

1 Probability

Sum Rule $P(X = x_i) = \sum_j p(X = x_i, Y = y_i)$
 Product rule $P(X, Y) = P(Y|X)P(X)$
 Independence $P(X, Y) = P(X)P(Y)$
 Bayes' Rule $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_i P(X|Y_i)P(Y_i)}$

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

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

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

$\text{Cov}(X, Y) := E_{x,y}[(X - E_x[X])(Y - E_y[Y])]$

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

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

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

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

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

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

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

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

$f_{\alpha,Y}(y) = \frac{1}{\alpha} f_Y(\frac{y}{\alpha})$

T: The moment generating function (MGF) $\psi_X(t) = E[e^{tX}]$ characterizes the distr. of r v.s.

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

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

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

$(H(X, p) = H(X) + KL(p, q) \geq 0, \text{ where } H \text{ uses } p$.

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

Com. Usually, q is the approximation of the unknown p .

D: (**Kullback-Leibler Divergence**)

For probability distributions p and q defined on the same probability space, the KL-divergence between p and q is defined as

$KL(p, q) = -\sum_{x \in X} p(x) \log\left(\frac{q(x)}{p(x)}\right) = \sum_{x \in X} p(x) \log\left(\frac{p(x)}{q(x)}\right) \geq 0$

$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{q} = \frac{1}{n} \sum_{i=1}^n (x - \hat{p})(x - \hat{p})^T$

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

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

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

T: $P(\left[\begin{smallmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{smallmatrix}\right] = \left[\begin{smallmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{smallmatrix}\right]) = N(\left[\begin{smallmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{smallmatrix}\right] | \left[\begin{smallmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{smallmatrix}\right], \left[\begin{smallmatrix} \sum_{i=1}^n \mathbf{x}_i^2 \\ \vdots \\ \sum_{i=1}^n \mathbf{x}_i^2 \end{smallmatrix}\right])$

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

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

T: (**Chébyshev**) 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**)

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

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

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

Special case: $E[X^2]^2 \leq E[X^2]E[Y^2]$.

D: (**Convex Set**) A set $S \subseteq \mathbb{R}^d$ is called *convex* if

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

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

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

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

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.

T: (**Properties of Convex Functions**)

• 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

T: (**Taylor-Lagrange Formula**)

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

T: (**Jensen-Shannon Divergence**)

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

$M = \frac{1}{2}(P + Q)$

C: The JSD is symmetric!

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

D: (**Activation Functions and Their Derivatives**)

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**)

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

D: (**Sigmoid/Logistic**)

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

• Knots: $0 = x_0 < x_1 < \dots < x_m - 1 < x_m = 1$

• 1st + 1 parameters: $a, b, c_1, \dots, c_{m-1}$

With the dimension lifting theorem we can lift this property from ID to nD.

Note that there's an alternative representation of the above through AVUs

$g(\mathbf{x}) = a' \mathbf{x} + b' + \sum_{i=1}^{m-1} -1 c'_i |x_i - x_i|$

T: Networks with one hidden layer of ReLUs are universal function approximators.

Com. We can thus use a restricted set of activation functions (ReLU). But still we don't know how many hidden units we need.

8.3.1 Piecewise Linear Functions and Half Spaces

So the ReLU and the AVU define a piecewise linear function with 2 pieces. Hereby, \mathbb{R}^d is partitioned into two open half spaces (and a border face):

• $H^+ := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b > 0 \} \subset \mathbb{R}^d$

• $H^- := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b < 0 \} \subset \mathbb{R}^d$

• $H^0 := \{ \mathbf{x} \mid \mathbf{w}^\top \mathbf{x} + b = 0 \} = \mathbb{R}^d - H^+ - H^- \subset \mathbb{R}^d$

Further note that:

• $g_{|+}(H^0) = g_{|+}(H^-) = 0$

• $g_{|-}(H^-) = 0$

• $g_{|+}(\mathbf{x}) = g_{|+}(\mathbf{v} - \mathbf{w})$ with $\mathbf{v} = -2b \frac{\mathbf{w}}{\|\mathbf{w}\|_2^2}$ (mirroring at \mathbf{w}): equivalent to subtracting the projection of \mathbf{x} onto \mathbf{w} twice from \mathbf{x}

Partitions of ReLUs go to infinity, even if we have no examples there → weak to extrapolation errors (adversarial examples). However, if you have enough data, then you can overcome this, because there won't be any new examples that lie outside of the training data regions.

