x} , $\mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{x}), \mathbf{h}(\mathbf{x}), u = v(\mathbf{x}), v = w(\mathbf{x})$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{u}(\mathbf{x})] = u(\mathbf{x}) \frac{\partial \mathbf{x}}{\partial \mathbf{x}} + \mathbf{f}(\mathbf{x}) \frac{\partial w(\mathbf{x})}{\partial \mathbf{x}}$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a} \odot \mathbf{b}] = \mathbf{a} \odot \mathbf{b}$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{A}] = \frac{\partial \mathbf{a}}{\partial \mathbf{x}}^\top \mathbf{A}$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{A}^\top \mathbf{b}] = \mathbf{a}^\top \mathbf{A}^\top \mathbf{b}$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{A}^\top \mathbf{A} \mathbf{b}] = \mathbf{a}^\top \mathbf{A}^\top \mathbf{A} \mathbf{b}$

T. The moment generating function (MGF) $\psi_X(t) = E[e^{tX}]$ characterizes the distr. of a rv
Be(p): $p^t + (1-p)^{1-t}$ $\mathcal{N}(\mu, \sigma^2): \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$
Bin(n, p): $(pe + (1-p))^n$ $\text{Gam}(\alpha, \beta): (\frac{\alpha}{\alpha+\beta t})^\alpha$ for $t < 1/\beta$
Pois(λ): $e^{\lambda(e^{t-1})}$

T. If X_1, \dots, X_n are indep. rvs with MGFs $M_{X_i}(t) = E[e^{tX_i}]$, 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 : $f_Z(z) = \int_X f_X(x) f_Y(z-x) dt = \int_X f_X(x) f_Y(z-x) dx$

N($\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}$) = $\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$
 $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top$

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

6 Approximation Theory

6.1 Compositions of Maps

D. (Linear Function) A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear function if the following properties hold

- $\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n: f(\mathbf{x} + \mathbf{x}') = f(\mathbf{x}) + f(\mathbf{x}')$
- $\forall \lambda \in \mathbb{R}: f(\lambda \mathbf{x}) = \lambda f(\mathbf{x})$

T. $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear $\iff f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ for some $\mathbf{w} \in \mathbb{R}^n$

6.2 Hyperplane

$H := \{ \mathbf{x} | \langle \mathbf{w}, \mathbf{x} - \mathbf{p} \rangle = 0 \} = \{ \mathbf{x} | \langle \mathbf{w}, \mathbf{x} \rangle = b \}$
where $b = \langle \mathbf{w}, \mathbf{p} \rangle$. \mathbf{w} =normal vector, \mathbf{p} points onto a point on the plane.

6.3 Universal Approximation

T. (Universal Approximation Theorem) MLPs with one hidden layer and a non-polynomial activation function can approximate any continuous function on a compact subset of \mathbb{R}^d to arbitrary accuracy $\epsilon > 0$. Let ϕ be a non-constant, bounded, and continuous activation function (e.g., Sigmoid, ReLU). There exist weights v, w, b and an integer N (number of neurons) such that the network output $F(\mathbf{x})$ satisfies the condition: $|f(\mathbf{x}) - F(\mathbf{x})| < \epsilon, \quad \forall \mathbf{x} \in \mathbb{R}^d$

6.4 Scalar-by-Matrix

$\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X} \mathbf{b}] = \mathbf{ab}^\top$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{ba}^\top$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X}^\top \mathbf{X} \mathbf{b}] = \mathbf{X}(\mathbf{ab}^\top + \mathbf{b} \mathbf{a}^\top)$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X}^\top \mathbf{X}^\top \mathbf{b}] = \mathbf{X}^\top(\mathbf{ab}^\top + \mathbf{b} \mathbf{a}^\top)$
 $\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X}^\top \mathbf{X}^\top \mathbf{X} \mathbf{b}] = \mathbf{X}^\top \mathbf{X}(\mathbf{ab}^\top + \mathbf{b} \mathbf{a}^\top)$

T. (Eckhart-Young) For $m \leq \min(n, k)$ and the objective

$\arg \min_{\mathbf{X}} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 = \mathbf{U}_m \mathbf{U}_m^\top \mathbf{V}_m^\top$

Uniform Approximation: This specific type of convergence guarantees that the maximum error across the entire domain K is bounded by ϵ . It is stricter than pointwise approximation.

T. (Weierstrass Approximation) Let f be a continuous real-valued function defined on a closed interval $[a, b]$. For every $\epsilon > 0$, there exists a polynomial $P(x)$ such that for all $x \in [a, b]$: $|f(x) - P(x)| < \epsilon$

6.5 Vector-by-Matrix (Generalized Gradient)

$\frac{\partial}{\partial \mathbf{x}}[\mathbf{a}^\top \mathbf{X} \mathbf{b}] = \mathbf{ab}^\top$

5 Information Theory

D. (Entropy) Let X be a random variable distributed according to P . Then the entropy of X
 $H(X) = E[-\log(P(X))] = -\sum_{x \in \mathcal{X}} p(x) \log(p(x)) \geq 0$. It describes the expected information content $I(X)$ of X .

D. (Cross-Entropy) For distributions p and q over a given set is $H(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log(q(x)) = E_{x \sim p}[-\log(q(x))] \geq 0$. $H(p, q) = H(X) + K L(p, q) \geq 0$, where H uses p .

