

1 Probability

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$B(e, p) = pe + (1-p)$

$N(\mu, \sigma) = \exp(\mu t + \frac{1}{2}\sigma^2 t^2)$

$\text{Bin}(n, p) = (pe + (1-p))^n$

$\text{Gam}(\alpha, \beta) = (\frac{1}{\alpha} \beta t)^{\alpha}$ for $t < 1/\beta$

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

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

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

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

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

T: $P(\left|\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}\right|) = \mathcal{N}\left(\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \mid \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|}, \frac{\boldsymbol{\Sigma}_{11}}{\|\mathbf{x}_i\|^2}\right) \left[\frac{\sum_{i=1}^n \mathbf{x}_i}{\|\mathbf{x}\|} \right]$

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

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

$P(a_2 | a_1 = \mathbf{z}) =$

$N(a_2 | \mathbf{u}_2 + \Sigma_{12} \Sigma_{11}^{-1}(\mathbf{z} - \mathbf{u}_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{21})$

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

2 Analysis

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

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

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

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

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

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

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

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

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

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)

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

$\cdot f'(x) \geq 0$

\cdot Local minima are global minima, strictly convex functions have a unique global minimum

\cdot If f, g are convex then $\alpha f + \beta g$ is convex for $\alpha, \beta \geq 0$

\cdot If f, g are convex then $\max(f, g)$ is convex

\cdot If f is convex and g is convex and non-decreasing then $g \circ f$ is convex

T: (Taylor-Lagrange Formula)

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

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

$f(\sum_i \lambda_i x_i) \leq \sum_i \lambda_i f(x_i)$

Special case: $E[E(X)] \leq E[X]$.

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

$x^{t+1} = x^t - \text{Hess}(f(x^t))^{-1} \nabla_x f(x^t)$.

3 Linear Algebra

Kernels are positive semi-definite matrices.

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

\cdot All eigenvalues $\lambda_i \geq 0$.

\cdot The trace $\text{Tr}(\mathbf{A}) \geq 0$ and determinant $\det(\mathbf{A}) \geq 0$.

\cdot Cholesky Decomposition exists: $\mathbf{A} = LL^T$.

T: (Sylvester Criterion) $A \times d$ 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: undefined.

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

Properties:

\cdot $\text{Tr}(\mathbf{A}) = \sum_i \lambda_i$ (sum of eigenvalues).

\cdot Cyclic property: $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$.

\cdot Linear: $\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$ and $\text{Tr}(cA) = c\text{Tr}(A)$.

\cdot $\text{Tr}(A) = \text{Tr}(A^T)$.

D: (Frobenius Norm) $\|A\|_F = \sqrt{\sum_{i,j} |A_{ij}|^2}$

Properties:

\cdot Relation to Trace: $\|A\|_F = \sqrt{\text{Tr}(A^T A)}$.

\cdot Invariant under orthogonal rotations: $\|QA\|_F = \|A\|_F$ for orthogonal Q .

\cdot Relation to Singular Values: $\|A\|_F = \sqrt{\sum_i \sigma_i^2}$.

4 Derivatives

4.1 Numerator and Denominator Convention

Jacobian Layout Convention

We use the denominator-layout convention for a vector-valued function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined

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

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

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

S: Regression

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

$$\ell(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}))^2$$

Analytically solvable: $w^* = (X^\top X)^{-1} X^\top y$

D: (Ridge Regression): $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ with ℓ_2 regularization

$$\ell(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}))^2 + \lambda \|\mathbf{w}\|^2$$

Analytically solvable: $w^* = (X^\top X + \lambda I)^{-1} X^\top y$

D: (Lasso Regression): $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ with ℓ_1 regularization

$$\ell(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (y_i - f(\mathbf{x}))^2 + \lambda \|\mathbf{w}\|_1$$

D: (Logistic Regression): $f(\mathbf{x}) = \sigma(w^\top \mathbf{x} + b)$

It minimizes the cross-entropy loss (negative log-likelihood), which acts as a surrogate loss for the 0/1 classification error.

D: (Scalar-by-Vector)

Denominator Convention

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

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

D: (Vector-by-Vector)

Goal: Compute the data into m-dim. ($m \leq d$) representation.

D: (Autoencoder)

tries to learn the identity function.

$$\ell(\mathbf{w}) = -\sum_{i=1}^n (y_i \log(\sigma(w^\top \mathbf{x}_i)) + (1-y_i) \log(1-\sigma(w^\top \mathbf{x}_i)))$$

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

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

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

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

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

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

This is what gave rise to the Inception module.