8.3.2 Linear Combinations of ReLUs

T: (**Zaslavsky, 1975**) By linearly combining m rectified units \mathbb{R}^d can be partitioned at most into $R(m)$ cells:

$R(m) \leq \min(m, n) \binom{m}{n}$

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

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

8 Approximation Theory

8.1 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

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

• $\forall \mathbf{x} \in \mathbb{R}^n: f(\alpha \mathbf{x}) = \alpha f(\mathbf{x})$

T: (**Comp. of Lin. Maps**) $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear $\iff f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$

D: (**Hyperplane**) $H = \{ \mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle - p = 0 \} = \{ \mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle = b \}$

where $b = \langle \mathbf{w}, \mathbf{p} \rangle$, \mathbf{w} = normal vector, \mathbf{p} points onto a point on the plane.

D: (**Level Sets**) of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is a one-parametric family of sets defined as

$\{ \mathbf{x} \mid f(\mathbf{x}) = c \} = f^{-1}(c) \subset \mathbb{R}^n$.

T: (**Comp. of Lin. Maps**) f is a Lin. Map/Unit

L_1, \dots, L_d be linear maps then $F = L_d \circ \dots \circ L_1$ is also a linear map.

C: Every L -layer NN or linear layer collapses to a 1-layer NN. Furthermore that hereby

Com. This reduces the growth of absolute value nesting to logarithmic growth, instead of linear growth.

11.13 Networks Similar to CNNs

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

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

Ex: input image $m \times n \times c$ (c = number of channels). K convolution kernels: $p \times q$ (valid padding and stride 1) output dimensions: $(m - p + 1) \times (n - q + 1) \times K$ #parameters CNN: $K(pq + 1)$ #parameters of fully-connected NN with same number of outputs as CNN: $mnc(m - p + 1)(n - q + 1) + 1$ #parameters of locally-connected NN with same connections as CNN: $pqc(m - p + 1)(n - q + 1) + 1$

12 Optimization

12.1 Learning as Optimization

Machine learning uses optimization, but it's *not equal* to optimization for two reasons:

- The empirical risk is only a *proxy* for the expected risk
- The loss function may only be a *surrogate*

12.2 Objectives as Expectations

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

The typical structure of a learning objective in a NN is a *large* *finite* sum (over all training instances). Accuracy-complexity trade-off: in practice we subsample terms in the sum, by using *mini-batches* of the training data (so we'll get something close to the true gradient - but not exactly). The idea behind it, is that everything will work out in expectation. Further, we favour cheap and imprecise computations over many datapoints rather than precise and expensive computations over a few datapoints.

12.3 Gradient Descent

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

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

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

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

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

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

Ex: So for our risk function \mathcal{R} , we say that the gradient of it is $\nabla_{\theta} \mathcal{R}: \Omega \rightarrow \Omega$ where $\Omega = \mathbb{R}^n$ is *L*-Lipschitz continuous if it holds that $\forall \theta_1, \theta_2 \in \Omega: \|\nabla_{\theta} \mathcal{R}(\theta_1) - \nabla_{\theta} \mathcal{R}(\theta_2)\| \leq L \|\theta_1 - \theta_2\|$

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

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

Continuously differentiable \subseteq Lipschitz continuous \subseteq (Uniformly) continuous

For everywhere differentiable function $f: \mathbb{R} \rightarrow \mathbb{R}$ is Lipschitz continuous (with $L = \sup |f'(x)|$) iff it has *bounded first derivatives*.

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

T. If R is convex with *L*-Lipschitz-continuous gradients then we have that

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

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

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

C. Convex combination of two points < evaluation of convex combination of two points.

Com. Another way to formulate that f is convex function is to say that the *epigraph* of f is a convex set.

