

COMPUTING GRADIENTS OF TENSOR NETWORKS

I Jacobian and backpropagation

• $f: \mathbb{R}^m \rightarrow \mathbb{R}$ Gradient $\nabla_{\theta} f = \begin{pmatrix} \partial f / \partial \theta_1 \\ \vdots \\ \partial f / \partial \theta_m \end{pmatrix}$ for each $\theta \in \mathbb{R}^m$

• $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ Jacobian $\frac{\partial f}{\partial \theta} = \left(\frac{\partial f_i(\theta)}{\partial \theta_j} \right)_{i,j} \in \mathbb{R}^{n \times m}$ for each $\theta \in \mathbb{R}^m$

• Backpropagation / Automatic Differentiation (reverse mode).

Computational Graph $\theta \rightarrow T^1 \rightarrow T^2 \rightarrow T^3 \rightarrow \mathcal{L}$ (where $X \rightarrow Y$ means that Y is a function of X)
 $\mathbb{R}^m \rightarrow \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_3} \rightarrow \mathbb{R}$

We want to compute $\frac{\partial \mathcal{L}}{\partial \theta}$.

Chain rule: $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial T_3} \frac{\partial T_3}{\partial T_2} \frac{\partial T_2}{\partial T_1} \frac{\partial T_1}{\partial \theta}$
 $\bar{\theta} \xrightarrow{\text{green arrow}} \frac{\partial \mathcal{L}}{\partial \theta}$
 $1 \times m \quad 1 \times d_3 \quad d_3 \times d_2 \quad d_2 \times d_1 \quad d_1 \times m$

We define the adjoint: $\bar{T} = \frac{\partial \mathcal{L}}{\partial T}$

$$\begin{aligned} \bar{T}_3 &= \frac{\partial \mathcal{L}}{\partial T_3} & \bar{T}_2 &= \frac{\partial \mathcal{L}}{\partial T_2} = \frac{\partial \mathcal{L}}{\partial T_3} \frac{\partial T_3}{\partial T_2} = \bar{T}_3 \frac{\partial T_3}{\partial T_2} \\ \bar{T}_1 &= \frac{\partial \mathcal{L}}{\partial T_1} = \frac{\partial \mathcal{L}}{\partial T_2} \frac{\partial T_2}{\partial T_1} = \bar{T}_2 \frac{\partial T_2}{\partial T_1} \\ \bar{\theta} &= \dots = \bar{T}_1 \frac{\partial T_1}{\partial \theta} \end{aligned} \quad \left\{ \begin{array}{l} \bar{T}^j = \bar{T}^{j+1} \frac{\partial T^{j+1}}{\partial T^j} \end{array} \right.$$

More complex graphs: $\emptyset \rightarrow T_1 \begin{matrix} \nearrow T_2 \\ \searrow T_3 \end{matrix} \rightarrow T_4 \rightarrow \mathcal{L}$

$$\overline{T_i} = \sum_{j: \text{child of } i} \overline{T_j} \frac{\partial T_j}{\partial T_i}$$

ex: $\overline{T_1} = \frac{\partial \mathcal{L}}{\partial T_1} = \frac{\partial \mathcal{L}}{\partial T_2} \frac{\partial T_2}{\partial T_1} + \frac{\partial \mathcal{L}}{\partial T_3} \frac{\partial T_3}{\partial T_1}$
 $= \overline{T_2} \frac{\partial T_2}{\partial T_1} + \overline{T_3} \frac{\partial T_3}{\partial T_1}$

$f: \mathbb{R}^{m_1 \times m_2 \times \dots \times m_N} \rightarrow \mathbb{R}^{n_1 \times n_2 \times \dots \times n_P} \quad (\cong \mathbb{R}^{m_1 m_2 \dots m_N} \rightarrow \mathbb{R}^{n_1 n_2 \dots n_P})$

Jacobian Tensor $\frac{\partial f(T)}{\partial T} = \left(\frac{\partial f(T)_{i_1 \dots i_P}}{\partial T_{j_1 \dots j_N}} \right)_{i_1 \dots i_P j_1 \dots j_N} \in \mathbb{R}^{n_1 \times \dots \times n_P \times m_1 \times \dots \times m_N}$

II Jacobian of Tensor Networks

Let W be a tensor given as a tensor network where G is a core tensor appearing only once. Then $\frac{\partial W}{\partial G}$ is simply obtained by deleting G in the tensor network.