D. (Dimension Reduction) m channels of a $1 \times 1 \times k$ convolution $m \leq k$:

$$x^T i_j = \sigma(\mathbf{W} \mathbf{x}_{ij}), \quad \mathbf{W} \in \mathbb{R}^{m \times k}.$$

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

12.12.5 — Google Inception Network

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

• dimensionality reduction through the inception modules

• convolutions at various sizes, as different filter sizes turned out to be useful

• a max-pooling of the previous layer, and a dimensionality reduction of the result

• 1×1 conv for dimension reduction before convolving with larger ones

• then all these informations are passed to the next layer

• gradient shortcuts: connect softmax layer at intermediate stages to have the gradient flow until the beginnings of the network

• decomposition of convolution kernels for computational performance

• all-in-all the dimensionality reductions improved the efficiency

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

12.13 — Networks Similar to CNNs

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

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

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

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

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

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

#parameters FC: $m \times n \times c$

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

Ex: Assume we have an $m \times n$ image (with one channel). And we convolve it with a filter $(2p+1) \times (2q+1)$

Then the convolved image has dimensions (assuming stride 1)

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

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

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

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

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

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

• the first term $-\eta \|\nabla_\theta \mathcal{R}(\theta)\|_H^2$ will obviously be negative, as it's a negative factor of a norm

• the second term will always be positive, as the hessian matrix is positive semi-definite. Fortunately, we're squaring η , which may already be small, so the term is small. However, if the Hessian is ill-conditioned (as in the cliff-situation (curvature)). Then we can have a very large positive value in the second term. So what then can happen is that

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

14.1.1 Bi-Linear Models
The first thing that we could do is to use an information theoretic quantity: the so-called *mutual information*. The mutual information is described in information theory as *how much information one random variable has about another random variable*. If two variables are independent, then, the mutual information will be zero.

So, if we put two words nearby, it's because they have to be related somehow in the *meaning* of the sentence. Hence, we expect them to have a larger mutual information.

D. (Pointwise Mutual Information)

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

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

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

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

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

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

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

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

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

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

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

$$C_{ij} = \underbrace{\text{one-hot}(w_i)}_{\mathbf{o}_i} \underbrace{\text{one-hot}(w_j)}_{\mathbf{o}_j} = \mathbf{o}_i \mathbf{V} \mathbf{A} \mathbf{V}^T \mathbf{o}_j \\ \approx \mathbf{o}_i \underbrace{\mathbf{V}_k \mathbf{A}_k \mathbf{V}_k^T}_{k \text{ PCs}} \mathbf{o}_j = \mathbf{o}_i \mathbf{V}_k \frac{1}{k} \mathbf{A}_k^2 \mathbf{V}_k^T \mathbf{o}_j$$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

For this reason we'll introduce the following function

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

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

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

and the opposite event is given by:

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

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

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

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

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

Given an observation sequence $\mathbf{x}_1, \dots, \mathbf{x}_T$. We want to identify the hidden activities \mathbf{h}^* with the state of a dynamical system. The discrete time evolution of the *hidden state sequence* is expressed as a HMM with a non-linearity.

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

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

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

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

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

For multi-output loss

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

There are two scenarios for producing outputs

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

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

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: $\alpha = \frac{3}{4}$). Then $p_n(w) = \alpha p(w)$ to increase the chance of true negatives. Even if by chance we might produce false negative examples this event is rather rare. The factor k is just to weight the negative examples a bit more in the loss. Usually $k \approx 2$ to 10 (oversampling factor).

As we can see in the last step, instead of computing these sums separately, we'll do it as follows: we'll compute the log likelihood for one concrete positive example (w_i, w_j) and produce some negative examples (w_i, 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: $\alpha = \frac{3}{4}$). Then $p_n(w) = \alpha p(w)$ to increase the chance of true negatives. Even if by chance we might produce false negative examples this event is rather rare. The factor k is just to weight the negative examples a bit more in the loss. Usually $k \approx 2$ to 10 (oversampling factor).

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)

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

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

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

D. (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:

Backpropagation in 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

$$\text{Compute } \frac{\partial \mathcal{R}}{\partial \mathbf{y}_1}, \frac{\partial \mathcal{R}}{\partial \mathbf{y}_2}, \dots, \frac{\partial \mathcal{R}}{\partial \mathbf{y}_T} \quad (= \frac{\partial \mathcal{R}}{\partial \mathbf{y}})$$

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

$$\frac{\partial \mathcal{R}}{\partial \mathbf{h}_1} \leftarrow \sum_{i=1}^T \frac{\partial \mathcal{R}}{\partial \mathbf{y}_i} \frac{\partial \mathbf{y}_i}{\partial \mathbf{h}_1}$$

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

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

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

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

$= \sum_{t=1}^T \frac{\partial \mathcal{R}}{\partial \mathbf{y}_i} \frac{\partial \mathbf{y}_i}{\partial \mathbf{w}_{ij}}$

