

A Deep Dive into $\frac{d}{dn}$ $\begin{bmatrix} \text{deep} \\ \text{declarative} \\ \text{networks} \end{bmatrix}$

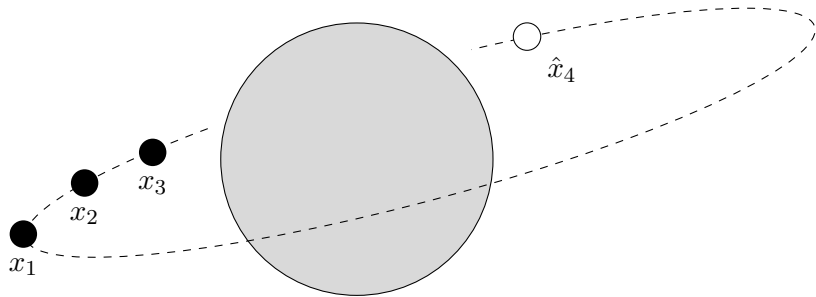
Stephen Gould

`stephen.gould@anu.edu.au`

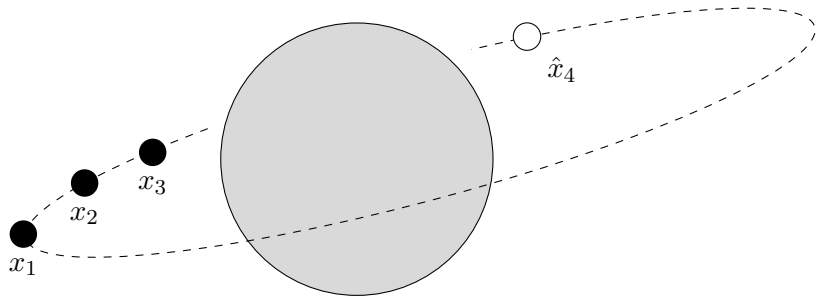
Robotic Vision Summer School (RVSS), 2024
Australian National University

9 February 2024

Discovery of Ceres



Discovery of Ceres



Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints
- ▶ **electrical engineering:** minimise the size of a transistor in a circuit subject to power and timing constraints

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints
- ▶ **electrical engineering:** minimise the size of a transistor in a circuit subject to power and timing constraints
- ▶ **logistics and planning:** find the cheapest way to distribute goods from suppliers to consumers across a transportation network

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints
- ▶ **electrical engineering:** minimise the size of a transistor in a circuit subject to power and timing constraints
- ▶ **logistics and planning:** find the cheapest way to distribute goods from suppliers to consumers across a transportation network
- ▶ **statistics/data science:** curve fitting and data visualisation

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints
- ▶ **electrical engineering:** minimise the size of a transistor in a circuit subject to power and timing constraints
- ▶ **logistics and planning:** find the cheapest way to distribute goods from suppliers to consumers across a transportation network
- ▶ **statistics/data science:** curve fitting and data visualisation
- ▶ **robotics:** optimise control parameters to achieve some goal state or trajectory

Optimisation is Everywhere

- ▶ **financial mathematics:** maximise profits or minimise costs subject to constraints on resources and budgets
- ▶ **mechanical engineering:** maximise the span of a bridge subject to load constraints; minimise drag on a vehicle subject to volume constraints
- ▶ **electrical engineering:** minimise the size of a transistor in a circuit subject to power and timing constraints
- ▶ **logistics and planning:** find the cheapest way to distribute goods from suppliers to consumers across a transportation network
- ▶ **statistics/data science:** curve fitting and data visualisation
- ▶ **robotics:** optimise control parameters to achieve some goal state or trajectory
- ▶ **machine learning and deep learning:** minimise loss functions with respect to the parameters of our model

Optimisation Problems

*find an assignment to variables that minimises
a measure of cost subject to some constraints¹*

¹In these lectures we will be concerned with continuous-valued variables

Optimisation Problems

minimize (over x) $\text{objective}(x)$
subject to $\text{constraints}(x)$

Optimisation Problems

$$\begin{array}{ll}\text{minimize} & f_0(x) \\ \text{subject to} & f_i(x) \leq 0, \quad i = 1, \dots, p \\ & h_i(x) = 0, \quad i = 1, \dots, q\end{array}$$

- ▶ $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ — optimisation variables
- ▶ $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ — objective (or cost or loss) function
- ▶ $f_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 1, \dots, p$ — inequality constraint functions
- ▶ $h_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 1, \dots, q$ — equality constraint functions

Least Squares

$$\text{minimize} \quad \|Ax - b\|_2^2$$

Least Squares

$$\text{minimize } \|Ax - b\|_2^2$$

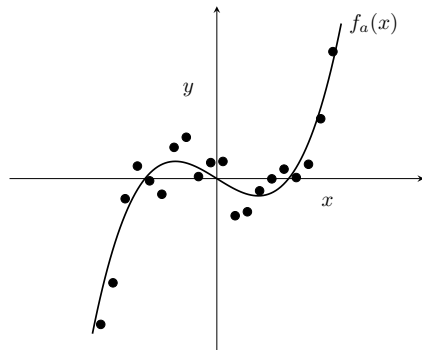
- ▶ unique solution if $A^T A$ is invertible, $x^\star = (A^T A)^{-1} A^T b$
- ▶ solution via SVD, $A = U \Sigma V^T$, if $A^T A$ not invertible, $x^\star = V \Sigma^{-1} U^T b$
 - ▶ in fact, $x^\star + w$ for any $w \in \mathcal{N}(A)$ also a solution
- ▶ solution via QR factorisation, $x^\star = R^{-1} Q^T b$
- ▶ solved in $O(n^2 m)$ time, less if structured
- ▶ typically use iterative solver (for large scale problems)

Example: Polynomial Curve Fitting

fit n -th order polynomial $f_a(x) = \sum_{k=0}^n a_k x^k$ to set of noisy points $\{(x_i, y_i)\}_{i=1}^m$
(here a are the variables, and x and y are the data)

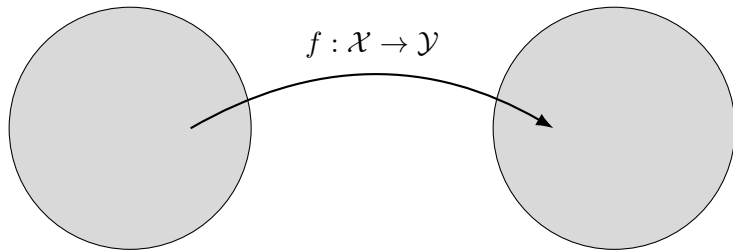
$$\text{minimize (over } a) \quad \sum_{i=1}^m (f_a(x_i) - y_i)^2$$

$$\text{minimize} \quad \left\| \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \right\|_2^2$$

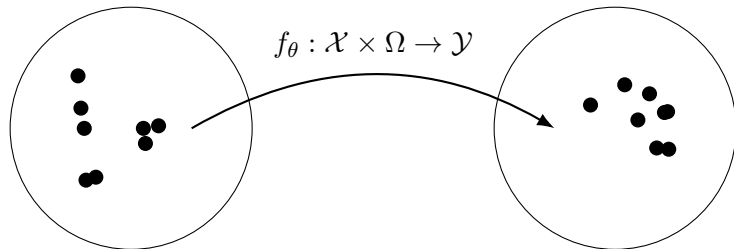


Part I. Machine Learning and Deep Learning

Machine Learning from 10,000ft



Machine Learning from 10,000ft

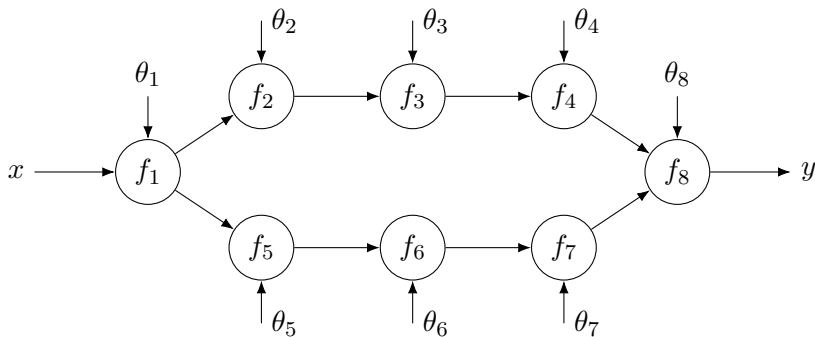


$$\text{minimize (over } \theta) \quad \sum_{(x,y) \sim \mathcal{X} \times \mathcal{Y}} L(f_\theta(x), y)$$