Ex: $W = \begin{matrix} d_1 & & d_3 \\ & \text{---} A \text{---} B \text{---} D & \\ & | & | & \\ d_2 & C & R_2 & d_4 \end{matrix}$

$\frac{\partial W}{\partial B} = \begin{matrix} d_1 & & d_3 \\ & \text{---} A \text{---} & R_3 & \\ & | & | & \\ d_2 & C & R_2 & d_4 \end{matrix}$

$\frac{\partial W}{\partial D} = \begin{matrix} d_1 & & d_3 \\ & \text{---} A \text{---} B \text{---} & \\ & | & | & \\ d_2 & C & R_2 & d_4 \end{matrix}$

Jacobian of $B \mapsto W$
 $R_1 \times R_2 \times R_3 \leftarrow d_1 \times d_2 \times d_3 \times d_4$

• Rederiving classical matrix/vector gradients:

$\frac{\partial \langle u, v \rangle}{\partial u} = \partial(u \cdot v) / \partial u = -v = v$

$\frac{\partial Ax}{\partial x} = \partial(-A \cdot x) / \partial x = -A = A$

$\frac{\partial x^T A x}{\partial A} = \partial(x \cdot A \cdot x) / \partial A = x \cdot \quad \cdot x = x x^T$

$\frac{\partial \text{Tr}(A)}{\partial A} = \partial(\text{Tr}(A)) / \partial A = \text{Tr} = I$

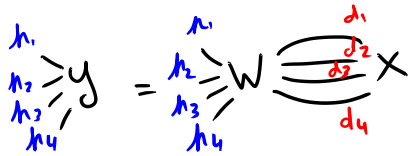
$\frac{\partial Ax}{\partial A} = \partial(-A \cdot x) / \partial A = -x \cong I \circ x$

If the core tensor G appears k times in the tensor network of W , then $\frac{\partial W}{\partial G}$ is given the sum of k copies of W where a different occurrence of G is deleted in each copy.

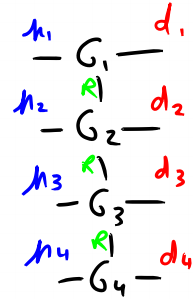
Ex: $\frac{\partial x^T A x}{\partial x} = \frac{\partial (x - A_2 x)}{\partial x} = -A_2 x + x - A_2 x = (A + A^T) x$

Applications

TENSORIZING
NEURAL
NETWORKS



where $\begin{pmatrix} W \end{pmatrix} =$



↳ We want to compute $\frac{\partial \mathcal{L}}{\partial G_2} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial G_2}$ ← Some loss depending on y



We need to compute $\frac{\partial y}{\partial G_2} \in \mathbb{R}^{h_1 \times h_2 \times h_3 \times h_4 \times \mathbb{R} \times h_2 \times d_2 \times \mathbb{R}}$

We have:

$$\frac{\partial y}{\partial G_2} = \frac{\partial \left(\begin{pmatrix} -G_1 \\ -G_2 \\ -G_3 \\ -G_4 \end{pmatrix} \right) x}{\partial G_2} = \begin{pmatrix} h_1 - G_1 \\ h_2 - G_2 \\ h_3 - G_3 \\ h_4 - G_4 \end{pmatrix} x$$

IV Generalized CP decomposition

GENERALIZED CANONICAL POLYADIC TENSOR DECOMPOSITION*

DAVID HONG¹, TAMARA G. KOLDA², AND JED A. DUERSCH³

↳ Likelihood maximization & generalized linear model.

Observations/Data : $X \in \mathbb{R}^{m_1 \times \dots \times m_d}$

Parameterized PDF: $X_{i_1, \dots, i_d} \sim p(X_{i_1, \dots, i_d} | \theta_{i_1, \dots, i_d})$
where $p(\theta_{i_1, \dots, i_d}) = M_{i_1, \dots, i_d}$

$\ell: \mathbb{R} \rightarrow \mathbb{R}$ is an invertible link function connecting the model parameters M_{i_1, \dots, i_d} to the "natural" parameter of the distribution.