Com. The minimizer of the cross-entropy is $q := p$, due to the second formulation.

D. (Kullback-Leibler Divergence) For probability distributions p and q defined on the same probability space, the KL-divergence between p and q is defined as $KL(p, q) = -\sum_{x \in \mathcal{X}} p(x) \log\left(\frac{q(x)}{p(x)}\right) \geq 0$. Special cases $KL(p, p) = 0$, $KL(p, q) = H(p) - H(q)$, where H uses p .

T. (Montufar (Linear Regions)) The number of linear regions carved out by a Deep ReLU Network grows exponentially with depth L , but only polynomially with width n . For a network with L layers and width ($n \geq d$):
 $\# \text{Regions} = O\left(\left(\frac{n}{d}\right)^{(L-1)d} n^d\right)$

Com. Deep networks are exponentially more expressive (in terms of complex decision boundaries) than shallow networks of the same parameter count.

T. (Shekhutman (ReLU Basis)) Any continuous piecewise linear function $g(x)$ on $[0, 1]$ (a polygonal line) with breakpoints (knobs) $t_1 < \dots < t_k$ can be represented exactly as a weighted sum of ReLU units:

$$g(x) = a + bx + \sum_{i=1}^k c_i \text{ReLU}(x - t_i)$$

Com. A single hidden layer of ReLUs acts as a universal basis for 1D splines. This establishes the theoretical link between splines and ReLU networks.

T. (Dimension Lifting (Leshto's Theorem)) By the universal approximation theorem, we know that we can approximate $C(\mathbb{R})$ with a single hidden layer of ReLUs. The lifting theorem then allows us to lift the dimension of the function space to a higher dimension, i.e. span($\phi(x+b): b \in \mathbb{R}$) universally approximates $C(\mathbb{R})$.

6.6 Regularization Approaches

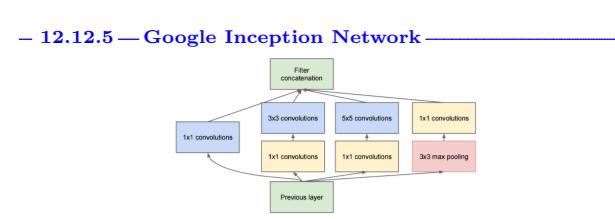
D. Weight Decay (L_2 Regularization): Adds a penalty to the loss function proportional to the square of the weights' magnitude to constrain complexity.

6.7 Universal Approximation

T. (Montufar (Linear Regions)) The number of linear regions carved out by a Deep ReLU Network grows exponentially with depth L , but only polynomially with width n . For a network with L layers and width ($n \geq d$):

$$\# \text{Regions} = O\left(\left(\frac{n}{d}\right)^{(L-1)d} n^d\right)$$

Com. Implications: Polynomials are dense



actually constant on the wall of the cliff. Let's have a look at this through some equations.

Now, let's evaluate what the risk is at some point, plus some gradient step. If we do the 2nd-order Taylor expansion of that, then we get

$$\mathcal{R}(\theta - \eta \nabla_\theta \mathcal{R}(\theta))$$

$$\stackrel{\text{Taylor}}{\approx} \mathcal{R}(\theta) - \eta \|\nabla_\theta \mathcal{R}(\theta)\|_2^2 + \frac{\eta^2}{2} \frac{\nabla_\theta \mathcal{R}(\theta)^T \mathbf{H} \nabla_\theta \mathcal{R}(\theta)}{\|\nabla_\theta \mathcal{R}(\theta)\|_2^2}$$

where

$$\mathbf{H} := \nabla^2 \mathcal{R}(\theta) \quad (\text{Hessian matrix})$$

Now, what we want is that the rest of the sum is negative. If that is the case, because then we're improving our cost function. If not, we're basically diverging.

So,

- the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_2^2$ will obviously be negative, as it's a negative factor of a norm
- the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may be already small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

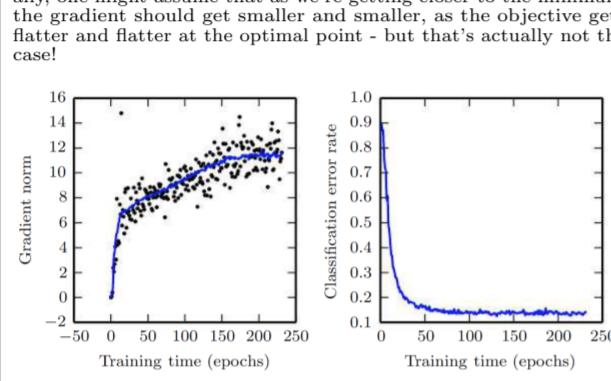
$$\frac{\eta}{2} \|\nabla_\theta \mathcal{R}(\theta)\|_2^2 \gtrsim \|\nabla_\theta \mathcal{R}(\theta)\|^2$$

So, the hessian term becomes much larger than the gradient. So we're not improving our cost function.

and the remaining terms, will be negative (as defined by the Taylor sum)

So a typical remedy for first-order methods is to take very small step sizes η .

However, things become even stranger because of the curvature. As we can see, the gradient norm gets larger and larger the more we train (can be checked empirically). And the gradient norm also tends to have larger fluctuations. And as we can see, starting at some point, the error just fluctuates around at a certain level. Actually, one might assume that as we're getting closer to the minimum, the gradient should get smaller and smaller, as the objective gets flatter and flatter at the optimal point - but that's actually not the case!