- ▶ loss L — what to do
- ▶ model f_θ — how to do it
- ▶ optimised by gradient descent (or variant thereof)

Deep Learning as an End-to-end Computation Graph

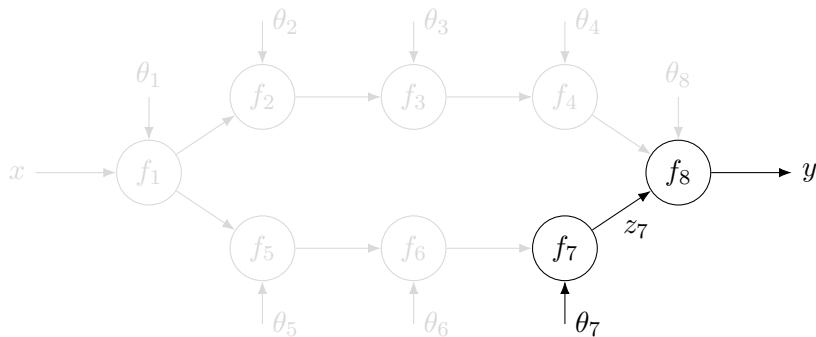
Deep learning does this by constructing the model f_θ (equiv. computation graph) as the composition of many simple parametrized functions (equiv. computation nodes).



$$y = f_8(f_4(f_3(f_2(f_1(x)))), f_7(f_6(f_5(f_1(x)))))$$

(parameters θ_i omitted for brevity)

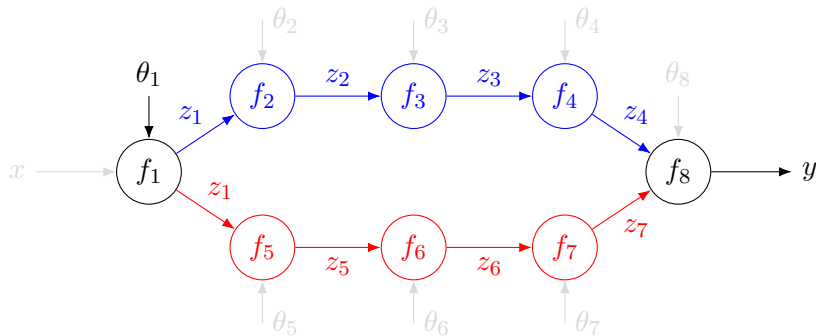
Backward Pass Gradient Calculation



Example 1.

$$\frac{\partial L}{\partial \theta_7} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_7} \frac{\partial z_7}{\partial \theta_7}$$

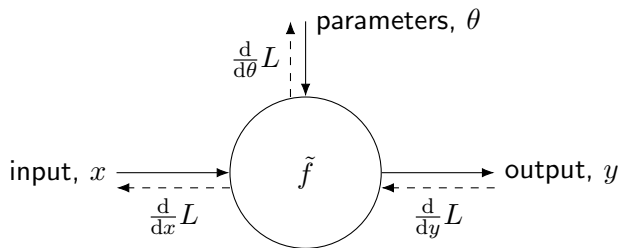
Backward Pass Gradient Calculation



Example 2.

$$\frac{\partial L}{\partial \theta_1} = \frac{\partial L}{\partial y} \left(\frac{\partial y}{\partial z_4} \frac{\partial z_4}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} + \frac{\partial y}{\partial z_7} \frac{\partial z_7}{\partial z_6} \frac{\partial z_6}{\partial z_5} \frac{\partial z_5}{\partial z_4} \right) \frac{\partial z_1}{\partial \theta_1}$$

Deep Learning Node



- **Forward pass:** compute output y as a function of the input x (and model parameters θ).

- **Backward pass:** compute the derivative of the loss with respect to the input x (and model parameters θ) given the derivative of the loss with respect to the output y .

Aside: Notation (Often Sloppy)

For scalar-valued functions:

total derivative: $\frac{df}{dx}$

partial derivative: $\frac{\partial f}{\partial x}$

Aside: Notation (Often Sloppy)

For scalar-valued functions:

total derivative: $\frac{df}{dx}$

partial derivative: $\frac{\partial f}{\partial x}$

For multi-dimensional vector-valued functions, $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

$$\frac{d}{dx} f(x) = \begin{bmatrix} \frac{df_1}{dx_1} & \cdots & \frac{df_1}{dx_n} \\ \vdots & \ddots & \vdots \\ \frac{df_m}{dx_1} & \cdots & \frac{df_m}{dx_n} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad \left(\frac{\partial}{\partial x} f(x, y) \text{ for partial} \right)$$

Sometimes D and D_X for $\frac{d}{dx}$ and $\frac{\partial}{\partial x}$, respectively.

Aside: Notation (Often Sloppy)

For scalar-valued functions:

total derivative: $\frac{df}{dx}$

partial derivative: $\frac{\partial f}{\partial x}$

For multi-dimensional vector-valued functions, $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

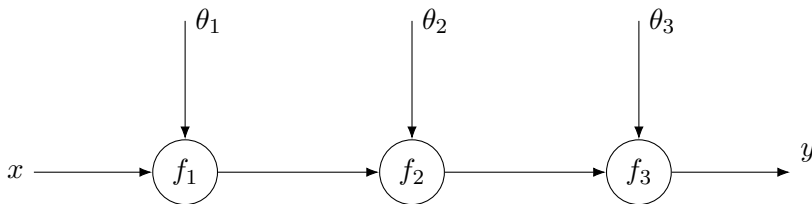
$$\frac{d}{dx} f(x) = \begin{bmatrix} \frac{df_1}{dx_1} & \cdots & \frac{df_1}{dx_n} \\ \vdots & \ddots & \vdots \\ \frac{df_m}{dx_1} & \cdots & \frac{df_m}{dx_n} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad \left(\frac{\partial}{\partial x} f(x, y) \text{ for partial} \right)$$

Sometimes D and D_X for $\frac{d}{dx}$ and $\frac{\partial}{\partial x}$, respectively.

Mathematically, derivatives with respect to (scalar-valued) loss functions are row vectors ($m = 1$).

