بسم الله الرحمن الرحيم

انهارده هنقفل كلام علي الباك بروب وهنخش في الديب ليرننج .. لما بنتكلم علي الديب ليرننج هنتكلم علي حاجتين CNNs and RNNs واللي محتاجين نعرفو هو ازاي ال 2 دول بيعملولهم تمرين .. فعلياً هو الباك بروب يعني .. المهم الدكتور مسميهم هم الاتنين decision functions ... المهم HW4 هيخلص انهارده ان شاء الله .. ربنا يعين ويوفق ... الدكتور مش بيطلع الامتحانات برا .. فبينزلو paper versions بس ... في سؤال الدكتور اتسألو ...

# Q&A

| Q: | Do I need to know Matrix Calculus to derive the backprop algorithms used in this class? |
|---|---|
| A: | No. We've carefully constructed our assignments so that you do **not** need to know Matrix Calculus. |

That said, it's kind of handy.

الدكتور بيقلك ان ده يعني بيجي في الإيدين كسكيلل يعني . بس اشطا ..



## Matrix Calculus

Let $y, x \in \mathbb{R}$ be scalars, $\mathbf{y} \in \mathbb{R}^M$ and $\mathbf{x} \in \mathbb{R}^P$ be vectors, and $\mathbf{Y} \in \mathbb{R}^{M \times N}$ and $\mathbf{X} \in \mathbb{R}^{P \times Q}$ be matrices

| Types of Derivatives | scalar |
|---|---|
| scalar | $\frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x}\right]$  1x1 matrix |
| vector | $\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_P} \end{bmatrix}$  List of deriv |
| matrix | $\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial X_{11}} & \frac{\partial y}{\partial X_{12}} & \cdots & \frac{\partial y}{\partial X_{1Q}} \\ \frac{\partial y}{\partial X_{21}} & \frac{\partial y}{\partial X_{22}} & \cdots & \frac{\partial y}{\partial X_{2Q}} \\ \vdots & & & \vdots \\ \frac{\partial y}{\partial X_{P1}} & \frac{\partial y}{\partial X_{P2}} & \cdots & \frac{\partial y}{\partial X_{PQ}} \end{bmatrix}$ |

# Matrix Calculus

| Types of Derivatives | scalar | vector |
|---|---|---|
| scalar | $\dfrac{\partial y}{\partial x} = \left[\dfrac{\partial y}{\partial x}\right]$ | $\dfrac{\partial \mathbf{y}}{\partial x} = \left[\begin{matrix}\dfrac{\partial y_1}{\partial x} & \dfrac{\partial y_2}{\partial x} & \cdots & \dfrac{\partial y_N}{\partial x}\end{matrix}\right]$ |
| vector | $\dfrac{\partial y}{\partial \mathbf{x}} = \left[\begin{matrix}\dfrac{\partial y}{\partial x_1}\\ \dfrac{\partial y}{\partial x_2}\\ \vdots \\ \dfrac{\partial y}{\partial x_P}\end{matrix}\right]$ | $\dfrac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left[\begin{matrix}\dfrac{\partial y_1}{\partial x_1} & \dfrac{\partial y_2}{\partial x_1} & \cdots & \dfrac{\partial y_N}{\partial x_1}\\ \dfrac{\partial y_1}{\partial x_2} & \dfrac{\partial y_2}{\partial x_2} & \cdots & \dfrac{\partial y_N}{\partial x_2}\\ \vdots & & & \\ \dfrac{\partial y_1}{\partial x_P} & \dfrac{\partial y_2}{\partial x_P} & \cdots & \dfrac{\partial y_N}{\partial x_P}\end{matrix}\right]$ |

*(handwritten annotations)* with respect to scalar — N entries — P entries

# Matrix Calculus

*Common Vector Derivatives*

Let $\dfrac{\partial f(\vec{x})}{\partial \vec{x}} = \nabla_x f(\vec{x})$ be the vector derivative of $f$, $B \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^m$

**Scalar Derivative**

$f(x) \rightarrow \dfrac{\partial f}{\partial x}$

$bx \rightarrow b$

$xb \rightarrow b$

$x^2 \rightarrow 2x$

$bx^2 \rightarrow 2bx$

**Vector Derivative**

$f(x) \rightarrow \dfrac{\partial f}{\partial x}$

$x^T b \rightarrow b$

$x^T x \rightarrow 2x$ ← vector x

$x^T B x \rightarrow 2Bx$ ← B symmetric

*Vector Chain Rule*

$\vec{x} \in \mathbb{R}^P \quad y \in \mathbb{R} \quad u \in \mathbb{R}^N$

$$\dfrac{\partial y}{\partial \vec{x}} = \left(\left(\dfrac{\partial y}{\partial \vec{u}}\right)^T \left(\dfrac{\partial \vec{u}}{\partial \vec{x}}\right)^T\right)^T$$

