

1 Probability

Sum Rule $P(X = x_i) = \sum_j p(x_i, Y = y_i)$
Product rule $P(X, Y) = P(Y|X)P(X)$
Independence $P(X, Y) = P(X)P(Y)$
Bayes' Rule $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_{i=1}^k P(X|Y_i)P(Y_i)}$
Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X|Y|Z) = P(X|Z)P(Y|Z)$
Cond. Ind. $X \perp\!\!\!\perp Y|Z \Rightarrow P(X|Y, Z) = P(X|Z)$
 $E[X] = \int_X t \cdot f_X(t) dt =: \mu_X$
 $\text{Cov}(X, Y) = E_{x,y}[(X - E_x[X])(Y - E_y[Y])]$
 $\text{Cov}(X, X) := \text{Var}(X) = \text{Var}[X]$
 $X, Y \text{ independent} \implies \text{Cov}(X, Y) = 0$
 $XX^T \geq 0$ (symmetric positive semidefinite)
 $\text{Var}[X] = E[X^2] - E[X]^2$
 $\text{Var}[\mathbf{A}\mathbf{X}] = \mathbf{A}\text{Var}[\mathbf{X}]\mathbf{A}^T \quad \text{Var}[a\mathbf{X} + b] = a^2 \text{Var}[\mathbf{X}]$
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
 $\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$
 $\frac{\partial}{\partial t} P(X \leq t) = \frac{\partial}{\partial t} F_X(t) = f_X(t)$ (derivative of c.d.f. is p.d.f.)
 $f_{a,Y}(y) = \frac{1}{\alpha} \delta(y - \frac{x}{\alpha})$
T. The moment generating function (MGF) $\psi_X(t) = E[e^{tX}]$ characterizes the distr. of a rv
 $B(p) := p e^{-t} + (1-p) \quad N(\mu, \sigma) := \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$
 $\text{Bin}(n, p) := (pe + (1-p))^n \quad \text{Gam}(\alpha, \beta) := (\frac{1}{\alpha} \beta t)^{\alpha}$ for $t < 1/\beta$
 $\text{Pois}(\lambda) := e^{\lambda(e^{-t}-1)}$
T. If X_1, \dots, X_n are indep. rvs with MGFs $M_{X_i}(t) = E[e^{tX_i}]$, 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 : $f_Z(z) = \int_R f_X(x) f_Y(z-x) dt = \int_R f_X(x-t) f_Y(t) dt$
 $N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(\boldsymbol{\Sigma})} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$
 $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top$
T. $N(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}) = N(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} | \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix})$
 $\mathbf{a}_1, \mathbf{u}_1 \in \mathbb{R}^e, \Sigma_{11} \in \mathbb{R}^{e \times e}$ p.s.d., $\Sigma_{12} \in \mathbb{R}^{e \times f}$ p.s.d.
 $\mathbf{a}_2, \mathbf{u}_2 \in \mathbb{R}^f, \Sigma_{22} \in \mathbb{R}^{f \times f}$ p.s.d., $\Sigma_{21} \in \mathbb{R}^{f \times e}$ p.s.d.
 $P(\mathbf{a}_2 | \mathbf{a}_1 = \mathbf{z}) = N(\mathbf{a}_2 | \mathbf{u}_2 + \Sigma_{21} \Sigma_{11}^{-1}(\mathbf{z} - \mathbf{u}_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{21})$
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) $\forall u, v \in V: |\langle u, v \rangle| \leq \|u\| \|v\|$
 $\forall u, v \in V: 0 \leq |\langle u, v \rangle| \leq \|u\| \|v\|$
Special cases: $(\sum_i x_i y_i)^2 \leq (\sum_i x_i^2)(\sum_i y_i^2)$. Special case: $E[XY]^2 \leq E[X^2]E[Y^2]$
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
 $\forall x, x' \in S, \lambda \in [0, 1]: f(\lambda x + (1-\lambda)x') \leq \lambda f(x) + (1-\lambda)f(x')$
Com. A function is strictly convex if the line segment between any two points on the graph of the function lies strictly above the graph. This guarantees that there is a unique global minimum.
T. (Properties of Convex Functions)