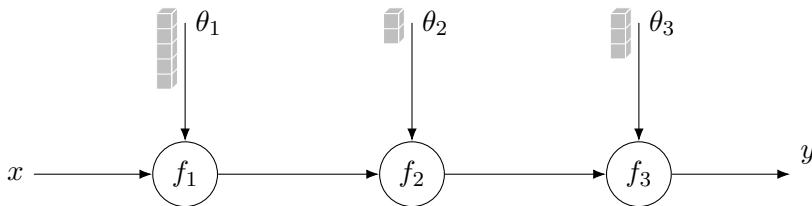
Concerning Memory

- data is often processed in batches ($B \times N \times \dots \times C$)



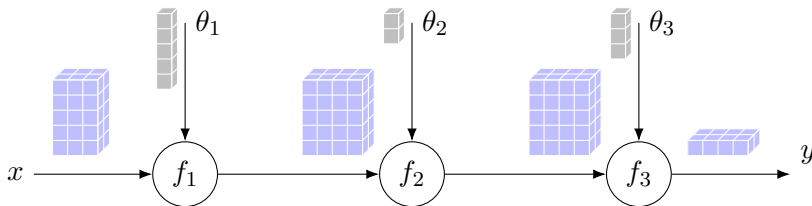
Concerning Memory

- ▶ data is often processed in batches ($B \times N \times \dots \times C$)



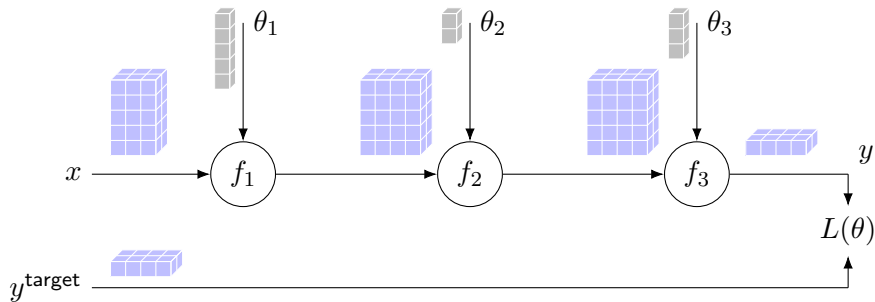
Concerning Memory

- ▶ data is often processed in batches ($B \times N \times \dots \times C$)



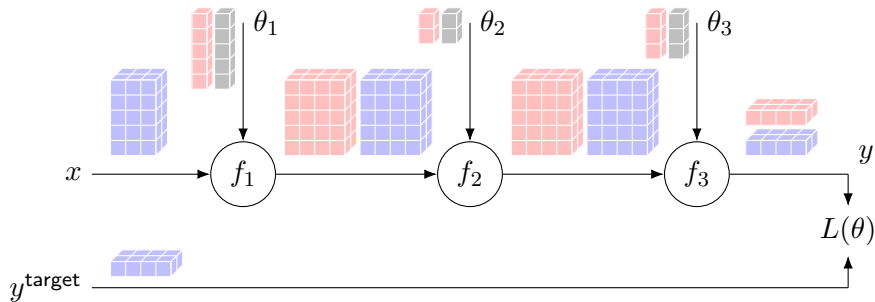
Concerning Memory

- ▶ data is often processed in batches ($B \times N \times \dots \times C$)



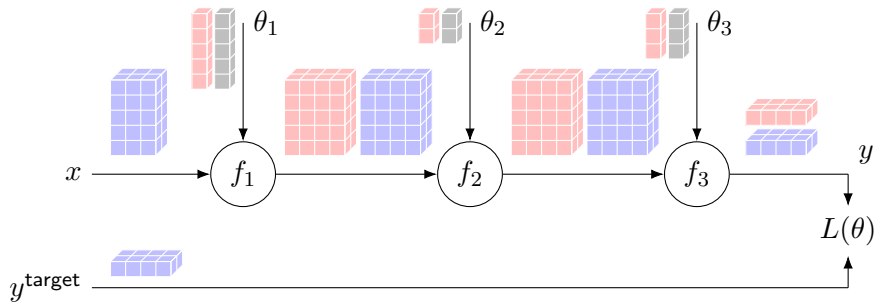
Concerning Memory

- ▶ data is often processed in batches ($B \times N \times \dots \times C$)



Concerning Memory

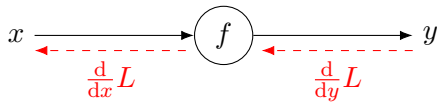
- ▶ data is often processed in batches ($B \times N \times \dots \times C$)



- ▶ parameters (usually) only take a small amount of memory (relative to data)
- ▶ derivatives take the same amount of space as the data **and stored transposed!**
- ▶ in-place operations may save memory in the forward pass
- ▶ re-using buffers may save memory in the backward pass
- ▶ at test time intermediate results are not stored

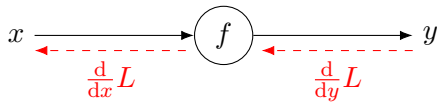
Quick Quiz

Quick Quiz



$$y = Ax$$

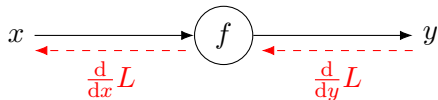
Quick Quiz



$$y = Ax$$

$$\begin{aligned}\frac{dL}{dx} &= \frac{dL}{dy} \frac{dy}{dx} \\ &= \frac{dL}{dy} A\end{aligned}$$

Quick Quiz

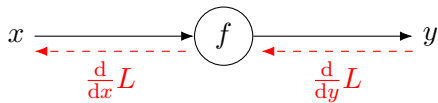


$$y = Ax$$

$$\begin{aligned}\frac{dL}{dx} &= \frac{dL}{dy} \frac{dy}{dx} \\ &= \frac{dL}{dy} A\end{aligned}$$

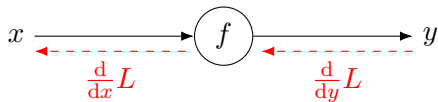
- ▶ forward pass $O(n^2)$, less if A is structured
- ▶ backward pass costs same as forward pass

Quick Quiz (2)



$$Ay = x$$

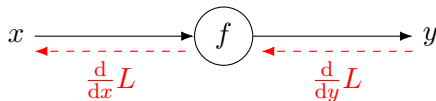
Quick Quiz (2)



$$Ay = x$$
$$\therefore y = A^{-1}x$$

$$\begin{aligned}\frac{dL}{dx} &= \frac{dL}{dy} \frac{dy}{dx} \\ &= \frac{dL}{dy} A^{-1}\end{aligned}$$

Quick Quiz (2)



$$Ay = x$$
$$\therefore y = A^{-1}x$$

$$\begin{aligned}\frac{dL}{dx} &= \frac{dL}{dy} \frac{dy}{dx} \\ &= \frac{dL}{dy} A^{-1}\end{aligned}$$

- ▶ forward pass $O(n^3)$, less if structured
- ▶ backward pass solves $w = A^T v$
 - ▶ **cheaper** than forward pass if decomposition of A is cached

Automatic Differentiation (AD)

- ▶ algorithmic procedure that produces code for computing **exact** derivatives
- ▶ assumes numeric computations are composed of a small set of elementary operations that we know how to differentiate
 - ▶ arithmetic, exp, log, trigonometric
- ▶ workhorse of modern machine learning that greatly reduces development effort
- ▶ roughly speaking, for each line of the forward pass code, `P, Q = foo(A, B, C)`, autodiff produces a line `dLdA, dLdB, dLdC = foo_vjp(dLdP, dLdQ)` in the backward pass code

Automatic Differentiation (AD)

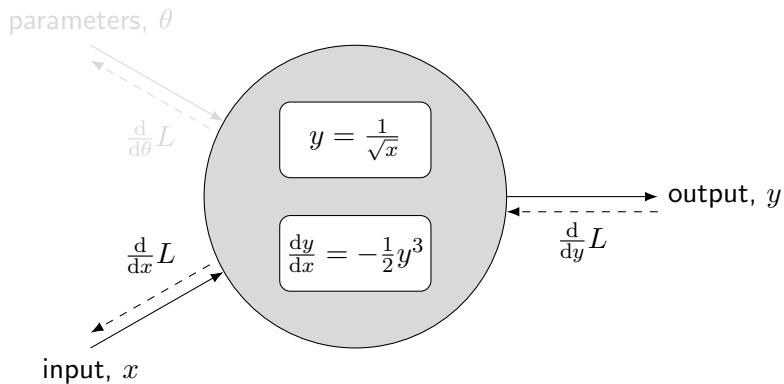
- ▶ algorithmic procedure that produces code for computing **exact** derivatives
- ▶ assumes numeric computations are composed of a small set of elementary operations that we know how to differentiate
 - ▶ arithmetic, exp, log, trigonometric
- ▶ workhorse of modern machine learning that greatly reduces development effort
- ▶ roughly speaking, for each line of the forward pass code, `P, Q = foo(A, B, C)`, autodiff produces a line `dLdA, dLdB, dLdC = foo_vjp(dLdP, dLdQ)` in the backward pass code
- ▶ but it doesn't always work (see point 2), and when it does work it can be slow and/or memory intensive