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

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

Then the convolved image has dimensions (assuming stride 1)

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

• 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

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

13.2 Gradient Descent

$$\theta(t+1) = \theta(t) - \eta_t \nabla_\theta \mathcal{R}$$

13.3 Gradient Descent: Classic Analysis

In classical machine learning we have a **convex** objective \mathcal{R} . And we denote

- \mathcal{R}^* as the minimum of \mathcal{R}
- θ^* as the optimal set of parameters (the minimizer of \mathcal{R})

So we have $\nabla_\theta \mathcal{R} \neq 0$: $\mathcal{R}^* = \mathcal{R}(\theta^*) \leq \mathcal{R}(\theta)$.

D. (Strictly Convex Objective) \rightarrow objective has only one (a unique) minimum.

$$\forall \theta \neq \theta^* : \mathcal{R}^* := \mathcal{R}(\theta^*) < \mathcal{R}(\theta)$$

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 $L \in \mathbb{R}_+$ (**Lipschitz constant**), such that

$$\forall x_1, x_2 \in \mathbb{R}^n : \|f(x_1) - f(x_2)\| \leq L \|x_1 - x_2\|.$$

For our risk function \mathcal{R} , we say that the gradient of it

$$\nabla_\theta \mathcal{R}: \Omega \rightarrow \Omega \quad \text{where } \Omega = \mathbb{R}^n$$

L -Lipschitz continuous if it holds that

$$\forall \theta_1, \theta_2 \in \Omega : \|\nabla_\theta \mathcal{R}(\theta_1) - \nabla_\theta \mathcal{R}(\theta_2)\| \leq L \|\theta_1 - \theta_2\|$$

Com. So, the L tells us how big the gradient could be.

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

T: An everywhere differentiable function $f: \mathbb{R} \rightarrow \mathbb{R}$ is Lipschitz continuous (with $L = \sup |f'(x)|$) iff it has **bounded first derivatives**.

T: In particular any continuously differentiable function is **locally Lipschitz continuous**. As continuous functions are bounded on an interval, so its gradient is locally bounded as well.

T: If \mathcal{R} is convex with L -Lipschitz-continuous gradients then we've that

$$\mathcal{R}(\theta(t)) - \mathcal{R}^* \leq \frac{2L}{t+1} \|\theta(0) - \theta^*\|^2 \in \mathcal{O}(t^{-1})$$

Com. So we have a polynomial (linear) convergence rate of θ towards the optimal parameter θ^* (note: just in the convex setting!). As we can see, the convergence time is bounded by a t that depends on our initial guess, and the Lipschitz constant L .

Com. Usually one value for η that people use in this setting is $\eta = \frac{1}{L}$ or $\eta := \frac{2}{L}$.

D. (Strongly Convex Function) A differentiable function f is **strongly convex** if the following inequality holds for all points x, y in its domain:

$$\nabla_\mathbf{x} \mathbf{y} : (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{y} - \mathbf{x}) \geq \mu \|\mathbf{y} - \mathbf{x}\|^2$$

where $\|\cdot\|$ is any norm. An equivalent condition is the following:

$$\forall \mathbf{x}, \mathbf{y} : f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Com. The concept of strong convexity extends and parametrizes the notion of strict convexity. A strongly convex function is also strictly convex, but not vice versa. Notice how the definition of strong convexity approaches the definition for strict convexity as $\mu \rightarrow 0$, and is identical to the definition of a convex function when $\mu = 0$. Despite this, functions exist that are strictly convex, but are not strongly convex for any $\mu > 0$.

T: Now, when \mathcal{R} is μ -strongly convex in θ and its gradient is L -Lipschitz continuous \Rightarrow

$$\mathcal{R}(\theta(t)) - \mathcal{R}^* \leq \left(1 - \frac{\mu}{L}\right)^t (\mathcal{R}(\theta(0)) - \mathcal{R}^*) \in \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^t\right)$$

So we have

• an exponential convergence ("linear rate")

• and the rate depends adversely on the condition number $\frac{L}{\mu}$. So we want the maximum gradient to be small, and we want the curvature to be large (which are somewhat contrary desires, but ideally the condition number is very close to 1).

T: If we use Nesterov acceleration (in the general case), then we get a polynomial convergence of $\mathcal{O}(t^{-2})$.

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

13.4 Optimization Challenges in NNs: Curvatures

A challenge in NNs arises when we have sharp non-linearities, then there are two approaches to solve this

• one is to be very conservative and only do small update steps by choosing a very small learning rate

• or we are courageous and due huge jumps as depicted in the image.

Typical approaches are to clip the gradient when it gets too large, or use a decreasing learning rate (in terms of time).

Now, the problem is not that the cliff is very steep. The problem is the curvature. Because when we take the gradient, the gradient is

actually constant on the wall of the cliff. Let's have a look at this through some equations.