- $f(y) \geq f(x) + \nabla f(x)^T(y - x) - \frac{1}{2}\|y - x\|^2$
- $f'(x) \geq 0$
- 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: $E[X] \leq E[f(X)]$
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 $L \in \mathbb{R}_0^+$ (*L*spschitz constant), such that
 $\forall x_1, x_2 \in X: \|f(x_1) - f(x_2)\| \leq L \cdot \|x_1 - x_2\|$.
Com. If the objective function is L -smooth a step size of $\eta = 1/L$ guarantees convergence.
D. (Lagrangian Formulation) $\text{arg max}_{x,y} f(x, y) \quad \text{s.t. } g(x, y) = C(x, y, \gamma) = f(x, y) - \gamma(g(x, y) - \gamma^*)$
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:

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 $x \in \mathbb{R}^n: x^T A x \geq 0$

- 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.

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

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

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

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

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

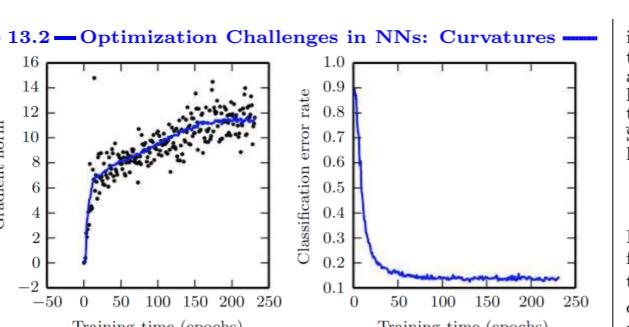
Remark. Sometimes the numerator-layout is used, where the Jacobian is defined as $(J_f)_{ij} = \frac{\partial f_i}{\partial x_j} / \frac{\partial x_j}{\partial x_j} = \frac{\partial f_i}{\partial x_j}$. The two conventions are related by transposition.