▶▶ example

Computing $1/\sqrt{x}$

```
1 float Q_rsqrt( float number )
2 {
3     long i;
4     float x2, y;
5     const float threehalfs = 1.5F;
6
7     x2 = number * 0.5F;
8     y = number;
9     i = * ( long * ) &y;          // evil floating point bit level hacking
10    i = 0x5f3759df - ( i >> 1 );   // what the f**k?
11    y = * ( float * ) &i;
12    y = y * ( threehalfs - ( x2 * y * y ) );    // 1st iter
13    // y = y * ( threehalfs - ( x2 * y * y ) ); // 2nd iter, can be removed
14
15    return y;
16 }
```

Separate Forward and Backward Operations



Part II. Differentiable Optimisation

Bi-level Optimisation: Stackelberg Games

Consider two players, a **leader** and a **follower**

- ▶ the market dictates the price it's willing to pay for some goods based on supply, i.e., quantity produced by both players, $P(q_1 + q_2)$

Bi-level Optimisation: Stackelberg Games

Consider two players, a **leader** and a **follower**

- ▶ the market dictates the price it's willing to pay for some goods based on supply, i.e., quantity produced by both players, $P(q_1 + q_2)$
- ▶ each player has a cost structure associated with producing goods, $C_i(q_i)$ and wants to maximize profits, $q_i P(q_1 + q_2) - C_i(q_i)$

Bi-level Optimisation: Stackelberg Games

Consider two players, a **leader** and a **follower**

- ▶ the market dictates the price it's willing to pay for some goods based on supply, i.e., quantity produced by both players, $P(q_1 + q_2)$
- ▶ each player has a cost structure associated with producing goods, $C_i(q_i)$ and wants to maximize profits, $q_i P(q_1 + q_2) - C_i(q_i)$
- ▶ the leader picks a quantity of goods to produce knowing that the follower will respond optimally. In other words, the leader solves

$$\begin{array}{ll} \text{maximize (over } q_1) & q_1 P(q_1 + q_2) - C_1(q_1) \\ \text{subject to} & q_2 \in \operatorname{argmax}_q q P(q_1 + q) - C_2(q) \end{array}$$

Bi-level Optimisation Problems in Machine Learning

- **quantities:** input x , output y , parameters θ

$$\begin{array}{ll} \text{minimize (over } \theta) & L(x, y; \theta) \\ \text{subject to} & y \in \operatorname{argmin}_{u \in C(x; \theta)} f(x, u; \theta) \end{array}$$

- lower-level is an optimisation problem parametrized by x and θ

Bi-level Optimisation Problems in Machine Learning

- **quantities:** input x , output y , parameters θ

$$\begin{array}{ll} \text{minimize (over } \theta) & L(x, y; \theta) \\ \text{subject to} & y \in \operatorname{argmin}_{u \in C(x; \theta)} f(x, u; \theta) \end{array}$$

- lower-level is an optimisation problem parametrized by x and θ
- **gradient descent:** compute gradient of lower-level solution y with respect to θ , and use the chain rule to get the total derivative,

$$\theta \leftarrow \theta - \eta \left(\frac{\partial L}{\partial \theta} + \frac{\partial L}{\partial y} \frac{dy}{d\theta} \right)$$

- by back-propagating *through* the optimisation problem

Differentiable Least Squares

Consider our old friend, the least-squares problem,

$$\text{minimize} \quad \|Ax - b\|_2^2$$

parameterized by A and b and with closed-form solution $x^\star = (A^T A)^{-1} A^T b$.

Differentiable Least Squares

Consider our old friend, the least-squares problem,

$$\text{minimize} \quad \|Ax - b\|_2^2$$

parameterized by A and b and with closed-form solution $x^\star = (A^T A)^{-1} A^T b$.

We are interested in derivatives of the solution with respect to the elements of A ,

$$\frac{dx^\star}{dA_{ij}} = \frac{d}{dA_{ij}} (A^T A)^{-1} A^T b \in \mathbb{R}^n$$

We could also compute derivatives with respect to elements of b (but not here).

Least Squares Backward Pass

The backward pass combines $\frac{dx^\star}{dA_{ij}}$ with $v^T = \frac{dL}{dx^\star}$ via the vector-Jacobian product. After some algebraic manipulation we get

$$\left(\frac{dL}{dA}\right)^T = wr^T - x^\star(Aw)^T \in \mathbb{R}^{m \times n}$$

where $w^T = v^T(A^T A)^{-1}$ and $r = b - Ax^\star$.

Least Squares Backward Pass

The backward pass combines $\frac{dx^*}{dA_{ij}}$ with $v^T = \frac{dL}{dx^*}$ via the vector-Jacobian product. After some algebraic manipulation we get

$$\left(\frac{dL}{dA}\right)^T = wr^T - x^*(Aw)^T \in \mathbb{R}^{m \times n}$$

where $w^T = v^T(A^T A)^{-1}$ and $r = b - Ax^*$.

- ▶ $(A^T A)^{-1}$ is used in both the forward and backward pass
- ▶ factored once to solve for x , e.g., into $A = QR$
- ▶ cache R and re-use when computing gradients

►► derivation

PyTorch Implementation: Forward Pass

```
1 class LeastSquaresFcn(torch.autograd.Function):
2     """PyTorch autograd function for least squares."""
3
4     @staticmethod
5     def forward(ctx, A, b):
6         B, M, N = A.shape
7         assert b.shape == (B, M, 1)
8
9         with torch.no_grad():
10             Q, R = torch.linalg.qr(A, mode='reduced')
11             x = torch.linalg.solve_triangular(R,
12                 torch.bmm(b.view(B, 1, M), Q).view(B, N, 1), upper=True)
13
14             # save state for backward pass
15             ctx.save_for_backward(A, b, x, R)
16
17             # return solution
18             return x
```

$$A = QR$$

$$x = R^{-1} (Q^T b)$$

(solves $Rx = Q^T b$)