T. Every local optimum of a convex function is a global optimum.

T. (Operations that Preserve Convexity)

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf{x})$
- composition with affine mapping: $f(Ax + b)$
- restriction to a line (of convex set domain)

D. (Strictly Convex Function) f is called *strictly convex* if

- f is concave if and only if f is convex
- non-negative weighted sums
- point/element-wise maximum $\max(f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$
- composition with non-decreasing function, e.g. $e^f(\mathbf$

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

D. (Pointwise Mutual Information)

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

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

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

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

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

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

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

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

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

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

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

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

de-embedding: $\mathbf{VA}_k^{-\frac{1}{2}}$ (then find nearest neighbour)

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

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

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

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

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

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

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

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

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

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

Now our approach to approximate the probability $P_\theta(w_i \mid w_j)$ as follows: it should be something that is related to the dot product of the embeddings, so $\mathbf{x}_w \cdot \mathbf{z}_w$ (note how we use two different embeddings as the conditional probability is asymmetric), but in order to make the probability positive we'll take the exponential of it and normalize. Further, for SGD it's always better to optimize a sum: so we'll optimize the log-likelihood of co-occurring words in our dataset $\mathbf{w} = (w_1, \dots, w_T)$:

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

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

$$\theta = (\mathbf{x}_w, \mathbf{z}_w) \in \mathbb{R}^{n \times 2}$$

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

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

Note that \mathbf{C} is actually symmetric (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 ~ 10^9 up to 10^7). And we'd have to do that every time we pass a new batch-sample (w_{t+1}, \dots, w_T) through the network.

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

For this reason we'll introduce the following function

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

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

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

and the opposite event is given by:

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

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

$$\theta^* = \arg \max_{\theta} \sum_{(i,j) \in D} \log(P_\theta(D_{w_i, w_j} = 1 \mid \mathbf{x}_{w_i}, \mathbf{z}_{w_j})) + \sum_{(i,j) \in \bar{D}} \log(1 - P_\theta(D_{w_i, w_j} = 0 \mid \mathbf{x}_{w_i}, \mathbf{z}_{w_j}))$$

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

Now an LSTM has two memories:

C is the *memory* in the LSTM (a vector). That's where we store all of our information. The memory we'll control through the gated units.

h_t is our previous output. So it's not our memory, but it's what the sequence model outputs at timestep t . This could then pass further to a prediction network.

For multi-output loss

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

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

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

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

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

There are two scenarios for producing outputs

$$\begin{aligned} L &= \sum_{t=1}^T \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \sum_{t=1}^T \left[\prod_{k=t}^{T-1} \frac{\partial h_k}{\partial h_{k-1}} \right] \frac{\partial h_t}{\partial \theta} \\ \theta^* &= \arg \min_{\theta} L \end{aligned}$$

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

The negative examples are just sampled from the negative sampling distribution $p_n(w)$: for that we determine relative frequencies of words $p(w)$ and dampen them with a factor $\alpha < 1$ (usually: $\alpha = \frac{1}{3}$). Then $p_n(w) = \alpha p(w)$ to increase the chance of true negatives. Even if by chance rare. The factor α just weights the negative examples a bit more in the loss. Usually $\alpha \approx 2$ to 10 (oversampling factor).

1. Only one output at the end:

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

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

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

• **Markov Property:** hidden state at time t depends on input of time $t-1$ as well as the previous hidden state (but we don't need the older hidden states).

• **Time-Invariance:** the state evolution function F is independent of t (it's just parameterized by θ).

Feedforward VS Recurrent Networks: RNNs process inputs in sequence, parameters shared between layers (same H and F at every timestep).

13.2 From Embedding Words to Embedding Sequences of Words

Question: Can we extend word embeddings to embeddings for sequences of words? So what we're after is *understanding the sentence*.

Why is this relevant? This is the fundamental question of *statistical language modeling* (cf. Shannon). So we'd like to estimate the probability of a sequence of words in a certain order:

13.2.1 Deep Recurrent Networks

The backpropagation is straightforward: we propagate the derivatives backwards through time. So, the parameter sharing leads to a sum over t when dealing with the derivatives of the weights:

Algorithm 2: Backpropagation in RNNs

(Blue terms only need to be comp. for multiple-output RNNs)

// Compute derivative w.r.t. outputs

$$\text{Compute } \frac{\partial \mathbf{y}^1}{\partial \mathbf{x}^1}, \frac{\partial \mathbf{y}^2}{\partial \mathbf{x}^2}, \dots, \frac{\partial \mathbf{y}^T}{\partial \mathbf{x}^T} \quad (= \frac{\partial \mathbf{y}^t}{\partial \mathbf{x}^t})$$