4.2 Scalar-by-Vector

Denominator Convection $\frac{\partial}{\partial x} [u(\mathbf{x})] = u(x) \frac{\partial x}{\partial x} + v(x) \frac{\partial v}{\partial x}$
 $\frac{\partial}{\partial x} [u(v(\mathbf{x}))] = u(\mathbf{x}) \frac{\partial v}{\partial x} + v(\mathbf{x}) \frac{\partial u}{\partial x}$
 $\frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{g}(\mathbf{x})] = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}) + \frac{\partial \mathbf{g}}{\partial \mathbf{x}}(\mathbf{x}) f(\mathbf{x}) = J_f(\mathbf{x}) + J_g(\mathbf{x})$
 $\frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{A} \mathbf{g}(\mathbf{x})] = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}) \mathbf{A} + \frac{\partial \mathbf{g}}{\partial \mathbf{x}}(\mathbf{x})^T f(\mathbf{x})$
 $\frac{\partial}{\partial x} [a^T \mathbf{x}] = \frac{\partial a}{\partial x} \mathbf{x}^T = a \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [x^T \mathbf{A}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \\ \frac{\partial}{\partial x} [x^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [b^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [a^T \mathbf{x}] = \frac{\partial a}{\partial x} \mathbf{x}^T = a \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [b^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{A}(\mathbf{x})^T \mathbf{C}(\mathbf{D}(\mathbf{x}) + \mathbf{e})] = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} + \mathbf{b}) + \mathbf{A}^T \mathbf{C}(\mathbf{D} + \mathbf{e}) \\ \frac{\partial}{\partial x} [\|f(\mathbf{x})\|_2^2] = \frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial x} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x}) \end{array} \right.$
 $\frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{A}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [(\mathbf{A} + \mathbf{b})^T \mathbf{C}(\mathbf{D}(\mathbf{x}) + \mathbf{e})] = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} + \mathbf{b}) + \mathbf{A}^T \mathbf{C}(\mathbf{D} + \mathbf{e}) \\ \frac{\partial}{\partial x} [\|f(\mathbf{x})\|_2^2] = \frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial x} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x}) \end{array} \right.$
 $\frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{A}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [(\mathbf{A} + \mathbf{b})^T \mathbf{C}(\mathbf{D}(\mathbf{x}) + \mathbf{e})] = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} + \mathbf{b}) + \mathbf{A}^T \mathbf{C}(\mathbf{D} + \mathbf{e}) \\ \frac{\partial}{\partial x} [\|f(\mathbf{x})\|_2^2] = \frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial x} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x}) \end{array} \right.$
 $\frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{A}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [(\mathbf{A} + \mathbf{b})^T \mathbf{C}(\mathbf{D}(\mathbf{x}) + \mathbf{e})] = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} + \mathbf{b}) + \mathbf{A}^T \mathbf{C}(\mathbf{D} + \mathbf{e}) \\ \frac{\partial}{\partial x} [\|f(\mathbf{x})\|_2^2] = \frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial x} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x}) \end{array} \right.$
 $\frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{A}] = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \\ \frac{\partial}{\partial x} [\mathbf{x}^T \mathbf{x}] = 2\mathbf{x} \\ \frac{\partial}{\partial x} [\mathbf{b}^T \mathbf{A} \mathbf{x}] = \mathbf{A}^T \mathbf{b} \end{array} \right. \quad \left| \begin{array}{l} \frac{\partial}{\partial x} [(\mathbf{A} + \mathbf{b})^T \mathbf{C}(\mathbf{D}(\mathbf{x}) + \mathbf{e})] = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} + \mathbf{b}) + \mathbf{A}^T \mathbf{C}(\mathbf{D} + \mathbf{e}) \\ \frac{\partial}{\partial x} [\|f(\mathbf{x})\|_2^2] = \frac{\partial}{\partial x} [f(\mathbf{x})^T \mathbf{f}(\mathbf{x})] = 2 \frac{\partial}{\partial x} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2J_f(\mathbf{x}) \end{array} \right.$
 $\frac{\partial}{\partial x} [\mathbf{a}^T \mathbf{x}] = \frac{\partial \mathbf{a}}{\partial x} \mathbf{x}^T = \mathbf{a} \quad \left| \begin{array}{l} \frac{\$

12.8 — Channels

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



13.2.1 — 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 decaying learning rate, s.t.

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

we get far enough our steps get always smaller

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

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

Then, we can get the following convergence rates:

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

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

14. Natural Language Processing

Similarities between text and image processing: local information. Differences between text and image processing: texts have various lengths, texts may have long-term interactions, language is a man-made conception on how to communicate with each other / pictures capture the reality, pictures capture the reality / sentences may mean different things in different contexts

14.1 — Word Embeddings

Basic Idea: Map symbols over a vocabulary V to a vector representation — embedding into an (euclidean) vector space (see lookup table in architecture overview).

embedding map: (vocabulary) $V \mapsto \mathbb{R}^d$ (embeddings)
 (symbolic) $v \mapsto \mathbf{x}_v$ (quantitative)

word $w \in V \rightarrow$ one-hot $w \in \{0, 1\}^{|V|} \rightarrow$ embedding \mathbf{x}_w .

$m := |V|$, usually $|V| = 10^5$

d = dimensionality of embedding, $d \ll m$

So for each of the m words in V we have a corresponding embedding in \mathbb{R}^d , which can be stored in a shared lookup table: $\mathcal{R}^{d \times m}$ shared lookup table

Any sentence of k words can then be represented as a $d \times k$ matrix (a sequence of k embedding vectors in \mathbb{R}^d).

Now, how should an embedding be? Ideally, the embedding carries the information/structure that we need in order to go from the input text to the question that we want to solve. Typical questions are:

- Clustering based on context (co-occurrence)
- Sentiment analysis (group words according to mood/feelings)
- Translation (group by meaning)
- Part-of-Speech tagging (understand the structure of text, e.g., location, time, actor, ... or, noun, verb, adjective, ...)