PyTorch Implementation: Backward Pass

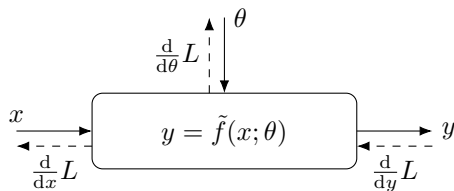
```
1  @staticmethod
2  def backward(ctx, dx):
3      # check for None tensors
4      if dx is None:
5          return None, None
6
7      # unpack cached tensors
8      A, b, x, R = ctx.saved_tensors
9      B, M, N = A.shape
10
11     dA, db = None, None
12
13     w = torch.linalg.solve_triangular(R,
14         torch.linalg.solve_triangular(torch.transpose(R, 2, 1),
15             dx, upper=False), upper=True)
16     Aw = torch.bmm(A, w)
17
18     if ctx.needs_input_grad[0]:
19         r = b - torch.bmm(A, x)
20         dA = torch.bmm(r.view(B,M,1), w.view(B,1,N)) - \
21             torch.bmm(Aw.view(B,M,1), x.view(B,1,N))
22     if ctx.needs_input_grad[1]:
23         db = Aw
24
25     # return gradients
26     return dA, db
```

$$\begin{aligned}w &= (A^T A)^{-1} v \\ &= R^{-1} (R^{-T} v) \\ r &= b - Ax\end{aligned}$$

$$\left(\frac{dL}{dA}\right)^T = rw^T - (Aw)x^T$$

$$\left(\frac{dL}{db}\right)^T = Aw$$

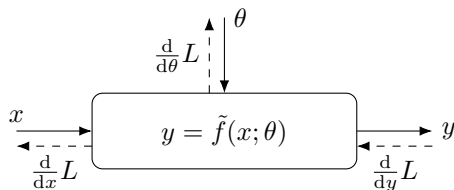
Imperative vs Declarative Nodes



- ▶ imperative node
- ▶ input-output relationship explicit,

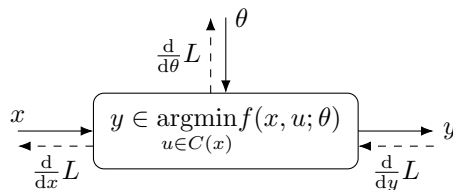
$$y = \tilde{f}(x; \theta)$$

Imperative vs Declarative Nodes



- ▶ imperative node
- ▶ input-output relationship explicit,

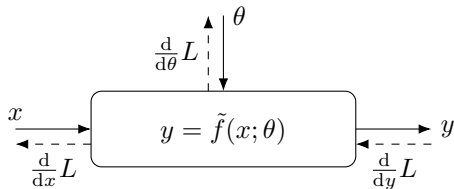
$$y = \tilde{f}(x; \theta)$$



- ▶ declarative node
- ▶ input-output relationship specified as solution to an optimisation problem,

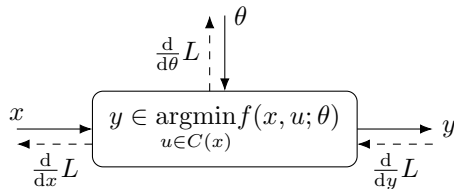
$$y \in \operatorname{argmin}_{u \in C(x)} f(x, u; \theta)$$

Imperative vs Declarative Nodes



- ▶ imperative node
- ▶ input-output relationship explicit,

$$y = \tilde{f}(x; \theta)$$



- ▶ declarative node
- ▶ input-output relationship specified as solution to an optimisation problem,

$$y \in \operatorname{argmin}_{u \in C(x)} f(x, u; \theta)$$

can co-exist in the same computation graph (network)

Average Pooling Example

$$\{x_i \in \mathbb{R}^m \mid i = 1, \dots, n\} \rightarrow \mathbb{R}^m$$

► imperative specification

$$y = \frac{1}{n} \sum_{i=1}^n x_i$$

► declarative specification

$$y = \operatorname{argmin}_{u \in \mathbb{R}^m} \sum_{i=1}^n \|u - x_i\|^2$$

Average Pooling Example

$$\{x_i \in \mathbb{R}^m \mid i = 1, \dots, n\} \rightarrow \mathbb{R}^m$$

- ▶ imperative specification

$$y = \frac{1}{n} \sum_{i=1}^n x_i$$

- ▶ declarative specification

$$y = \operatorname{argmin}_{u \in \mathbb{R}^m} \sum_{i=1}^n \|u - x_i\|^2$$

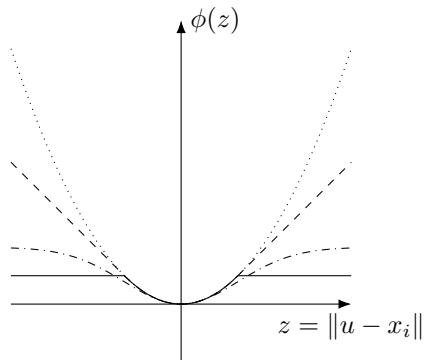
- ▶ can be easily varied, e.g., made robust

$$y = \operatorname{argmin}_{u \in \mathbb{R}^m} \sum_{i=1}^n \phi(u - x_i)$$

for some penalty function ϕ

Average Pooling Example

$$\{x_i \in \mathbb{R}^m \mid i = 1, \dots, n\} \rightarrow \mathbb{R}^m$$



- declarative specification

$$y = \operatorname{argmin}_{u \in \mathbb{R}^m} \sum_{i=1}^n \|u - x_i\|^2$$

- can be easily varied, e.g., made robust

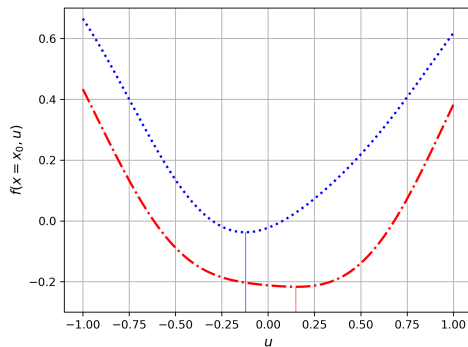
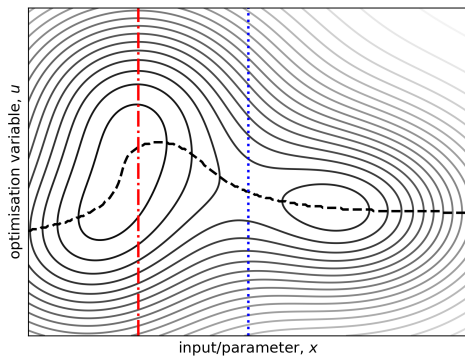
$$y = \operatorname{argmin}_{u \in \mathbb{R}^m} \sum_{i=1}^n \phi(u - x_i)$$

for some penalty function ϕ

Parametrized Optimisation Re-cap

Think of y and an implicit function of x (wlog we'll ignore θ from here on),

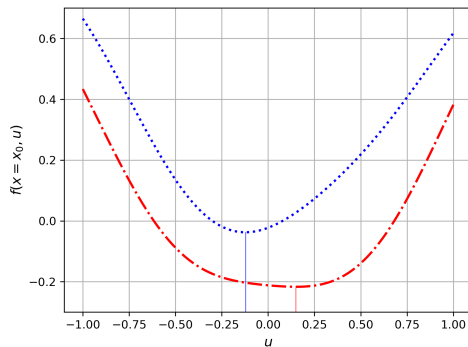
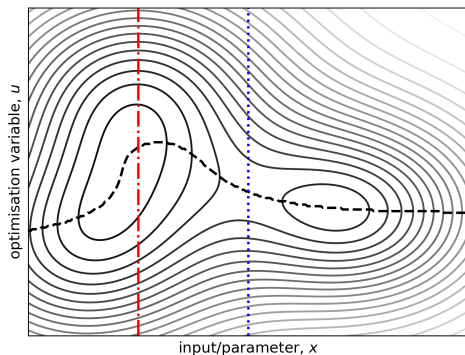
$$y(x) = \operatorname{argmin}_{u \in C(x)} f(x, u)$$



Parametrized Optimisation Re-cap

Think of y and an implicit function of x (wlog we'll ignore θ from here on),

$$y(x) = \operatorname{argmin}_{u \in C(x)} f(x, u)$$



Main question: How do we compute $\frac{d}{dx} \operatorname{argmin}_{u \in C(x)} f(x, u)$?