$P \times 1 \quad N \times 1 \quad P \times N$

$1 \times N \quad N \times P$

$P \times 1$

$\left[\dfrac{\partial y}{\partial \vec{x}}\right]_P \quad \left[\dfrac{\partial y}{\partial \vec{u}}\right]_n \quad \left[\dfrac{\partial \vec{u}}{\partial \vec{x}}\right]_{pn}$

طيب هنا احنا بنبص علي الصوره بتاعت الباك بروب ألجورذم في الجينيرال فيرجن .. عندك كمان الفورود بروب .. بتقول ان عندك فانكشن f .. اللي انت عاوز تحسباه فانت عاوز تكتب الجورذم يحسبها .. الألجورذم هيعرف directed acyclic graph .. ده اللي هو ال computational graph .. ال forward computational هيقلك visit each node in topological order اللي معناها visit all the children of a node before visiting the node itself .. فلفاريبل اسمو ui .. عندو انبوتس اسمها v1 ... vN كإنها هتحسب ال ui كإنها intermediate function و تخزن النواتج في النوود اللي عندك . عشان هتحتاجها في الباك بروب ... في الباك بروب انت بت initialize dy/dy ل 1 .. و بعدين تبدأ ت visit each node in reverse topological order ...

الدكتور رسم الصوره وقال ان كل ال Us فانكشن في كل ال Vs ... انت هنا في الباك بروب دايماً هت visit the parent before the children ...



في فيرجن مختلف في الباك بروب .. التعديل بس هييقا في انك هتعامل ال intermediate derivatives كإنهم intermediate variables موجودين عندك من البدايه ... احنا هن ن initialize them to 0 .. بس الطريقه اللي هنحسب بيها هتبقا مختلفه ... ال computational graph لسه متغيرش ... لما هتخش علي ال ui انت عارف كل ال children بتوعو اللي هو كل ال Vs اللي عندك ..

**Automatic Differentiation – Reverse Mode (aka. Backpropagation)**

<u>Forward Computation</u>
1. Write an **algorithm** for evaluating the function y = f(**x**). The algorithm defines a **directed acyclic graph**, where each variable is a node (i.e. the "**computation graph**")
2. Visit each node in **topological order**.
   For variable $u_i$ with inputs $v_1,\ldots, v_N$
   a. Compute $u_i = g_i(v_1,\ldots, v_N)$
   b. Store the result at the node

<u>Backward Computation (Version B)</u>
1. **Initialize** all partial derivatives $dy/du_j$ to 0 and $dy/dy = 1$.
2. Visit each node in **reverse topological order.**
   For variable $u_i = g_i(v_1,\ldots, v_N)$
   a. We already know $dy/du_i$
   b. Increment $dy/dv_j$ by $(dy/du_i)(du_i/dv_j)$
      (Choice of algorithm ensures computing $(du_i/dv_j)$ is easy)

**Return** partial derivatives $dy/du_i$ for all variables

14

---

*Why is the backpropagation algorithm efficient?*
1. Reuses **computation from the forward pass** in the backward pass
2. Reuses **partial derivatives** throughout the backward pass (*but only if the algorithm reuses shared computation in the forward pass*)

   (Key idea: partial derivatives in the backward pass should be thought of as variables stored for reuse)

*Example: 1-Hidden Layer Neural Network*

**Algorithm 1** Stochastic Gradient Descent (SGD)

1: **procedure** SGD(Training data $\mathcal{D}$, test data $\mathcal{D}_t$)
2:     Initialize parameters $\alpha, \beta$   *randomly, zero*
3:     **for** $e \in \{1, 2, \ldots, E\}$ **do**   ← *epochs*
4:       **for** $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$ **do**
5:        Compute neural network layers:
6:        $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{b}, \mathbf{z}, \hat{\mathbf{y}}, J) = \text{NNFORWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta)$
7:        Compute gradients via backprop:
8:        $\left.\begin{array}{l} \mathbf{g}_\alpha = \nabla_\alpha J \\ \mathbf{g}_\beta = \nabla_\beta J \end{array}\right\} = \text{NNBACKWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta, \mathbf{o})$
9:        Update parameters:
10:        $\alpha \leftarrow \alpha - \gamma \mathbf{g}_\alpha$
11:        $\beta \leftarrow \beta - \gamma \mathbf{g}_\beta$
12:       Evaluate training mean cross-entropy $J_{\mathcal{D}}(\alpha, \beta)$
13:       Evaluate test mean cross-entropy $J_{\mathcal{D}_t}(\alpha, \beta)$
14:     **return** parameters $\alpha, \beta$

16

المهم الدكتور حط مثال كدا وقال ان احنا فعلا ممكن نستخدم ال finite difference method عشان نتشيك الحل بتاعنا في الباكبروب .. انا بحب الدكتور ده جداً ..... يلا نخش في الديب ليرننج

اول حاجه هنتكلم علي ال CNN بس الأول هنسال ليه الاهتمام ده ... عشان فيه فلوس كتيره الناس بتستثمر فيها .. بقيت المحاضره جري كتير .. لو مهتم اقراها تاني