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)

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

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

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 looking at co-occurrences of words within a window size R (instead of looking at subsequences of some length n 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. (Dilated Convolutions) To handle long sequences without losing resolution (pooling), *dilation* introduces gaps between kernel elements.

Receptive Field: Grows exponentially with the dilation factor d ($1, 2, 4, 8, \dots$), allowing the network to capture long-range dependencies with few layers.

D. (WaveNet [2016]) A deep generative model for raw audio waveforms.

Causal Convolutions: Strict ordering ensures the prediction at time t only depends on samples x_{t-1}, \dots, x_1 (cannot see the future).

Dilated Causal Convolutions: Stacks layers with increasing receptive field of the previous time step (millions of audio) to generate realistic high-fidelity sound structure.

Skip Connections: Uses residual and parameterized skip connections to speed up convergence and allow training of very deep networks.

14.2.1 — 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: $K(p+1)$

#parameters of fully-connected NN with same number of outputs as CNN: $mnc(m-p+1)(n-q+1)+1K$

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

13. Optimization

13.1 — Objectives as Expectations

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

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

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

13.2 — Optimization Challenges in NNs: Curvatures

The main stages in the “ConvNet” architecture aka “time-delay NN with max-over-time pooling” are:

1. **Embed Words:** Map word sequence (n words) to a sequence of embedded vectors in \mathbb{R}^d . Store these into the *sentence matrix* $\mathcal{R}^{d \times N}$.

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

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

Now our approach to approximate the probability $P_\theta(w_i | w_j)$ as follows: it should be something that is related to the dot product of the embeddings, $\mathbf{x}_w^T \mathbf{z}_{w_j}$ (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 $\mathbf{w} = (w_1, \dots, w_T)$:

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

partition function

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

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

14.2.1 — Convergence Rates

Under certain conditions SGD converges to the optimum:

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

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

we get far enough our steps get always smaller

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

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

Then, we can get the following convergence rates:

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

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

14.2.2 — Dynamic CNNs

Kalchbrenner et al. suggested Dynamic CNNs in 2014 (as an alternative to ConvNets). They are exactly the same as ConvNets except for one thing: before doing the max-pooling over time (to get a fixed-size representation), they do a dynamic max-pooling (dynamic since it depends on the input size) over the sentence. Note that \mathbf{C} is actually symmetric (as it represents the joint probabilities), but the probabilities that we're computing are asymmetric (conditional probability).

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

14.2.3 — Recurrent Networks (RRNs)

Disadvantage of CNNs: capture better the time component, lossy memorization of past in hidden state.

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 with the state of a dynamical system. The discrete time evolution of the *hidden state sequence* is expressed as a HMM with a non-linearity.

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

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

and the opposite event is given by:

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

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

$$\begin{aligned} \theta^* = \arg \max_{\theta} & \sum_{(i,j) \in D} \log(P_\theta(w_i | w_j = 1 \mid \mathbf{x}_w, \mathbf{z}_{w_j})) \\ & + \sum_{(i,j) \in \bar{D}} \log(1 - P_\theta(w_i | w_j = 0 \mid \mathbf{x}_w, \mathbf{z}_{w_j})) \\ = \arg \max_{\theta} & \sum_{(i,j) \in D} \log(\sigma(\mathbf{x}_w^T \mathbf{z}_{w_j})) + \sum_{(i,j) \in \bar{D}} \log(\sigma(-\mathbf{x}_w^T \mathbf{z}_{w_j})) \\ = \arg \max_{\theta} & \sum_{(i,j) \in C_R} \log(\sigma(\mathbf{x}_w^T \mathbf{z}_{w_j})) + k \cdot \mathbb{E}_{w_i \sim p_n} [\log(\sigma(-\mathbf{x}_w^T \mathbf{z}_{w_i}))] \end{aligned}$$

positive examples negative examples

As we can see in the last step, instead of computing these sums separately, we'll do it as follows: we'll compute the log likelihood for one concrete positive example (w_i, w_j) and produce some negative examples (w_i, w_j), so we'll approximate the expectation below and weight it a bit higher (oversampling factor).