Computing $\frac{d}{dx} \operatorname{argmin}_{u \in C(x)} f(x, u)$

- ▶ explicit from closed-form solution
 - ▶ e.g., least-squares

Computing $\frac{d}{dx} \operatorname{argmin}_{u \in C(x)} f(x, u)$

- ▶ explicit from closed-form solution
 - ▶ e.g., least-squares
- ▶ automatic differentiation of forward pass code
 - ▶ e.g., unrolling gradient descent (next)

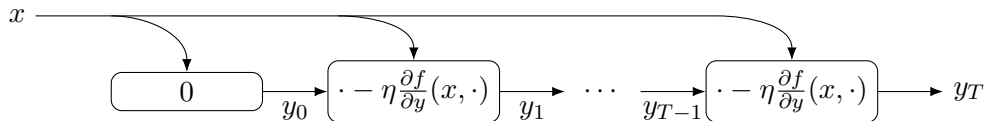
Computing $\frac{d}{dx} \operatorname{argmin}_{u \in C(x)} f(x, u)$

- ▶ explicit from closed-form solution
 - ▶ e.g., least-squares
- ▶ automatic differentiation of forward pass code
 - ▶ e.g., unrolling gradient descent (next)
- ▶ implicit differentiation of optimality conditions (later)
 - ▶ allows non-differentiable steps in the forward pass
 - ▶ no need to store intermediate calculations

Unrolling Gradient Descent

repeat until convergence:

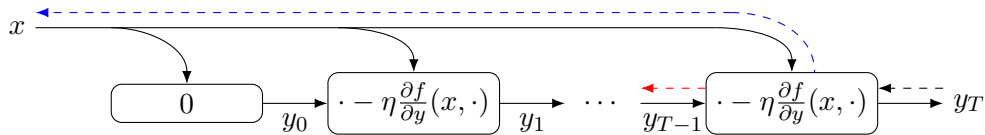
$$y_t \leftarrow y_{t-1} - \eta \frac{\partial f}{\partial y}(x, y_{t-1})$$



Unrolling Gradient Descent

repeat until convergence:

$$y_t \leftarrow y_{t-1} - \eta \frac{\partial f}{\partial y}(x, y_{t-1})$$



$$\frac{dy_t}{dx} = \frac{\partial y_t}{\partial x} + \frac{\partial y_t}{\partial y_{t-1}} \frac{dy_{t-1}}{dx} = -\eta \frac{\partial^2 f}{\partial x \partial y}(x, y_{t-1}) + \left(I - \eta \frac{\partial^2 f}{\partial y^2}(x, y_{t-1}) \right) \frac{dy_{t-1}}{dx}$$

Dini's Implicit Function Theorem

Consider the solution mapping associated with the equation $f(x, u) = 0$,

$$Y : x \mapsto \{u \in \mathbb{R}^m \mid f(x, u) = 0\} \text{ for } x \in \mathbb{R}^n.$$

We are interested in how elements of $Y(x)$ change as a function of x .

Dini's Implicit Function Theorem

Consider the solution mapping associated with the equation $f(x, u) = 0$,

$$Y : x \mapsto \{u \in \mathbb{R}^m \mid f(x, u) = 0\} \text{ for } x \in \mathbb{R}^n.$$

We are interested in how elements of $Y(x)$ change as a function of x .

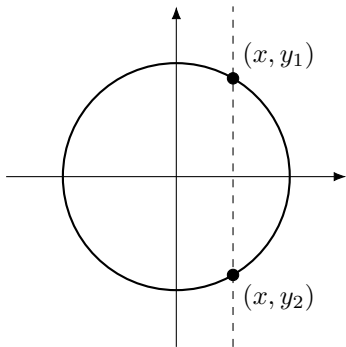
Theorem

Let $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ be differentiable in a neighbourhood of (x, u) and such that $f(x, u) = 0$, and let $\frac{\partial}{\partial u} f(x, u)$ be nonsingular. Then the solution mapping Y has a single-valued localization y around x for u which is differentiable in a neighbourhood \mathcal{X} of x with Jacobian satisfying

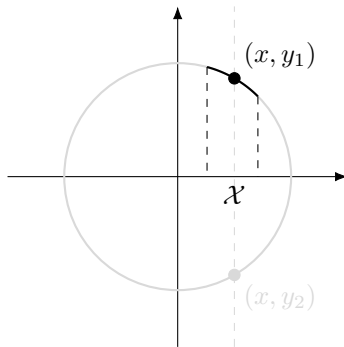
$$\frac{dy(x)}{dx} = - \left(\frac{\partial f(x, y(x))}{\partial y} \right)^{-1} \frac{\partial f(x, y(x))}{\partial x}$$

for every $x \in \mathcal{X}$.

Unit Circle Example



$$y = \pm\sqrt{1-x^2}$$
$$\frac{dy}{dx} = \frac{\mp 2x}{2\sqrt{1-x^2}} = -\frac{x}{y}$$



$$f(x, y) = x^2 + y^2 - 1$$
$$\frac{dy}{dx} = -\left(\frac{\partial f}{\partial y}\right)^{-1}\left(\frac{\partial f}{\partial x}\right)$$
$$= -\left(\frac{1}{2y}\right)(2x) = -\frac{x}{y}$$

Differentiating Unconstrained Optimisation Problems

Let $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be twice differentiable and let

$$y(x) \in \operatorname{argmin}_u f(x, u)$$

then for non-zero Hessian

$$\frac{dy(x)}{dx} = - \left(\frac{\partial^2 f}{\partial y^2} \right)^{-1} \frac{\partial^2 f}{\partial x \partial y}.$$

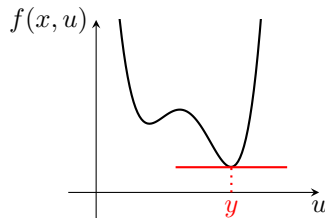
Differentiating Unconstrained Optimisation Problems

Let $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be twice differentiable and let

$$y(x) \in \operatorname{argmin}_u f(x, u)$$

then for non-zero Hessian

$$\frac{dy(x)}{dx} = - \left(\frac{\partial^2 f}{\partial y^2} \right)^{-1} \frac{\partial^2 f}{\partial x \partial y}.$$

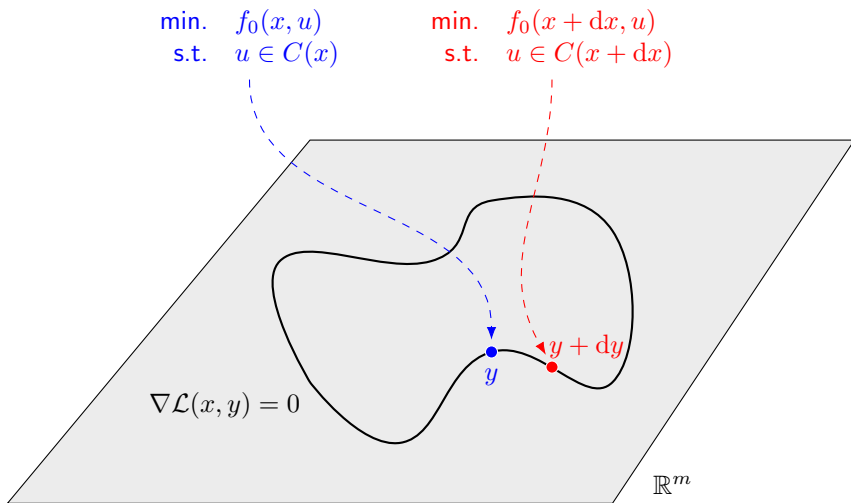


Proof. The derivative of f vanishes at (x, y) , i.e., $y \in \operatorname{argmin}_u f(x, u) \implies \frac{\partial f(x, y)}{\partial y} = 0$.

$$\begin{aligned} \text{LHS : } \frac{d}{dx} \frac{\partial f(x, y)}{\partial y} &= \frac{\partial^2 f(x, y)}{\partial x \partial y} + \frac{\partial^2 f(x, y)}{\partial y^2} \frac{dy}{dx} \\ \text{RHS : } \frac{d}{dx} 0 &= 0 \end{aligned}$$

Equating and rearranging gives the result. Or directly from Dini's implicit function theorem on $\frac{\partial f(x, y)}{\partial y} = 0$.

