

## 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  $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_{i=1}^n P(X|Y_i)P(Y_i)}$

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

$E[X] = f_X(t) \cdot dt: \mu_X$   
 $\text{Cov}(X, Y) = E_{x,y}[X - E_x[X](Y - E_y[Y])]$   
 $\text{Cov}(X, X) := \text{Var}(X)$   
 $X, Y \text{ independent} \Rightarrow \text{Cov}(X, Y) = 0$   
 $\mathbf{XX}^T \geq 0$  (symmetric positive semi-definite)

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

**T.** The moment generating function (MGF)  $\psi_X(t) = E[e^{tX}]$  characterizes the distr. of a rv  
 $B(p) = p^t + (1-p) \quad N(\mu, \sigma) = \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$   
 $\text{Bin}(n, p) = (pe + (1-p))^n \quad \text{Gam}(\alpha, \beta) = (\frac{\alpha}{\alpha + \beta t})^\alpha$  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  $Y$ :  $f_Z(z) = \int_{\mathbb{R}} f_X(t) f_Y(z-t) dt = \int_{\mathbb{R}} f_X(z-t) f_Y(t) dt$

$N(\mathbf{x}, \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu))$   
 $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \Sigma = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$

**T.**  $\text{P}(|\mathbf{x}|_F) = \mathcal{N}(\frac{\mathbf{x}_1}{\|\mathbf{x}\|}, \frac{\mathbf{x}_2}{\|\mathbf{x}\|}, \dots, \frac{\mathbf{x}_d}{\|\mathbf{x}\|})$   
 $\mathbf{a}_1, \mathbf{u}_1 \in \mathbb{R}^d, \Sigma_{11} \in \mathbb{R}^{d \times d}$  p.s.d.,  $\Sigma_{12} \in \mathbb{R}^{d \times d}$  p.s.d.  
 $\mathbf{a}_2, \mathbf{u}_2 \in \mathbb{R}^d, \Sigma_{22} \in \mathbb{R}^{d \times d}$  p.s.d.,  $\Sigma_{21} \in \mathbb{R}^{d \times d}$  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)  $\forall u, v \in V: |\langle u, v \rangle| \leq \|u\| \|v\|$ .

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

**T.** Special cases:  $(\sum x_i y_i)^2 \leq (\sum x_i^2)(\sum y_i^2)$ . 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  $\forall \mathbf{x}, \mathbf{x}' \in S, \forall \lambda \in [0, 1]: \lambda \mathbf{x} + (1 - \lambda)\mathbf{x}' \in S$ .  
 Conv. Any point on the line between two points is within the set.  $\mathbb{R}^d$  is convex.

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

- $f(y) \geq f(x) + \nabla f(x)^T(y - x)$
- $f'(x) \geq 0$
- Local minima are global minima, strictly convex functions have a unique global minimum
- If  $f, g$  are convex then  $\alpha f + \beta g$  is convex for  $\alpha, \beta \geq 0$
- If  $f, g$  are convex then  $\max(f, g)$  is convex
- If  $f$  is convex and  $g$  is convex and non-decreasing then  $g \circ f$  is convex

**D.** (Strongly Convex Function) A function  $f$  is  $\mu$ -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  $\kappa$  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.** (Kullback-Leibler Divergence) Given two metric spaces  $(X, d_X)$  and  $(Y, d_Y)$ , a function  $f: X \rightarrow Y$  is called Lipschitz continuous, if there exists a real constant  $L \in \mathbb{R}_0^+$  (Lipschitz constant), such that  $\forall x_1, x_2 \in X: \|f(x_1) - f(x_2)\| \leq L \cdot \|x_1 - x_2\|$ .  
 Com. If the objective function is  $L$ -smooth a step size of  $\eta = 1/L$  guarantees convergence.

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

**D.** (PL Condition) A differentiable function  $f(x)$  with global minimum  $f^*$  satisfies the  $\mu$ -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  $\mathbf{x} \in \mathbb{R}^n: \mathbf{x}^T \mathbf{Ax} \geq 0$  Properties:
 

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

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

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

Otherwise: indefinite.