$\frac{\partial \mathcal{R}}{\partial \mathbf{b}_i} \leftarrow \sum_{t=1}^T \frac{\partial \mathcal{R}}{\partial \mathbf{h}_i} \frac{\partial \mathbf{h}_i}{\partial \mathbf{b}_i}$

where \mathbf{x}_w is the word embedding

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

There are three main approaches to construct *sequence embeddings*:

1

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

and the empirical co-variance matrix as

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

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

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

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

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

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

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

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

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

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

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

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

Now, using the chain rule we get:

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

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

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

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

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

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

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

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

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

– Refresher on MGFs and Gaussians

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

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

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

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

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

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

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

– 16.3 Latent Variable Models

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

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

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

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

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

– 16.3.2 Latent Variable Models

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

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

– 16.3.3 Dimensionality Reduction

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

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

where

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

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

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

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

– 16.3.4 Implicit Models

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

Deep implicit models

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

IV. Chain Rule and Jacobians for Tensors

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

D. (Tensor Multiplication)

$$\mathbf{T} = \underset{\in \mathbb{R}^{(a+b)}}{\mathbf{P}} \times_b \underset{\in \mathbb{R}^{(b+c)}}{\mathbf{Q}}$$

$$\mathbf{T} = \underset{\in \mathbb{R}^{(a \times \cdots \times c)}}{\mathbf{P}} \times \underset{\in \mathbb{R}^{(c \times \cdots \times b)}}{\mathbf{Q}}$$

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

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

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

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

T. (Tensor Chain Rule)

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

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

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

18 Generative Models

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

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

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

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

Taxonomy of Generative Models



18.1 Variational Autoencoders

Recall, that with autoencoders, we had defined a concatenation of two differentiable (non-linear) mappings $x \xrightarrow{E} z \xrightarrow{D} \hat{x}$ (an encoder E and a decoder D) and trained it with the following loss $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ (approximating identity function) in order to learn some compressed representation \mathbf{z} of \mathbf{x} which just contains the essence of \mathbf{x} according to some meaningful feature-dimensions.

Further, we could then use E to bootstrap a classification model, by first applying E and then a classification network C and fine-tuning them jointly on the cross-entropy loss.

So the lower-dimensional features \mathbf{z} capture the factors of variation in the data. And we can reconstruct an \mathbf{x} from its compressed representation \mathbf{z} .

Now the question is can we use a similar kind of setup to use new images?

VAEs define an **intractable** density function $p_{\text{model}}(\mathbf{x})$ with a latent \mathbf{z} . Having a latent variable allows us to build a network similar to an autoencoder. A tractable lower bound for the intractable density function is then derived and optimized

• **explicit density estimation:** explicitly define and solve for $p_{\text{model}}(\mathbf{x})$

• tractable we can compute $p_{\text{model}}(\mathbf{x})$

• approximate we approximate $p_{\text{model}}(\mathbf{x})$ in some way

• **implicit density estimation:** learn a model that can sample from $p_{\text{model}}(\mathbf{x})$ without explicitly defining it.

18.1.1 Variational Autoencoders (VAEs)

Relation Autoencoders

Recall, that with autoencoders, we had defined a concatenation of two differentiable (non-linear) mappings $x \xrightarrow{E} z \xrightarrow{D} \hat{x}$ (an encoder E and a decoder D) and trained it with the following loss $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ (approximating identity function) in order to learn some compressed representation \mathbf{z} of \mathbf{x} which just contains the essence of \mathbf{x} according to some meaningful feature-dimensions.

Further, we could then use E to bootstrap a classification model, by first applying E , and then a classification network C and fine-tuning them jointly on the cross-entropy loss.

So the lower-dimensional features \mathbf{z} capture the factors of variation in the data. And we can reconstruct an \mathbf{x} from its compressed representation \mathbf{z} .

Now the question is can we use a similar kind of setup to use new images?

VAEs define an **intractable** density function $p_{\text{model}}(\mathbf{x})$ with a latent \mathbf{z} . Having a latent variable allows us to build a network similar to an autoencoder. A tractable lower bound for the intractable density function is then derived and optimized

• **explicit density estimation:** explicitly define and solve for $p_{\text{model}}(\mathbf{x})$

• tractable we can compute $p_{\text{model}}(\mathbf{x})$

• approximate we approximate $p_{\text{model}}(\mathbf{x})$ in some way

• **implicit density estimation:** learn a model that can sample from $p_{\text{model}}(\mathbf{x})$ without explicitly defining it.

18.1.2 Deep Latent Gaussian Models

18.1.3 Generative Adversarial Networks (GANs)

GANs do not try to model a density function but directly aim to build a function to generated data (implicit generative method).

The whole optimization motivated by a game-theoretic approach in a