Differentiable Optimisation: Big Picture Idea



Differentiating (Unconstrained) Optimisation Problems

Consider functions $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$. Let

$$y(x) \in \arg \min_{u \in \mathbb{R}^m} f(x, u)$$

Assume that $y(x)$ exists and that f is twice differentiable in the neighbourhood of $(x, y(x))$. Then for H non-singular

$$\frac{dy(x)}{dx} = -H^{-1}B$$

where

$$B = \frac{\partial^2 f(x, y)}{\partial x \partial y} \in \mathbb{R}^{m \times n} \quad H = \frac{\partial^2 f(x, y)}{\partial y^2} \in \mathbb{R}^{m \times m}$$

Differentiating (Unconstrained) Optimisation Problems

Consider functions $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$. Let

$$y(x) \in \arg \min_{u \in \mathbb{R}^m} f(x, u)$$

Assume that $y(x)$ exists and that f is twice differentiable in the neighbourhood of $(x, y(x))$. Then for H non-singular

$$\frac{dy(x)}{dx} = -H^{-1}B$$

where

$$B = \frac{\partial^2 f(x, y)}{\partial x \partial y} \in \mathbb{R}^{m \times n} \quad H = \frac{\partial^2 f(x, y)}{\partial y^2} \in \mathbb{R}^{m \times m}$$

This result can be extended to constrained optimisation problems by differentiating optimality conditions, e.g., $\nabla \mathcal{L} = 0$.

► result

Automatic Differentiation for Differentiable Optimisation

(assuming a closed-form optimal solutions does not exist)

- ▶ At one extreme we can try back propagate through the optimisation algorithm (i.e., unrolling the optimisation procedure using automatic differentiation)
- ▶ At the other extreme we can use the implicit differentiation result to hand-craft efficient backward pass code
- ▶ There are also options in between, e.g.,
 - ▶ use automatic differentiation to obtain quantities in expression for $\frac{dy(x)}{dx}$ from software implementations of the objective and (active) constraint functions
 - ▶ implement the optimality condition $\nabla \mathcal{L} = 0$ in software and automatically differentiate that

Vector-Jacobian Product

For brevity consider the unconstrained optimisation case. The backward pass computes

$$\begin{aligned}\frac{dL}{dx} &= \frac{dL}{dy} \frac{dy}{dx} \\ &= \underbrace{(v^T)}_{\mathbb{R}^{1 \times m}} \underbrace{(-H^{-1}B)}_{\mathbb{R}^{m \times n}}\end{aligned}$$

$$\text{evaluation order:} \quad -v^T (H^{-1}B) \qquad (-v^T H^{-1}) B$$

$$\text{cost}^\dagger: \quad O(m^2n + mn) \qquad O(m^2 + mn)$$

[†] assumes H^{-1} is already factored (in $O(m^3)$ if unstructured, less if structured)

Summary and Open Questions

- ▶ optimisation problems can be embedded *inside* deep learning models
- ▶ back-propagation by either unrolling the optimisation algorithm or implicit differentiation of the optimality conditions
 - ▶ the former is easy to implement using automatic differentiation but memory intensive
 - ▶ the latter requires that solution be strongly convex locally (i.e., invertible H)
 - ▶ but does not need to know how the problem was solved, nor store intermediate forward-pass calculations
 - ▶ computing H^{-1} may be costly

Summary and Open Questions

- ▶ optimisation problems can be embedded *inside* deep learning models
- ▶ back-propagation by either unrolling the optimisation algorithm or implicit differentiation of the optimality conditions
 - ▶ the former is easy to implement using automatic differentiation but memory intensive
 - ▶ the latter requires that solution be strongly convex locally (i.e., invertible H)
 - ▶ but does not need to know how the problem was solved, nor store intermediate forward-pass calculations
 - ▶ computing H^{-1} may be costly
- ▶ active area of research and many open questions
 - ▶ Are declarative nodes slower?
 - ▶ Do declarative nodes give theoretical guarantees?
 - ▶ How best to handle non-smooth or discrete optimization problems?
 - ▶ What about problems with multiple solutions?
 - ▶ What if the forward pass solution is suboptimal?
 - ▶ Can problems become infeasible during learning?
 - ▶ ...

Part III. Applications

Differentiable Eigen Decomposition

Finding the eigenvector corresponding to the maximum eigenvalue of a real symmetric matrix $X \in \mathbb{R}^{m \times m}$ can be formulated as

$$\begin{array}{ll} \text{maximize (over } u \in \mathbb{R}^m) & u^T X u \\ \text{subject to} & u^T u = 1 \end{array}$$

which has applications in, for example, back propagating through normalized cuts.

Differentiable Eigen Decomposition

Finding the eigenvector corresponding to the maximum eigenvalue of a real symmetric matrix $X \in \mathbb{R}^{m \times m}$ can be formulated as

$$\begin{array}{ll} \text{maximize (over } u \in \mathbb{R}^m) & u^T X u \\ \text{subject to} & u^T u = 1 \end{array}$$

which has applications in, for example, back propagating through normalized cuts. Optimality conditions (for solution y) are

$$Xy = \lambda_{\max} y \quad \text{and} \quad y^T y = 1.$$

Differentiable Eigen Decomposition

Finding the eigenvector corresponding to the maximum eigenvalue of a real symmetric matrix $X \in \mathbb{R}^{m \times m}$ can be formulated as

$$\begin{array}{ll} \text{maximize (over } u \in \mathbb{R}^m) & u^T X u \\ \text{subject to} & u^T u = 1 \end{array}$$

which has applications in, for example, back propagating through normalized cuts. Optimality conditions (for solution y) are

$$Xy = \lambda_{\max} y \quad \text{and} \quad y^T y = 1.$$

Taking derivatives with respect to components of X we get,

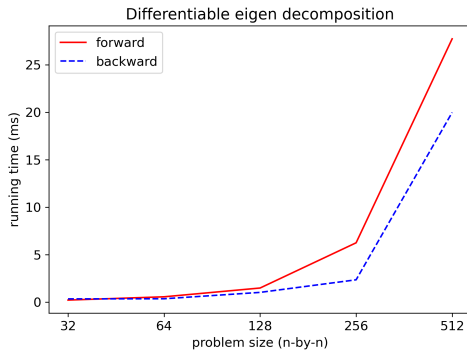
$$\frac{dy}{dX_{ij}} = -\frac{1}{2}(X - \lambda_{\max} I)^\dagger (E_{ij} + E_{ji})y \in \mathbb{R}^m$$

► derivation

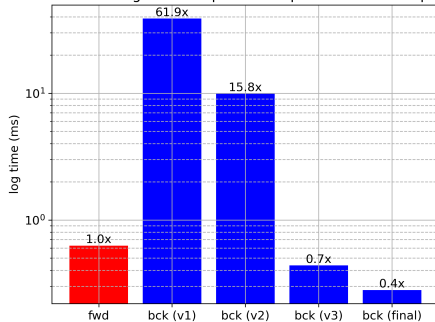
PyTorch Implementation

```
1 class EigenDecompositionFcn(torch.autograd.Function):
2     """PyTorch autograd function for eigen decomposition."""
3
4     @staticmethod
5     def forward(ctx, X):
6         B, M, N = X.shape
7
8         # use torch's eigh function to find the eigenvalues and eigenvectors of a symmetric matrix
9         with torch.no_grad():
10             lmd, Y = torch.linalg.eigh(0.5 * (X + X.transpose(1, 2)))
11
12             ctx.save_for_backward(lmd, Y)
13             return Y
14
15     @staticmethod
16     def backward(ctx, dJdY):
17         lmd, Y = ctx.saved_tensors
18         B, M, N = Y.shape
19
20         # compute all pseudo-inverses simultaneously
21         L = lmd.view(B, 1, M) - lmd.view(B, M, 1)
22         L = torch.where(torch.abs(L) < eps, 0.0, 1.0 / L)
23
24         # compute full gradient over all eigenvectors
25         dJdX = torch.bmm(torch.bmm(Y, L * torch.bmm(Y.transpose(1, 2), dJdY)), Y.transpose(1, 2))
26         dJdX = 0.5 * (dJdX + dJdX.transpose(1, 2))
27
28         return dJdX
```