Now, how should an embedding be? Ideally, the embedding carries the information/structure that we need in order to go from the input text to the question that we want to solve. Typical questions are:

- Clustering based on context (co-occurrence)
- Sentiment analysis (group words according to mood/feelings)
- Translation (group by meaning)
- Part-of-Speech tagging (understand the structure of text, e.g., location, time, actor, ... or, noun, verb, adjective, ...)

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)

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

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

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)
- useful for neural translation
- sizes of sentences might not be the same
- outputted words might have slightly different order
- Note that if 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 $\mathcal{O}(\log(n))$). Gives a single output at the root.

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

$$\mathbf{x}^n = F(\mathbf{h}^{n\text{left}}, \mathbf{h}^{n\text{right}})$$

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 slowly. 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:

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

Finding the normalization constant c_θ might be really complicated, so we can only evaluate probabilities. Further, here we cannot use the log-likelihood, because scaling up \bar{p}_θ leads to an unbounded likelihood.

So the question still is: is there an alternative *estimation method* for unnormalized models?

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 use. This is already in use 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 without knowing what our final goal is.

Finally, Hyvonen 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 (Hyvonen 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

D. (Factor Analysis)

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

Minimize the criterion

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

or equivalently (by eliminating ψ by integration by parts)

$$J(\theta) = \mathbb{E} \left[\nabla_{\mathbf{A}} \log (\det(\mathbf{A})) \right] \stackrel{!}{=} 0 \iff \mathbf{S} \mathbf{A}^{-1} = \mathbf{I}.$$

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

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:

$$\nabla_{\mathbf{W}} \mathbf{A} = 2\mathbf{W}, \quad \nabla_{\Sigma} \mathbf{A} = \mathbf{I}$$

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

$$\mathbf{S}(\Sigma + \mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} = \mathbf{W}.$$

In general, finding \mathbf{W} is not easy. However, a special case is if we assume that $\Sigma = \sigma^2 \mathbf{I}$ (isotropic reconstruction error/noise) and $\mathbf{W}^\top \mathbf{W} = \text{diag}(\rho_i^2)$, then by Woodbury's formula we have simplified to:

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

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*.

– Linear Factor Analysis

The idea of linear factor analysis is to explain the data through some low-dimensional isotropic gaussian. And the data is mapped/reconstructed through some linear map to/from the lower-dimensional space. The reconstruction is done via a linear map \mathbf{W} and then different gaussian noises are added to the reconstructed vector (via η).

So the *latent variable prior* is $\mathbf{z} \in \mathbb{R}^m$ where

$$\mathbf{z} \sim \mathcal{N}(\mathbf{o}, \mathbf{I})$$

and we have a linear *observation model* for $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = \mu + \mathbf{W} \mathbf{z} + \eta, \quad \eta \sim \mathcal{N}(\mathbf{o}, \Sigma), \quad \Sigma := \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$$

Further note that

- μ and \mathbf{z} are independent
- typically $m \ll n$ (fewer factors than features)

– 16.3 — Latent Variable Models

– 16.3.1 — DeFinetti's Theorem

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*

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

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

Later we'll put our Bayesian priors onto 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*

T. The distribution of the observation model is

$$\mathbf{x} \sim \mathcal{N}(\mu, \mathbf{W} \mathbf{W}^\top + \Sigma).$$

– Non-Identifiability of Factors

Now this seems to be nice, but again we have the *non-identifiability problem*, since there exist an infinite amount of solutions for any \mathbf{W} that is a solution. Just let \mathbf{Q} be an orthogonal $m \times m$ -matrix. Then $\mathbf{W}\mathbf{Q}$ is also a solution, because

$$(\mathbf{W}\mathbf{Q})(\mathbf{W}\mathbf{Q})^\top = \mathbf{W}\mathbf{W}^\top \mathbf{Q}^\top \mathbf{Q} = \mathbf{W}\mathbf{W}^\top = \Sigma.$$

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)

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

– 16.4 — Dimensionality Reduction

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

$$\mathbf{x} = f(\mathbf{z}\mathbf{B})$$

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 "psomolu" or Bernoulli model, and bot \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})

– 16.4.4 — 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. $\pi(\mathbf{z}) = \mathcal{N}(\mathbf{o}, \mathbf{I})$
- parametrized mechanism: $F: \mathbb{R}^d \rightarrow \mathbb{R}^m$
- induced distribution $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{x} \sim p_\theta(\mathbf{x})$
- sampling is easy: random vector + forward propagation.