Model parameters : $M \in \mathbb{R}^{m_1 \times \dots \times m_d}$

Goal: Find M that maximizes the likelihood of the observed entries of X .
↳ $\Omega \subseteq [m_1] \times \dots \times [m_d]$

$$\max_M \mathcal{L}(M; X, \Omega) \equiv \prod_{(i_1, \dots, i_d) \in \Omega} p(X_{i_1, \dots, i_d} | \theta_{i_1, \dots, i_d})$$

$$\text{where } \ell(\theta_{i_1, \dots, i_d}) = M_{i_1, \dots, i_d}$$

→ We only have (at most) one observation for each entry of X !

We will assume a low rank structure of M (interdependence between the θ_{i_1, \dots, i_d})
↳ low CP rank of M .

Reformulation of optimization problem:

★
$$\min_M F(M; X, \Omega) = \sum_{(i_1, \dots, i_d) \in \Omega} f(X_{i_1, \dots, i_d}, M_{i_1, \dots, i_d}) \text{ s.t. } \text{rank}_{CP}(M) \leq R$$

where $f(x, m) = -\log p(x | \ell^{-1}(m))$

↳ "loss function" (assumed differentiable)

Examples of data distribution

① Gaussian distribution

(Gaussian MLE \Leftrightarrow Mean squared error)

$$X_{i, \dots, d} = M_{i, \dots, d} + \varepsilon_{i, \dots, d} \text{ where } \varepsilon_{i, \dots, d} \sim \mathcal{N}(0, \sigma^2)$$

↑ constant

Equivalently: $X_{i, \dots, d} \sim \mathcal{N}(\theta_{i, \dots, d}, \sigma^2)$ where $\theta_{i, \dots, d} = M_{i, \dots, d}$

the link function is the identity $\ell(\theta) = \theta$

Simple derivation:

$$f(x, m) = -\log \mathcal{P}(x; m, \sigma^2) = \frac{(x - m)^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)$$

↳ ignoring σ^2 and constants $f(x, m) = (x - m)^2$

\Rightarrow ★ is then equivalent to CP decomposition.

② Bernoulli distribution.

If $X_{i, \dots, d} \in \{0, 1\}$ (binary data), we can assume $X_{i, \dots, d} \sim \mathcal{B}(\theta_{i, \dots, d})$

$$p(x|\theta) = \theta^x (1-\theta)^{1-x}, \quad x \in \{0, 1\}, \theta \in [0, 1]$$

link function?

- identity \leadsto we need to constrain $m \in [0, 1]$
- odds ratio: $m = \ell(\theta) = \frac{\theta}{1-\theta}$

\leadsto only need to enforce $m \geq 0$

- log odds ratio:

$$m = \ell(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$$

$$\hookrightarrow \theta = \frac{e^m}{1+e^m}, \quad 1-\theta = \frac{1}{1+e^m}$$

$$\dots f(x, m) = -\log p(x|\theta) = \underbrace{-x \log \theta - (1-x) \log (1-\theta)}_{\text{cross-entropy}}$$

$$f(x, m) = \log(1+e^m) - xm$$

③ Other distributions

Table 1: Statistically-motivated loss functions. Parameters in blue are assumed to be constant. Numerical adjustments are indicated in red.

Distribution	Link function	Loss function	Constraints
$\mathcal{N}(\mu, \sigma)$	$m = \mu$	$(x - m)^2$	$x, m \in \mathbb{R}$
Gamma(k, σ)	$m = k\sigma$	$x/(m+\epsilon) + \log(m+\epsilon)$	$x > 0, m \geq 0$
Rayleigh(θ)	$m = \sqrt{\pi/2}\theta$	$2\log(m+\epsilon) + (\pi/4)(x/(m+\epsilon))^2$	$x > 0, m \geq 0$
Poisson(λ)	$m = \lambda$	$m - x \log(m+\epsilon)$	$x \in \mathbb{N}, m \geq 0$
	$m = \log \lambda$	$e^m - xm$	$x \in \mathbb{N}, m \in \mathbb{R}$
Bernoulli(ρ)	$m = \rho / (1 - \rho)$	$\log(m+1) - x \log(m+\epsilon)$	$x \in \{0, 1\}, m \geq 0$
	$m = \log(\rho / (1 - \rho))$	$\log(1 + e^m) - xm$	$x \in \{0, 1\}, m \in \mathbb{R}$
NegBinom(r, ρ)	$m = \rho / (1 - \rho)$	$(r+x) \log(1+m) - x \log(m+\epsilon)$	$x \in \mathbb{N}, m \geq 0$