Now, let's evaluate what the risk is at some point, plus some gradient step. If we do the 2nd-order Taylor expansion of that, then we get

$$\mathcal{R}(\theta - \eta \nabla_\theta \mathcal{R}(\theta)) \stackrel{\text{Taylor}}{\approx} \mathcal{R}(\theta) - \eta \|\nabla_\theta \mathcal{R}(\theta)\|_2^2 + \frac{\eta^2}{2} \frac{\nabla_\theta \mathcal{R}(\theta)^T \mathbf{H} \nabla_\theta \mathcal{R}(\theta)}{\|\nabla_\theta \mathcal{R}(\theta)\|_2^2}$$

Note that when computing the derivative we've used the following trace differentiation rules (cf. wikipedia, The Matrix Cookbook)

$$\nabla_\mathbf{A} \text{Tr}(\mathbf{AA}^\top) = 2\mathbf{A} \quad \nabla_\mathbf{A} \text{Tr}(\mathbf{AB}) = \mathbf{B}^\top$$

Note that one could have solved the solution also through the following way, by recognizing that $\mathbf{A}(t)$ follows the following differential equation

$$\dot{\mathbf{A}}(t) = -\eta \nabla_\mathbf{A} \mathcal{R}(\mathbf{A}(t)) = 2\eta(\mathbf{I} - \mathbf{A}(t)) = 2\eta\mathbf{I} - 2\eta\mathbf{A}(t).$$

So rearranging the equation for $\mathbf{A}(t)$ we get

$$\mathbf{A}(t) = -\frac{1}{2} \eta \dot{\mathbf{A}}(t) + \mathbf{I}$$

And now, since we're using gradient descent, we'll converge, and the gradients will go to zero, hence

$$\lim_{t \rightarrow \infty} \mathbf{A}(t) = -\frac{1}{2} \eta \left(\lim_{t \rightarrow \infty} \dot{\mathbf{A}}(t) \right) + \mathbf{I} = \mathbf{I}$$

So $\mathbf{A}(t)$ will converge to \mathbf{I} .

13.7 Least Squares: Two-Layer Lin. Netw.

Now, the question is what happens, when we build a two-layer linear network (again with the squared error), with no nonlinearity. So we'll have a linear mapping (that is composed of two linear mappings)

$$\mathbf{F}(\mathbf{x}) = \mathbf{Ax} = \mathbf{QWx}$$

Now, from what we've seen before, we can express the risk (due to trace identities, trace linearity, etc.) just by replacing $\mathbf{A} = \mathbf{QW}$,

$$\mathcal{R}(\mathbf{Q}, \mathbf{W}) = \text{const.} + \text{Tr}((\mathbf{QW})(\mathbf{QW})^\top) - 2\text{Tr}(\mathbf{QW}^\top)$$

Now, taking the derivatives w.r.t. the parameters, we get (using the chain rule)

$$\nabla_\mathbf{Q} \mathcal{R}(\mathbf{Q}, \mathbf{W}) = \frac{\partial \mathcal{R}(\mathbf{A})}{\partial \mathbf{Q}} \frac{\partial \mathbf{A}}{\partial \mathbf{Q}}$$

$$\nabla_\mathbf{W} \mathcal{R}(\mathbf{Q}, \mathbf{W}) = \frac{\partial \mathcal{R}(\mathbf{A})}{\partial \mathbf{W}} \frac{\partial \mathbf{A}}{\partial \mathbf{W}}$$

Which in the end gives us

$$\nabla_\mathbf{Q} \mathcal{R}(\mathbf{Q}, \mathbf{W}) = 2\mathbf{QW}^\top - 2\mathbf{FW}^\top = 2(\mathbf{A} - \mathbf{I})\mathbf{W}^\top$$

$$\nabla_\mathbf{W} \mathcal{R}(\mathbf{Q}, \mathbf{W}) = 2\mathbf{Q}^\top \mathbf{W} - 2\mathbf{Q}^\top \mathbf{F} = 2\mathbf{Q}^\top (\mathbf{A} - \mathbf{I})$$

Now we compute the (partial) row sums of \mathbf{G} (note: not the gradient norms! \rightarrow rows!)

$$\gamma_i^2(t) := \sum_{s=1}^t g_{is}^2$$

And then we adapt the gradient stepsize for each dimension as follows:

$$\theta_i(t+1) = \theta_i(t) - \frac{\eta}{\delta + \gamma_i^2(t)} \nabla_\theta \mathcal{R}(\theta_i(t), \mathbf{W}). \quad \delta > 0 \text{ (small)}$$

This will transform the gradient such as if the loss landscape would be in a more isometric shape. It will scale the gradient appropriately to each dimension. So instead of having a valley, we'll have a nice round hole again. This avoids this typical situation where the gradient descent bounces left and right in the valley, instead of walking down the valley.

What is not very clear is why batch-normalization works.

The original paper about batch-normalization (BN) said that BN reduces the internal covariance shift of the data. What they meant by this is that: let's say that we have a very simple classifier, that will basically classify everything that is negative to one class, and everything that is positive to another class. Then, when we just shift the data by a constant vector, then, without batch-normalization we'd shift all the datapoints into one class. However, with BN since the mean is removed we'll remove that constant shift the BN layer and it will work out. So BN reduces the covariance shift. That was the main idea.

However, it turns out that some other people came later on and said the following: They didn't negate the effect of the covariance-shift reduction, but the reason they said that BN works is that it makes the landscape of the loss more smooth. Hence, the optimization works better and gives better results.

No one really knows why BN works so well.