So, if we know \mathbf{W} and Σ is assumed to be isotropic with the error going to zero the encoding distribution gets very easy to compute.

– Maximum Likelihood Estimation

Now, how do we estimate \mathbf{W} and Σ ? The idea is fairly simple. Now, the question is what strategies can we use for more complex models.

– 17. Chain-Rule and Jacobians for Tensors

T. (Chain-Rule and Jacobians for Tensors) $\mathbf{T} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_k}$

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

D. (Tensor Multiplication)

Let's assume that $\mathbf{x}_1, \dots, \mathbf{x}_k$ i.i.d. $\mathcal{N}(\mathbf{o}, \mathbf{A})$. Further let's define the data matrix \mathbf{X} as

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

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:

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

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 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:

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

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 use. This is already in use 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 without knowing what our final goal is.

Finally, Hyvonen 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 (Hyvonen 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)

$$J(\theta) = \mathbb{E} \left[\nabla_{\mathbf{A}} \log (\det(\mathbf{A})) \right] \stackrel{!}{=} 0 \iff \mathbf{S} \mathbf{A}^{-1} = \mathbf{I}.$$

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

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:

$$\nabla_{\mathbf{W}} \mathbf{A} = 2\mathbf{W}, \quad \nabla_{\Sigma} \mathbf{A} = \mathbf{I}$$

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

$$\mathbf{S}(\Sigma + \mathbf{W} \mathbf{W}^\top)^{-1} \mathbf{W} = \mathbf{W}.$$

In general, finding \mathbf{W} is not easy. However, a special case is if we assume that $\Sigma = \sigma^2 \mathbf{I}$ (isotropic reconstruction error/noise) and $\mathbf{W}^\top \mathbf{W} = \text{diag}(\rho_i^2)$, then by Woodbury's formula we have simplified to:

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

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*.

– Linear Factor Analysis

The idea of linear factor analysis is to explain the data through some low-dimensional isotropic gaussian. And the data is mapped/reconstructed through some linear map to/from the lower-dimensional space. The reconstruction is done via a linear map \mathbf{W} and then different gaussian noises are added to the reconstructed vector (via η).

So the *latent variable prior* is $\mathbf{z} \in \mathbb{R}^m$ where

$$\mathbf{z} \sim \mathcal{N}(\mathbf{o}, \mathbf{I})$$

and we have a linear *observation model* for $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = \mu + \mathbf{W} \mathbf{z} + \eta, \quad \eta \sim \mathcal{N}(\mathbf{o}, \Sigma), \quad \Sigma := \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$$

Further note that

- μ and \mathbf{z} are independent
- typically $m \ll n$ (fewer factors than features)

– 16.3 — Latent Variable Models

– 16.3.1 — DeFinetti's Theorem

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*

T. The distribution of the observation model is

$$\mathbf{x} \sim \mathcal{N}(\mu, \mathbf{W} \mathbf{W}^\top + \Sigma).$$

– Non-Identifiability of Factors

Now this seems to be nice, but again we have the *non-identifiability problem*, since there exist an infinite amount of solutions for any \mathbf{W} that is a solution. Just let \mathbf{Q} be an orthogonal $m \times m$ -matrix. Then $\mathbf{W}\mathbf{Q}$ is also a solution, because

$$(\mathbf{W}\mathbf{Q})(\mathbf{W}\mathbf{Q})^\top = \mathbf{W}\mathbf{W}^\top \mathbf{Q}^\top \mathbf{Q} = \mathbf{W}\mathbf{W}^\top = \Sigma.$$

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)

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

– 16.4 — Dimensionality Reduction

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

$$\mathbf{x} = f(\mathbf{z}\mathbf{B})$$

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