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

**Properties:**

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

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

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

**Properties:**

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

## 4 Derivatives

**4.1 Numerator and Denominator Convention**

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

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

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

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

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

**4.2 Scalar-by-Vector**

**Denominator Convection**

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

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

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

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

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

$$\frac{\partial \mathbf{x}^T \mathbf{A}}{\partial \mathbf{x}} = \left( \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \mathbf{A} \right)^T = (\mathbf{A}^T \mathbf{x})$$

$$\frac{\partial \mathbf{x}^T \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{x}$$

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

$$\frac{\partial \mathbf{b}^T \mathbf{x}}{\partial \mathbf{x}} = \left( \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \mathbf{b} \right)^T = (\mathbf{b}^T \mathbf{x})$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}^T \mathbf{C}(\mathbf{D} \mathbf{x} + \mathbf{e})}{\partial \mathbf{x}} = \mathbf{D}^T \mathbf{C}^T (\mathbf{A} \mathbf{x} + \mathbf{b}) + \mathbf{A}^T \mathbf{C} (\mathbf{D} \mathbf{x} + \mathbf{e})$$

$$\frac{\partial \mathbf{a}^T \mathbf{f}(\mathbf{x})^T \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = 2 \frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x}) = 2 \mathbf{J}_f \mathbf{f}(\mathbf{x})$$

**4.3 Vector-by-Vector**

**4.4 Scalar-by-Matrix**

**4.5 Vector-by-Matrix (Generalized Gradient)**

**5 Information Theory**

**T.** (Entropy) Let  $X$  be a random variable distributed according to  $p(\mathbf{x})$ . Then the entropy of  $X$

$$H(X) = E[-\log(p(\mathbf{x}))] = -\sum_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \log(p(\mathbf{x})) \geq 0.$$

It describes the expected information content  $I(X)$  of  $X$ .

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

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

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

$$|f(x) - P(x)| < \epsilon, \quad \forall x \in K$$

**6 Activation Functions**

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

**D.** (Tanh)

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

**D.** (ReLU)

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

**D.** (Sigmoid/Logistic)

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

**D.** (Softmax)

$$\text{softmax}(\mathbf{x})_i = f(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

**7 Connectionism**

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

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

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

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

## 8 Regression

**8.1 Linear Regression**

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

**Assymptotic Shattering:**

$$C(s, n) = \begin{cases} 1 - \mathcal{O}(e^{-n}) & \text{if } s \leq n \\ \frac{1}{2} & \text{if } n < s < 2n \\ \frac{1}{2} & \text{if } s = 2n \\ \mathcal{O}(e^{-n}) & \text{otherwise} \end{cases}$$

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

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

## 9 Approximation Theory

**9.1 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

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

**D.** (Scalar-by-Map)

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

$$\frac{\partial \mathbf{a}^T \mathbf{x}^T \mathbf{b}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{b} = \mathbf{a}^T \frac{\partial \mathbf{x}^T}{\partial \mathbf{x}} \mathbf{b} = \mathbf{a}^T \mathbf{b}$$

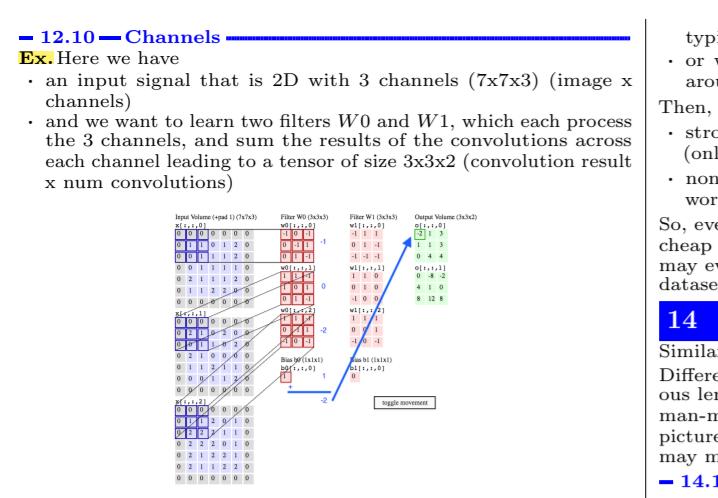
$$\frac{\partial \mathbf{a}^T \mathbf{x}^T \mathbf{A}(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{A}(\mathbf{x}) + \mathbf{A}^T \frac{\partial \mathbf{x}^T}{\partial \mathbf{x}} \mathbf{x} = \mathbf{a}^T \mathbf{b}$$