13.8.8 Other Heuristics

• Curriculum learning and non-uniform sampling of data points \rightarrow focus on the most relevant examples (Bengio, Louradour, Collobert, Weston, 2009) (DL-Book: 8.7.6) Or increase hardness of tasks (corner-cases) as NN improves

• Continuation methods: define a family of (simpler) objective functions and track solutions, gradually change hardness of loss (DL-Book: 8.7.6)

• Heuristics for initialization (DL-Book: 8.4) scale the weights of

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.

D. (Pointwise Mutual Information)

$$\text{pmi}(v, w) = \log \left(\frac{P(v, w)}{P(v)P(w)} \right) = \log \left(\frac{P(v|w)}{P(v)} \right) \approx \mathbf{x}_v^T \mathbf{x}_w + \text{const}$$

Com. As you can see this metric is bi-linear.

So we interpret the vectors as *latent variables* and link them to the observable probabilistic model. So the pointwise mutual information is related to the inner product between the latent vectors (the more related, the more co-linear the latent representations have to be).

Now, how do we compute the pointwise mutual information? One thing that we could do is to just look for words that are nearby and compute these probabilities empirically. This leads us to the idea of skip-grams.

D. (Skip Grams) The skip-gram approach is an approach to look at co-occurrences of words within a window size R (instead of looking at subsequences of some length n as with n grams). So we're only interested in the co-occurrence within some window size of words R , rather than a precise sequence.

D. (Co-Occurrence Set) Here we look at the *pairwise occurrences* of words in a *context window* of size R . So, if we have a long sequence of words $\mathbf{w} = (w_1, \dots, w_T)$, then the co-occurrence index set is defined as

$$C_R := \{(i, j) \mid 1 \leq |i - j| \leq R\}.$$

D. (Co-Occurrence Matrix) Note that in order to get an (empirical) idea of the co-occurrence frequencies one could compute the co-occurrence matrix

$$\mathbf{C} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}, \quad \text{where } C_{ij} = \#\text{of co-occurrences of } w_i \text{ and } w_j \text{ within window size } R.$$

Properties: $\mathbf{C} = \mathbf{C}^T$ (symmetric), peaky, sparse.

One approach for embeddings: do PCA of \mathbf{C} and use k eigenvectors corresponding to largest eigenvalues of \mathbf{C} . Note that we have

$$C_{ij} = \underbrace{\text{one-hot}(w_i)}_{\mathbf{o}_i} \underbrace{\text{one-hot}(w_j)}_{\mathbf{o}_j} = \mathbf{o}_i \mathbf{V} \mathbf{A} \mathbf{V}^T \mathbf{o}_j \\ \approx \mathbf{o}_i \underbrace{\mathbf{V}_k \mathbf{A}_k \mathbf{V}_k^T}_{k \text{ PCs}} \mathbf{o}_j = \mathbf{o}_i \mathbf{V}_k \frac{1}{k} \mathbf{A}_k^{\frac{1}{k}} \mathbf{V}_k^T \mathbf{o}_j$$

de-embedding: $\mathbf{V} \mathbf{A}_k^{-\frac{1}{k}}$ (then find nearest neighbour)

Problem: \mathbf{C} is huge ($|\mathcal{V}|^2$), hence matrix-factorization becomes prohibitively expensive!

Solution: Use skip-gram approach to avoid computing \mathbf{C} at all!

The solution to this is pretty simple: we train a model that tries to predict for one word w_t the preceding and following words:

$$w_{t-c} \rightarrow \dots \rightarrow w_{t-1}, w_t, w_{t+1}, \dots, w_t + c - 1, w_{t+c}$$

Here's an illustration of the model for $t = 3$: input $w_t \rightarrow$ projection $\rightarrow w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}$ output

Note that the assumption (or simplification) of this model is that it assumes that the words W_i, W_j within the window $+c, -c$ of W_t are conditionally independent of each other given W_t .

$$W_t \perp W_j \mid W_t \quad (i \neq j, i \neq t, j \neq t)$$

That might be too much of an assumption but you can see that sometimes when we're talking about something we may change the order of the words and still mean the same thing (e.g., "I was born in 1973." "1973 is the year I was born"). So in a way we're just trying to capture the meaning of W_t with this. So this gives us an idea of the context of W_t and might relieve the structure we're looking for. So, it's not as optimal as computing \mathbf{C} , but it's a way to start.

So actually, what we want to do is want to maximize the likelihood of the co-occurrences in our dataset:

$$\theta^* = \arg \max_{\theta} \prod_{(i,j) \in C_R} P_{\theta}(w_i \mid w_j)$$

Now our approach to approximate the probability $P_{\theta}(w_i \mid w_j)$ as follows: it should be something that is related to the dot product of the embeddings, so $\mathbf{x}_w^T \mathbf{z}_{w_j}$ (note how we use two different embeddings as the conditional probability is asymmetric), but in order to make the probability positive we'll take the exponential of it and normalize. Further, for SGD it's always better to optimize a sum: so we'll optimize the log-likelihood of co-occurring words in our dataset $\mathbf{w} = (w_1, \dots, w_T)$:

$$\begin{aligned} &= \arg \max_{\theta} \sum_{(i,j) \in C_R} \log \left(\frac{\exp(\mathbf{x}_{w_i}^T \mathbf{z}_{w_j})}{\sum_{u \in \mathcal{V}} \exp(\mathbf{x}_{w_i}^T \mathbf{z}_{w_u})} \right) \\ &= \arg \max_{\theta} \sum_{(i,j) \in C_R} \mathbf{x}_{w_i}^T \mathbf{z}_{w_j} - \log \left(\sum_{u \in \mathcal{V}} \exp(\mathbf{x}_{w_i}^T \mathbf{z}_u) \right) \end{aligned}$$

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 prob.)

$$\theta = (\mathbf{x}_w, \mathbf{z}_w) \in \mathbb{R}^n$$

$\cdot \mathbf{x}_w$ is used to predict w 's conditional probability, and $\cdot \mathbf{z}_w$ is used to use w as an evidence in the cond. prob.