// Compute the gradient w.r.t. all hidden states

$$\frac{\partial \mathbf{h}^1}{\partial \mathbf{h}^0} \leftarrow \sum_{i=1}^T \frac{\partial \mathbf{y}^i}{\partial \mathbf{h}^0} \frac{\partial \mathbf{h}^0}{\partial \mathbf{h}^1}$$

for $t = (T-1)$ down to 1 do

$$\frac{\partial \mathbf{h}^i}{\partial \mathbf{h}^{i-1}} \leftarrow \sum_{t=1}^{i-1} \frac{\partial \mathbf{y}^t}{\partial \mathbf{h}^{i-1}} \frac{\partial \mathbf{h}^{i-1}}{\partial \mathbf{h}^i} + \sum_{t=i+1}^T \frac{\partial \mathbf{y}^t}{\partial \mathbf{h}^{i-1}}$$

// Do back-propagation over time for weights and biases

$$\frac{\partial \mathbf{R}}{\partial \mathbf{w}_{ij}} \leftarrow \sum_{t=1}^T \frac{\partial \mathbf{y}^t}{\partial \mathbf{w}_{ij}} \frac{\partial \mathbf{h}^{t-1}}{\partial \mathbf{w}_{ij}}$$

In practice we often use 5-grams, i.e., $k = 4$ (last 4 words).

2. Modern Approach: Create language models via embeddings

log($P(w_t \mid \mathbf{w} := w_{t-1}, \dots, w_1)$) = $\mathbf{x}_w^\top \mathbf{z}_w + \text{const}$

where \mathbf{x}_w is the word embedding

\mathbf{z}_w is the sequence embedding (predicts the context)

There are three main approaches to construct sequence embeddings:

1. CNNs + conceptually simple, fast to train

2. RNNs + active memory management via gated units + more difficult to optimize, larger datasets needed

3. Recursive networks (in combination with parsers)

13.2.1 ConvNets: Word Representations

de-embedding: $\mathbf{VA}_k^{-\frac{1}{2}}$ (then find nearest neighbour)

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

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

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

<math display

T. Given the SVD of the data $\mathbf{X} = \mathbf{U}\text{diag}(\sigma_1, \dots, \sigma_n)\mathbf{V}^\top$. The choice $\mathbf{C} = \mathbf{U}_m^\top$ and $\mathbf{D} = \mathbf{U}_m$ minimizes the squared reconstruction error of a two-layer linear auto-encoder with m hidden units. **Proof.**

$$\mathbf{DCX} = \mathbf{U}_m \mathbf{U}_m^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_m [\mathbf{I}_m \quad \mathbf{0}] \Sigma \mathbf{V}^\top = \mathbf{U}_m [\Sigma_m \quad \mathbf{0}] \mathbf{V}^\top = \mathbf{U}_m \Sigma_m \mathbf{V}^\top (\mathbf{WQ})^\top = M \mathbf{WQ} \mathbf{Q}^\top \mathbf{W}^\top = \mathbf{WW}^\top.$$

And as we know from the Eckart-Young theorem $\tilde{\mathbf{X}} = \mathbf{U}_m \Sigma_m \mathbf{V}_m^\top$ is the best m -dimensional approximation of the original data \mathbf{X} . Now, since $\mathbf{C} = \mathbf{U}_m^\top$ and $\mathbf{D} = \mathbf{U}_m$ that means that we can do weight sharing between the decoder and encoder network, since $\mathbf{C} = \mathbf{D}^\top$.

Another thing to note is that the solution is *not unique!* For any invertible matrix $\mathbf{A} \in GL(m)$ we have

$$\underbrace{(\mathbf{U}_m \mathbf{A}^{-1})}_{\mathbf{D}} (\mathbf{A} \mathbf{U}_m^\top)_{\mathbf{C}} = \mathbf{U}_m \mathbf{U}_m^\top.$$

Now, restricting through weight sharing that $\mathbf{D} = \mathbf{C}^\top$ will enforce that $\mathbf{A}^{-1} = \mathbf{A}^\top$, hence $\mathbf{A} \in O(m)$ (orthogonal group, rotation matrices). Then the mapping $\mathbf{x} \rightarrow \mathbf{z}$ is determined (up some rotation of the components).