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

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

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

$$\frac{\partial \$$



- typically  $\eta_t = Ct^{-\alpha}$ ,  $\frac{1}{2} < \alpha < 1$  (c.f. harmonic series.)  
• or we use iterate (Polyak) averaging (once we start jumping around, we average the solutions over time).  
Then, we can get the following convergence rates:  
• strongly-convex case: can achieve a  $\mathcal{O}(1/t)$  suboptimality rate (only polynomial convergence)  
• non-strongly convex case:  $\mathcal{O}(1/\sqrt{t})$  suboptimality rate (even worse than polynomial convergence)

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

## 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).

So the idea is to use two different latent vectors  $\mathbf{x}_w$  and  $\mathbf{z}_w$  per word (to allow for the asymmetry for the conditional probability)  $\theta = (\mathbf{x}_w, \mathbf{z}_w)_{w \in V}$

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

Note that  $C$  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  $V$  (can be  $\sim 10^5$  up to  $10^7$ ). And we'd have to do this every time we pass a new batch sample ( $w_{i+1}, w_i$ ) through the network.

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

For this reason we'll introduce the following function

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

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

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

and the opposite event is given by:

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

One thing to note about this architecture is that most of the parameters were in the embedding.

### 14.2.2 Dynamic CNNs

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

### 14.3 Recurrent Networks (RNNs)

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

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

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

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

$\mathbf{h}^t = F(\mathbf{h}^{t-1}, \mathbf{x}^t; \theta), \quad \mathbf{h}^0 = \mathbf{o}$

$F := \sigma \circ \bar{F}, \quad \sigma \in \{\text{logistic, tanh, ReLU, } \dots\}$

$\bar{F}(\mathbf{h}, \mathbf{x}; \theta) := \mathbf{W}\mathbf{h} + \mathbf{U}\mathbf{x} + \mathbf{b}$ ,

$\mathbf{y}^t = H(\mathbf{h}^t; \theta) := \sigma(\mathbf{V}\mathbf{h}^t + \mathbf{c})$ ,

There are two scenarios for producing outputs

1. Only one output at the end:

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

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

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

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

This ensures that the gradient norm is never greater than  $\gamma_{\max}$ . However, the hidden state sequence is  $\mathbf{h}^1, \dots, \mathbf{h}^T$ , which means that the RNN is forgetting the past after a few timesteps, and usually then the RNN is not performing very well. This is harder to fix, and we'll see later how this is solved with LSTM.

14.3.1 Backprop Over Time

For multi-output loss

$L = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{y}_t} \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t} \sum_{k=t+1}^T \left[ \prod_{h=k+1}^T \frac{\partial \mathbf{h}_k}{\partial \mathbf{h}_{k-1}} \right] \frac{\partial \mathbf{h}_t}{\partial \theta}$

14.3.2 Bi-Directional RNNs

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

$\mathbf{g}^t := G(\mathbf{x}^t, \mathbf{g}^{t+1}; \theta) = \sigma(\mathbf{Px}^t + \mathbf{Qg}^{t+1} + \mathbf{d}), \quad \text{with } \mathbf{g}^T = \mathbf{o}$

the function  $F$  to compute the normal order sequence stays the same, and the output transform  $H$  becomes now a function of both hidden states:

$\mathbf{y}^t = H(\mathbf{h}^t, \mathbf{g}^t; \theta)$

The nice thing is that we can compute both sequences (the forward and backward sequence) in parallel (or independently) - just verify this in the graph. The back-propagation through time is done in reverse order for the reverse order sequence.

14.3.3 Deep Recurrent Networks

Deep recurrent networks (DRNNs) just use a deeper network for the evolution function. So we have hierarchical hidden states. Note, that this can also be combined with bi-directionality.

14.3.4 Probability Distributions over Sequences

Goal: Define a conditional probability distribution over output sequence  $\mathbf{y}^{1:T}$ , given input sequence  $\mathbf{x}^{1:T}$ .