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

Problem: Note however that it's too expensive to compute the *partition function* as we have to do a full sum over \mathcal{V} (can be ~ 10^6 up to 10^7). And we'd have to do that every time we pass a new batch-sample (w_{t+1}, w_t) through the network.

Brilliant idea of skip-grams: instead of computing the partition function, turn the problem of determining $P_{\theta}(w_i \mid w_j)$ 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

$$P_{\theta}(w_i \mid w_j) \approx P_{\theta}(D_{w_i, w_j} = 1 \mid \mathbf{x}_{w_i}, \mathbf{z}_{w_j}) = \sigma(\mathbf{x}_{w_i}^T \mathbf{z}_{w_j})$$

and the opposite event is given by:

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

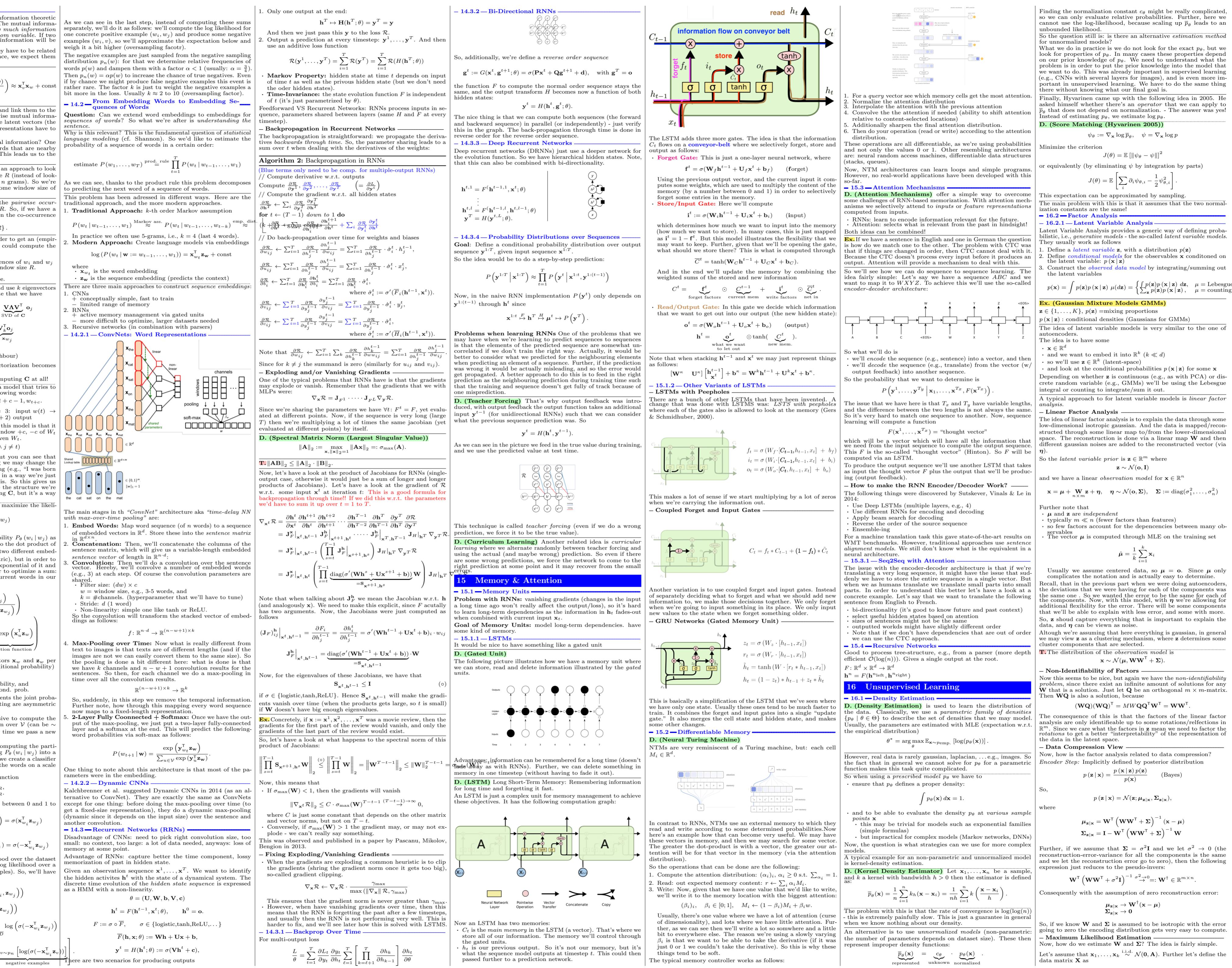
Further, instead of just maximizing the likelihood over the dataset of co-occurrences D , we'll also maximize the log likelihood over a dataset of non-co-occurrences \bar{D} (negative samples). So, we'll have the following maximization problem