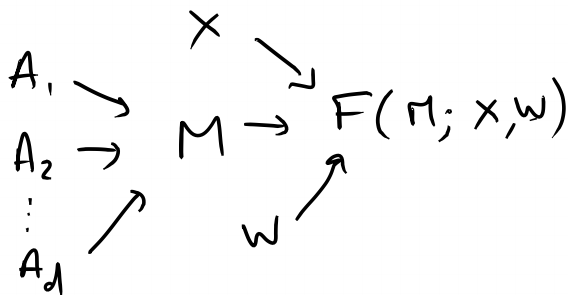
Gradient based optimization

(if Ω is the set of observed entries;
 $W_{i, \dots, id} = 1$ if $(i, \dots, id) \in \Omega$ and 0 o.w.)

★ : $\min_{A_1, \dots, A_d} F(M; X, W) \equiv \sum_{i, \dots, id} W_{i, \dots, id} f(x_{i, \dots, id}, M_{i, \dots, id})$

s.t. $M = \llbracket A_1, A_2, \dots, A_d \rrbracket$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ m_1 \times R & m_2 \times R & m_d \times R \end{matrix}$



We want to show

$$\frac{\partial F(M; X, W)}{\partial A_k} = -y_{m_k} \begin{pmatrix} A_1 \\ \vdots \\ A_{k-1} \\ A_{k+1} \\ \vdots \\ A_d \end{pmatrix}^T I_R \quad \text{where } y \in \mathbb{R}^{m_1 \times \dots \times m_d}$$

$\begin{matrix} 1 \times m_k \times R \end{matrix}$

$$= y_{(k)} (A_d \otimes \dots \otimes A_{k+1} \otimes A_{k-1} \otimes \dots \otimes A_1)$$

$y_{i, \dots, id} = W_{i, \dots, id} \frac{\partial f(x_{i, \dots, id}, M_{i, \dots, id})}{\partial M_{i, \dots, id}}$
 $\begin{matrix} 1 \times 1 \end{matrix}$
 ↳ of the same size as X but it is sparse if W is sparse.

$$\frac{\partial F(M; X, W)}{\partial A_k} = \overbrace{\frac{\partial F(M; X, \Pi)}{\partial M}}^y \frac{\partial M}{\partial A_k}$$

$1 \times m_k \times R$

$1 \times m_1 \times \dots \times m_d \quad m_1 \times \dots \times m_d \times m_k \times R$

$$\begin{aligned} \frac{\partial F(M; X, \Pi)}{\partial M}_{i_1, \dots, i_d} &= \frac{\partial \left(\sum_{j_1, \dots, j_d} W_{j_1, \dots, j_d} f(X_{j_1, \dots, j_d}, M_{j_1, \dots, j_d}) \right)}{\partial M_{i_1, \dots, i_d}} \\ &= W_{i_1, \dots, i_d} \frac{\partial f(X_{i_1, \dots, i_d}, M_{i_1, \dots, i_d})}{\partial M_{i_1, \dots, i_d}} := y_{i_1, \dots, i_d} \end{aligned}$$

$$\frac{\partial M}{\partial A_k} = \frac{\partial \left(\begin{matrix} -A_1 \\ -A_2 \\ \vdots \\ -A_d \end{matrix} \right) \cdot I}{\partial A_k} = \begin{matrix} m_1 & - & A_1 \\ & \vdots & \\ m_{k-1} & - & A_{k-1} \\ m_k & - & m_k \quad R \\ & \vdots & \\ m_{k+1} & - & A_{k+1} \\ & \vdots & \\ m_d & - & A_d \end{matrix} \cdot I$$

$$\begin{aligned} \frac{\partial F(M; X, W)}{\partial A_k} &= \frac{\partial F(M; X, \Pi)}{\partial M} \frac{\partial M}{\partial A_k} \\ &= \begin{matrix} & & A_1 \\ & & \vdots \\ & & A_{k-1} \\ m_k & - & y \\ & & A_{k+1} \\ & & \vdots \\ & & A_d \end{matrix} \cdot I \cdot R = y_{(k)} (A_d \otimes \dots \otimes A_{k+1} \otimes A_{k-1} \otimes \dots \otimes A_1) \end{aligned}$$

Follow up paper:

Stochastic Gradients for Large-Scale Tensor Decomposition*

Tamara G. Kolda[†] and David Hong[‡]