So the idea would be to do a step-by-step prediction:

$P(\mathbf{y}^{1:T} | \mathbf{x}^{1:T}) \approx \prod_{t=1}^T P(\mathbf{y}^t | \mathbf{x}^{1:t}, \mathbf{y}^{1:(t-1)})$

Now, in the naive RNN implementation  $P(\mathbf{y}^t)$  only depends on  $\mathbf{y}^{1:(t-1)}$  through  $\mathbf{h}^t$  since

$\mathbf{x}^{1:t} \xrightarrow{f} \mathbf{h}^t \xrightarrow{h} \mathbf{y}^t \mapsto P(\mathbf{y}^t)$

which determines how much we want to input into the memory so far. In many cases, this is just mapped as  $\mathbf{i}^t = 1 - \mathbf{f}^t$ . But this model illustrates the flexibility that we may want to keep. Further, given that we'll be opening the gate, what should we store there? This is what is computed through

$\mathbf{C}^t = \tanh(\mathbf{W}_C \mathbf{h}^{t-1} + \mathbf{U}_C \mathbf{x}^t + \mathbf{b}_C)$

Using the previous output vector, and the current input it computes some weights, which are used to multiply the content of the memory (by a number between 0 and 1) in order to selectively forget some entries in the memory.

• **Store/Input Gate:** Here we'll compute

$\mathbf{i}^t := \sigma(\mathbf{W}_I \mathbf{h}^{t-1} + \mathbf{U}_I \mathbf{x}^t + \mathbf{b}_I)$  (input)

which determines how much we want to input into the memory so far. However, no real-world applications have been developed with this so-far.

15.3 Attention Mechanisms

D. (Attention Mechanisms) offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

• RNNs: learn to encode information relevant for the future.

• Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

• **Ex:** If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

• **Store/Input Gate:** Here we'll compute

$\mathbf{i}^t := \sigma(\mathbf{W}_I \mathbf{h}^{t-1} + \mathbf{U}_I \mathbf{x}^t + \mathbf{b}_I)$  (input)

These operations are all differentiable, as we're using probabilities and not only the values 0 or 1. Other resembling architectures are: neural random access machines, differentiable data structures (stacks, queues).

Now, NMT architectures can learn loops and simple programs.

However, no real-world applications have been developed with this so-far.

15.3.4 Attention Mechanisms

D. (Attention Mechanisms) offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

• RNNs: learn to encode information relevant for the future.

• Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

• **Ex:** If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

• **Store/Input Gate:** Here we'll compute

$\mathbf{i}^t := \sigma(\mathbf{W}_I \mathbf{h}^{t-1} + \mathbf{U}_I \mathbf{x}^t + \mathbf{b}_I)$  (input)

These operations are all differentiable, as we're using probabilities and not only the values 0 or 1. Other resembling architectures are: neural random access machines, differentiable data structures (stacks, queues).

Now, NMT architectures can learn loops and simple programs.

However, no real-world applications have been developed with this so-far.

15.3.5 Attention Mechanisms

D. (Attention Mechanisms) offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

• RNNs: learn to encode information relevant for the future.

• Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

• **Ex:** If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

• **Store/Input Gate:** Here we'll compute

$\mathbf{i}^t := \sigma(\mathbf{W}_I \mathbf{h}^{t-1} + \mathbf{U}_I \mathbf{x}^t + \mathbf{b}_I)$  (input)

These operations are all differentiable, as we're using probabilities and not only the values 0 or 1. Other resembling architectures are: neural random access machines, differentiable data structures (stacks, queues).

Now, NMT architectures can learn loops and simple programs.

However, no real-world applications have been developed with this so-far.

15.3.6 Attention Mechanisms

D. (Attention Mechanisms) offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

• RNNs: learn to encode information relevant for the future.

• Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

• **Ex:** If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

• **Store/Input Gate:** Here we'll compute

$\mathbf{i}^t := \sigma(\mathbf{W}_I \mathbf{h}^{t-1} + \mathbf{U}_I \mathbf{x}^t + \mathbf{b}_I)$  (input)

These operations are all differentiable, as we're using probabilities and not only the values 0 or 1. Other resembling architectures are: neural random access machines, differentiable data structures (stacks, queues).