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \sum_{(i,j) \in D} \log(P_{\theta}(D_{w_i, w_j} = 1 \mid \mathbf{x}_{w_i}, \mathbf{z}_{w_j})) \\ &\quad + \sum_{(i,j) \in \bar{D}} \log(1 - P_{\theta}(D_{w_i, w_j} = 0 \mid \mathbf{x}_{w_i}, \mathbf{z}_{w_j})) \\ &= \arg \max_{\theta} \sum_{(i,j) \in D} \log(\sigma(\mathbf{x}_{w_i}^T \mathbf{z}_{w_j})) + \sum_{(i,j) \in \bar{D}} \log(\sigma(-\mathbf{x}_{w_i}^T \mathbf{z}_{w_j})) \\ &= \arg \max_{\theta} \sum_{(i,j) \in C_R} \underbrace{\log(\sigma(\mathbf{x}_{w_i}^T \mathbf{z}_{w_j}))}_{\text{positive examples}} + \underbrace{\log(\sigma(-\mathbf{x}_{w_i}^T \mathbf{z}_{w_j}))}_{\text{negative examples}} \end{aligned}$$

There are two scenarios for producing outputs

1. For multi-output loss

$$\frac{L}{\theta} = \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \sum_{t=1}^T \left[\prod_{k=t+1}^T \frac{\partial h_k}{\partial h_{k-1}} \right] \frac{\partial h_t}{\partial \theta}$$



$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_k \end{bmatrix}$$

and the empirical co-variance matrix as

$$\mathbf{S} := \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{k} \mathbf{X} \mathbf{X}^\top.$$

Then, the log-likelihood of the data \mathbf{X} , given \mathbf{A} can be written as:

$$\log(P(\mathbf{X}; \mathbf{A})) = -\frac{k}{2} \left(\text{Tr}(\mathbf{S} \mathbf{A}^{-1}) - \log(\det(\mathbf{A})) \right) + \underset{\text{i.T. of } \mathbf{A}}{\text{const}}$$

Note: this can be verified by using the definition of \mathbf{S} , the cyclic property of the trace, and then just write down the matrix-product as block-matrices and see what is the diagonal of the resulting matrix.

Now, let's compute the matrix gradients w.r.t. \mathbf{A} to know the equations that we need to compute the maximum likelihood:

$$\nabla_{\mathbf{A}} \text{Tr}(\mathbf{S} \mathbf{A}^{-1}) = -\mathbf{A}^{-1} \mathbf{S}$$

$$\nabla_{\mathbf{A}} \log(\det(\mathbf{A})) = \mathbf{A}^{-1}$$

Now, setting the gradient of the log-likelihood to zero gives us the following condition:

$$\nabla_{\mathbf{A}} \log(P(\mathbf{X}; \mathbf{A})) = 0 \iff \mathbf{S} \mathbf{A}^{-1} = \mathbf{I}.$$

So, the MLE for \mathbf{A} is just \mathbf{S} .

But recall, that what we want is not \mathbf{A} , but we want \mathbf{W} and Σ . However, we know that \mathbf{A} is just the empirical covariance matrix, and \mathbf{W} will be the mapping to the low-dimensional space and Σ is the reconstruction error.

$$\mathbf{A} = \mathbf{W} \mathbf{W}^\top + \Sigma$$

Now, using the chain rule we get:

$$\nabla_{\mathbf{W}} \mathbf{A} = 2\mathbf{W}$$

$$\nabla_{\Sigma} \mathbf{A} = \mathbf{I}$$

This gives us the following stationary condition for \mathbf{W} given Σ :

$$\mathbf{S}(\mathbf{W} + \mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} = \mathbf{W} \text{diag}\left(\frac{1}{\sigma^2 + \rho_i^2}\right).$$

Putting this back into the stationary condition, for each column \mathbf{w}_i of \mathbf{W} we get an eigenvector equation:

$$\mathbf{Sw}_i = (\sigma^2 + \rho_i^2)\mathbf{w}_i \quad \mathbf{Sw} = \text{diag}(\lambda)\mathbf{W}.$$

Then, if \mathbf{u}_i is the i -th eigenvector of \mathbf{S} , then

$$\mathbf{w}_i = \rho_i \mathbf{u}_i, \quad \rho_i^2 = \max\{0, \lambda_i - \sigma^2\}.$$

This gives us the probabilistic interpretation PCA and showed us how we can derive the PCA as a special case for $\sigma^2 \rightarrow 0$ (Tipping & Bishop, 1999).

– 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

$$M_X : \mathbb{R}^n \rightarrow \mathbb{R}, \quad t \mapsto \mathbb{E}_{\mathbf{X}} [e^{t^\top \mathbf{X}}].$$

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

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

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, MGPs 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 $M_{\mathbf{X}+\mathbf{Y}} = M_{\mathbf{X}} M_{\mathbf{Y}}$.