– Principal Component Analysis

A way to solve this problem is through PCA. First, we center the data (pre-processing) as follows:

$$\mathbf{x}_i \mapsto \mathbf{x}_i - \sum_{i=1}^k \mathbf{x}_i$$

Then we define

$$\mathbf{S} := \mathbf{XX}^\top$$

which is the sample covariance matrix. And then, in order to get \mathbf{U} we just do the singular value decomposition of \mathbf{S} . If we relate it to the SVD of \mathbf{X} we can see that

$$\mathbf{S} = \mathbf{U}\Sigma\mathbf{V}^\top \mathbf{V}\mathbf{U}^\top = \mathbf{U}\Sigma^2\mathbf{U}^\top.$$

So, the column vectors of \mathbf{U} are the eigenvectors of the covariance matrix. And $\mathbf{U}_m \mathbf{U}_m^\top$ is the orthogonal projection onto m principal components of \mathbf{S} . Note that if we wanted to get \mathbf{V} we'd just do the PCA with $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$.

– 15.2.2 Non-Linear Autoencoders

Non-linear autoencoders allow us to learn powerful non-linear generalizations of the PCA.

D. (Non-Linear Autoencoder) contains many hidden layers with nonlinear-activation functions as we want (as long as there's a bottleneck layer) and train the parameters via MLE.

– 15.2.3 Regularized Autoencoders

One may also regularize the code \mathbf{z} via a regularizers $\Omega(\mathbf{z})$. This will give us a regularized autoencoder.

There are various flavours of regularization:

- standard L_2 penalty: ability to learn “overcomplete” codes
- **D. (Code Sparseness)** e.g., via $\Omega(\mathbf{z}) = \lambda \|\mathbf{z}\|_1$.
- **D. (Contractive Autoencoders)** $\Omega(\mathbf{z}) = \lambda \|\frac{\partial \mathbf{z}}{\partial \mathbf{x}}\|^2$. This penalizes the Jacobian and generalizes weight decay (cf. Rifai et al., 2011)

– 15.2.4 Denoising Autoencoders

Autoencoders also allow us to separate the signal from noise: Denoising autoencoders aim to learn features of the original data representation that are robust under noise.

D. (Denoising Autoencoder) we perturb the inputs $\mathbf{x} \mapsto \mathbf{x}_\eta$, where η is a random noise vector, e.g., additive (white) noise

$$\mathbf{x}_\eta = \mathbf{x} + \eta, \quad \eta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

and instead of the original objective, we minimize the following

$$\mathbb{E}_\mathbf{x} [\mathbb{E}_\eta [l(\mathbf{x}, (H \circ G)(\mathbf{x}_\eta))]]$$

The hope is that we'll achieve *de-noising*, which happens if

$$\|\mathbf{x} - H(G(\mathbf{x}_\eta))\|^2 < \|\mathbf{x} - \mathbf{x}_\eta\|^2$$

So this would mean that the reconstruction error of the noisy data is less than the error we created by the noise we've added (then the de-noising works).

– 15.3 Factor Analysis

– 15.3.1 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

$$p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) \mu(d\mathbf{z}) = \left\{ \int p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) dz \right\} \mu = \text{Lebesgue measure}$$

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

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

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

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{W}\mathbf{z} + \boldsymbol{\eta}, \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \Sigma), \quad \Sigma := \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$$

Further note that

- $\boldsymbol{\mu}$ and \mathbf{z} are *independent*
- typically $m \ll n$ (fewer factors than features)
- so few factors account for the dependencies between many observables
- The vector $\boldsymbol{\mu}$ is computed through MLE on the training set

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

Usually we assume centered data, so $\boldsymbol{\mu} = \mathbf{0}$. Since $\boldsymbol{\mu}$ only complicates the notation and is actually easy to determine.

Recall, that in the previous part when we were doing autoencoders, the decomposition we were after was to be the same for each of the components. Now we're doing this with model \mathbf{y} we're 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 $\boldsymbol{\eta}$ 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

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

– 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{WQ} is also a solution, because

$$\mathbf{DCX} = \mathbf{U}_m \mathbf{U}_m^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_m [\mathbf{I}_m \quad \mathbf{0}] \Sigma \mathbf{V}^\top = \mathbf{U}_m [\Sigma_m \quad \mathbf{0}] \mathbf{V}^\top = \mathbf{U}_m \Sigma_m \mathbf{V}^\top (\mathbf{WQ})^\top = M \mathbf{WQ} \mathbf{Q}^\top \mathbf{W}^\top = \mathbf{WW}^\top.$$

And as we know from the Eckart-Young theorem $\tilde{\mathbf{X}} = \mathbf{U}_m \Sigma_m \mathbf{V}_m^\top$ is the best m -dimensional approximation of the original data \mathbf{X} . Now, since $\mathbf{C} = \mathbf{U}_m^\top$ and $\mathbf{D} = \mathbf{U}_m$ that means that we can do weight sharing between the decoder and encoder network, since $\mathbf{C} = \mathbf{D}^\top$.