Now, NMT architectures can learn loops and simple programs.

However, no real-world applications have been developed with this so-far.

15.3.7 Attention Mechanisms

D. (Attention Mechanisms) offer a simple way to overcome some challenges of RNN-based memorization. With attention mechanisms we selectively attend to *inputs* or *feature representations* computed from inputs.

• RNNs: learn to encode information relevant for the future.

• Attention: selects what is relevant from the past in hindsight!

Both ideas can be combined!

• **Ex:** If we have a sentence in English and one in German the question is how do we match one to the other. The problem with CTC was that if things are changed in order, then CTC cannot deal with it. Because the CTC doesn't process every input before it produces an output. Attention will provide a mechanism to deal with this.

### 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)
- learned hidden states based on attention
- size of words 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  $O(\log(n))$ ). Gives a single output at the root.

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

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

## 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 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  $p_\theta$  for a parametric function makes this task quite complicated.

So when using a *prescribed model*  $p_\theta$  we have to

- ensure that  $p_\theta$  defines a proper density:

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

and to be able to evaluate the density  $p_\theta$  at various sample points  $\mathbf{x}$

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

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

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

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

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

The problem with this is that the rate of convergence is  $\log(\log(n))$  - this is extremely painfully slow. This is just a guarantee in general when we know nothing about our density.

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

$$\underbrace{\bar{p}_\theta(\mathbf{x})}_{\text{represented unknowns}} = \underbrace{c_\theta}_{\text{normalization}} \cdot \underbrace{p_\theta(\mathbf{x})}_{\text{}}$$

Finding the normalization constant  $c_\theta$  might be really complicated, so we can only evaluate relative probabilities. Further, here we cannot use the log-likelihood, because scaling up  $\bar{p}_\theta$  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  $p_\theta$ , but we look for properties of  $p_\theta$ . In many cases these properties depend on our prior knowledge of  $p_\theta$ . We need to understand what the problem is in order to put the prior knowledge into the model that we want to do. This was already important in supervised learning (e.g., CNNs with several layers for images), and is even more important in unsupervised learning. We have to do the same thing there without knowing what our final goal is.

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

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

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

#### 16.2.1 - Latent Variable Analysis

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

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

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

### Ex. (Gaussian Mixture Models GMMs)

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

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

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

The idea is to have some

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

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

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

### 16.2.2 - 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  $\eta$ ).

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

- $\mu$  and  $\mathbf{z}$  are *independent*
- typically  $m \ll n$  (fewer factors than features)
- few factors account for the dependencies between many observables

• The vector  $\mu$  is computed through MLE on the training set

$$\hat{\mu} = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i$$

Usually we assume centered data, so  $\mu = \mathbf{o}$ . Since  $\mu$  only complicates the notation and is actually easy to determine. Recall, that in the previous part when we were doing autoencoders, the deviations that we were having for each of the components was the same. So we wanted the error to be the same for each of the components. Now, with this model, with  $\eta$  we are allowing for additional flexibility for the error. There will be some components that we'll be able to explain with less error, and some with more. So,  $\mathbf{z}$  should capture everything that is important to explain the data, and  $\eta$  can be viewed as noise.

Although we're assuming that here everything is gaussian, in general we may view  $\mathbf{z}$  as a clustering mechanism, where  $\mathbf{z}$  determines some cluster components that are selected.

T. The distribution of the *observation model* is

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

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 - 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*

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

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.

### 16.3.1 - 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

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

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

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

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

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.

### 16.3.3 - Dimensionality Reduction

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

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

where

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

and

$$\mathbf{z} | \mathbf{x} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

so

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

and to be able to evaluate the density  $p_\theta$  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:

$$\mathbf{W}^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^\top = \mathbf{W}^\top \mathbf{W}.$$

Consequently with the assumption of zero reconstruction error:

$$\mathbf{z} | \mathbf{x} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

so

$$\mathbf{z} | \mathbf{x} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

and

$$\mathbf{z} | \mathbf{x} \sim \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}}),$$

and to be able to evaluate the density  $p_\theta$  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.