– 16.3 Latent Variable Models

There's another way of looking at latent variable models which is by the DeFinetti exchangeable 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

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

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

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

– 16.3.3 Dimensionality Reduction

One of the recurring things that we see in all of these models is dimensionality reduction. So we have that

$$\mathbf{X} = f(\mathbf{ZB})$$

where

- \mathbf{X} is $N \times D$,
- \mathbf{Z} is $N \times K$,
- \mathbf{B} is $K \times D$, and
- $K \ll D$.

So we have the data \mathbf{X} that we're trying to understand. We'll try to understand this data by a tall matrix \mathbf{Z} and a fat matrix \mathbf{B} . The tall matrix are the latent factors that we've talking about (how do we summarize the information of each sample). And the matrix \mathbf{B} is telling us how we can recover the original data from the summary. Most of the unsupervised algorithms can be captured in this general framework.

Depending on $f(\cdot)$, \mathbf{Z} and \mathbf{B} , we arrive at different models:

- Principal Component Analysis / Factor Analysis (f linear)
- Nonnegative Matrix Factorization (f "psomol" or Bernoulli model, and both \mathbf{Z} and \mathbf{B} have to be a nonnegative matrix).
- LLE/Isomap/GPLVM (here we also try to do PCA or Factor analysis with nonlinear components (with p.w. linear components))
- Restricted Boltzmann Machine (the idea is that \mathbf{Z} is discrete)
- Dirichlet Process (aka Chinese Restaurant Process)
- Beta Process (aka Indian Buffet Process)
- Implicit Models (e.g., Generative Adversarial Networks) (here all the information is moved to the function f instead of computing the matrices \mathbf{B} and \mathbf{C})

– 16.3.4 Implicit Models

Here we develop statistical models via: generating stochastic mechanism or simulation process.

Deep implicit models

- latent code $\mathbf{z} \in \mathbb{R}^d$, $\mathbf{z} \sim \pi(\mathbf{z})$, e.g. $\pi(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- parametric mechanism: $F_{\theta}: \mathbb{R}^d \rightarrow \mathbb{R}^m$
- induced distribution $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- sampling is easy: random vector + forward propagation.

IV. Chain Rule and Jacobians for Tensors

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

D. (Tensor Multiplication)

$$\mathbf{T} = \underset{\in \mathbb{R}^{(a+b)}}{\mathbf{P}} \times_b \underset{\in \mathbb{R}^{(b+c)}}{\mathbf{Q}} = \underset{\in \mathbb{R}^{(a \times c)}}{\mathbf{P}} \times_{\mathbf{s}_1, \dots, \mathbf{s}_b} \underset{\in \mathbb{R}^{(c \times b)}}{\mathbf{Q}} \times_{\mathbf{s}_1, \dots, \mathbf{s}_b} \cdots \times_{\mathbf{s}_1, \dots, \mathbf{s}_b} \mathbf{Q}$$

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

Note that this is just the sum of the multiplications of two numbers which are corresponding locations in \mathbf{P} and \mathbf{Q} . Essentially, it's the dot product across the dimensions s_1, \dots, s_b .

Note how this tensor-tensor-multiplication isomorphic to some matrix-matrix product:

$$T_{i_1, \dots, i_a, k_1, \dots, k_c} := \sum_{j_1, \dots, j_b} P_{i_1, \dots, i_a, j_b} Q_{j_b, \dots, j_b, k_1, \dots, k_c}$$

T. (Tensor Chain Rule)

$$y(W) : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^{d_3 \times d_4}, L(y) : \mathbb{R}^{d_3 \times d_4} \rightarrow \mathbb{R}$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial y} \times_{d_3, d_4} \frac{\partial y}{\partial W}$$

then we have: $T_{i,j,k,l} = \frac{\partial y_{i,j}}{\partial W_{k,l}}$

18 Generative Models

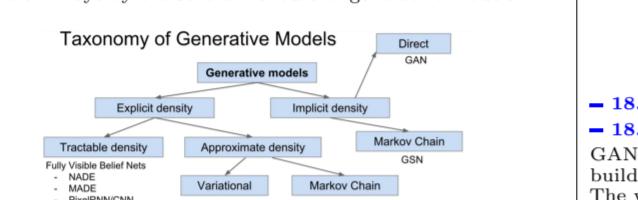
In unsupervised learning our goal is to learn some underlying hidden structure of the data (clustering, dimensionality reduction, feature learning, density estimation). Now, generative modeling has the following goal:

Goal: given data D , generate new samples from the same distribution p_{data} . We want to learn p_{model} similar to p_{data} .

The nice thing is that the training data is cheap, as we need no labels. However, it's a hard task.

In some way or another, any generative model has to cope with density estimation (which is the hard task). This problem is tackled in different ways by the several flavours of generative models:

Taxonomy of Generative Models



18.1 Variational Autoencoders

Recall, that with autoencoders, we had defined a concatenation of two differentiable (non-linear) mappings $x \xrightarrow{E} z \xrightarrow{D} \hat{x}$ (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 use 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 a 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

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

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).

Now the question is can we use a similar kind of setup to use new images? 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 decoder network modeling $p_{\theta}(\mathbf{x} | \mathbf{z})$, we define an additional encoder network $q_{\theta}(\mathbf{z} | \mathbf{x})$ that approximates $p_{\theta}(\mathbf{z} | \mathbf{x})$.

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

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

Enc. Netw. E: models $q_{\theta}(\mathbf{z} | \mathbf{x})$ 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)

$$\log(p$$