Another thing to note is that the solution is *not unique!* For any invertible matrix $\mathbf{A} \in GL(m)$ we have

$$\underbrace{(\mathbf{U}_m \mathbf{A}^{-1})}_{\mathbf{D}} (\mathbf{A} \mathbf{U}_m^\top)_{\mathbf{C}} = \mathbf{U}_m \mathbf{U}_m^\top.$$

Now, restricting through weight sharing that $\mathbf{D} = \mathbf{C}^\top$ will enforce that $\mathbf{A}^{-1} = \mathbf{A}^\top$, hence $\mathbf{A} \in O(m)$ (orthogonal group, rotation matrices). Then the mapping $\mathbf{x} \rightarrow \mathbf{z}$ is determined (up some rotation of the components).

– Principal Component Analysis

Now, how is the factor analysis related to data compression?

Encoder Step: Implicitly defined by posterior distribution

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{z}) p(\mathbf{z})}{p(\mathbf{x})} \quad (\text{Bayes})$$

So,

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

where

$$\boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}} = \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} (\mathbf{x} - \boldsymbol{\mu}),$$

$$\boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}} = \mathbf{I} - \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} \mathbf{W}$$

Further, if we assume that $\Sigma = \sigma^2 \mathbf{I}$ and we let $\sigma^2 \rightarrow 0$ (the reconstruction-error-variance for all the components is the same and we let the reconstruction error go to zero), then the following expression just reduces to the pseudo-inverses:

$$\mathbf{W}^\top (\mathbf{WW}^\top + \sigma^2 \mathbf{I})^{-1} \xrightarrow{\sigma^2 \rightarrow 0} \mathbf{W}^\top \in \mathbb{R}^{m \times n}.$$

Consequently with the assumption of zero reconstruction error:

$$\boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{I} - \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} \mathbf{W}$$

Further, if we assume that $\Sigma = \sigma^2 \mathbf{I}$ and we let $\sigma^2 \rightarrow 0$ (the reconstruction-error-variance for all the components is the same and we let the reconstruction error go to zero), then the following expression just reduces to the pseudo-inverses:

$$\mathbf{W}^\top (\mathbf{WW}^\top + \sigma^2 \mathbf{I})^{-1} \xrightarrow{\sigma^2 \rightarrow 0} \mathbf{W}^\top \in \mathbb{R}^{m \times n}.$$

Consequently with the assumption of zero reconstruction error:

$$\boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{I} - \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} \mathbf{W}$$

Further, if we assume that $\Sigma = \sigma^2 \mathbf{I}$ and we let $\sigma^2 \rightarrow 0$ (the reconstruction-error-variance for all the components is the same and we let the reconstruction error go to zero), then the following expression just reduces to the pseudo-inverses:

$$\mathbf{W}^\top (\mathbf{WW}^\top + \sigma^2 \mathbf{I})^{-1} \xrightarrow{\sigma^2 \rightarrow 0} \mathbf{W}^\top \in \mathbb{R}^{m \times n}.$$

Consequently with the assumption of zero reconstruction error:

$$\boldsymbol{\mu}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}_{\mathbf{z} | \mathbf{x}} \approx \mathbf{I} - \mathbf{W}^\top (\mathbf{WW}^\top + \Sigma)^{-1} \mathbf{W}$$

Further, if we assume that $\Sigma = \sigma^2 \mathbf{I}$ and we let $\sigma^2 \rightarrow 0$ (the reconstruction-error-variance for all the components is the same and we let the reconstruction error go to zero), then the following expression just reduces to the pseudo-inverses:

$$\mathbf{W}^\top (\mathbf{WW}^\top + \sigma^2 \mathbf{I})^{-1} \xrightarrow{\sigma^2 \rightarrow 0} \mathbf{W}^\top \in \mathbb{R}^{m \times n}.$$

Consequently with the assumption of zero reconstruction error:

$$\boldsymbol{\mu}_{\mathbf{z} |$$