Experiment



Differentiable eigen decomposition implementation comparison



Optimal Transport

One view of optimal transport is as a matching problem

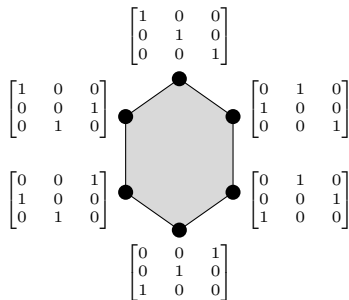
- ▶ from an m -by- n cost matrix M
- ▶ to an m -by- n probability matrix P ,

often formulated with an entropic regularisation term,

$$\begin{aligned} & \text{minimize} && \langle M, P \rangle + \frac{1}{\gamma} \langle P, \log P \rangle \\ & \text{subject to} && P \mathbf{1} = r \\ & && P^T \mathbf{1} = c \end{aligned}$$

with $\mathbf{1}^T r = \mathbf{1}^T c = 1$.

The row and column sum constraints ensure that P is a doubly stochastic matrix (lies within the convex hull of permutation matrices).



Solving Entropic Optimal Transport

Solution takes the form

$$P_{ij} = \alpha_i \beta_j e^{-\gamma M_{ij}}$$

and can be found using the Sinkhorn algorithm,

- ▶ Set $K_{ij} = e^{-\gamma M_{ij}}$ and $\alpha, \beta \in \mathbb{R}_{++}^n$
- ▶ Iterate until convergence,

$$\alpha \leftarrow r \oslash K\beta$$

$$\beta \leftarrow c \oslash K^T \alpha$$

where \oslash denotes componentwise division

- ▶ Return $P = \mathbf{diag}(\alpha) K \mathbf{diag}(\beta)$

Differentiable Optimal Transport

- ▶ Option 1: back-propagate through Sinkhorn algorithm

Differentiable Optimal Transport

- ▶ Option 1: back-propagate through Sinkhorn algorithm
- ▶ Option 2: use the implicit differentiation result

$$\underbrace{\frac{dL}{dM}}_{m\text{-by-}n} = \underbrace{\frac{dL}{dP}}_{m\text{-by-}n} \underbrace{\frac{dP}{dM}}_{m\text{-by-}n\text{-by-}m\text{-by-}n}$$

Differentiable Optimal Transport

- ▶ Option 1: back-propagate through Sinkhorn algorithm
- ▶ Option 2: use the implicit differentiation result

$$\underbrace{\frac{dL}{dM}}_{1\text{-by-}mn} = \underbrace{\frac{dL}{dP}}_{1\text{-by-}mn} \underbrace{\frac{dP}{dM}}_{mn\text{-by-}mn} \quad (\text{think of vectorising } M \text{ and } P)$$

Optimal Transport Gradient

Derivation of the optimal transport gradient is quite tedious (see notes). The result:

$$\begin{aligned}\frac{dL}{dM} &= \frac{dL}{dP} \left(H^{-1} \mathbf{A}^T (A H^{-1} A^T)^{-1} \mathbf{A} H^{-1} - H^{-1} \right) B \\ &= \gamma \frac{dL}{dP} \mathbf{diag}(P) \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}^T \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \mathbf{diag}(P) - \gamma \frac{dL}{dP} \mathbf{diag}(P)\end{aligned}$$

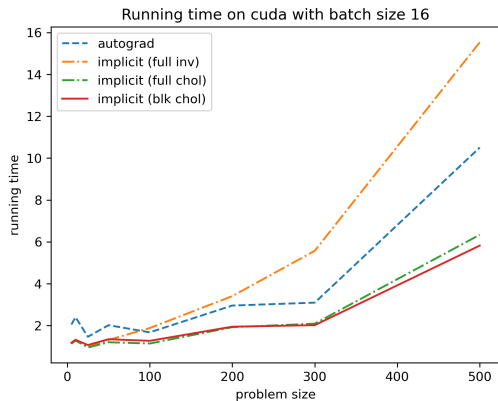
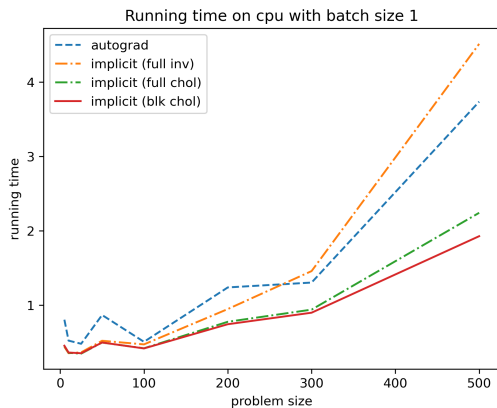
where

$$\begin{aligned}\begin{bmatrix} A_1 \\ A_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{0}_n^T & \mathbf{1}_n^T & \cdots & \mathbf{0}_n^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_n^T & \mathbf{0}_n^T & \cdots & \mathbf{1}_n^T \\ I_{n \times n} & I_{n \times n} & \cdots & I_{n \times n} \end{bmatrix} & (A H^{-1} A^T)^{-1} &= \frac{1}{\gamma} \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{12}^T & \Lambda_{22} \end{bmatrix} \\ & & &= \frac{1}{\gamma} \begin{bmatrix} \mathbf{diag}(r_{2:m}) & P_{2:m,1:n} \\ P_{2:m,1:n}^T & \mathbf{diag}(c) \end{bmatrix}^{-1}\end{aligned}$$

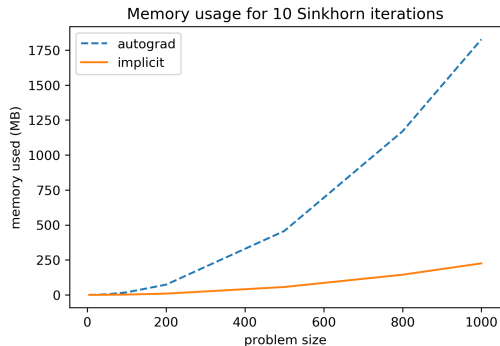
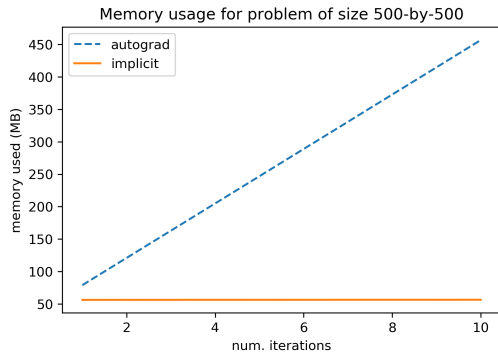
Implementation

```
1 @staticmethod
2 def backward(ctx, dJdP)
3     # unpacked cached tensors
4     M, r, c, P = ctx.saved_tensors
5     batches, m, n = P.shape
6
7     # initialize backward gradients  $(-v^T H^{-1} B)$ 
8     dLdM = -1.0 * gamma * P * dLdP
9
10    # compute  $[vHAt1, vHAt2] = -v^T H^{-1} A^T$ 
11    vHAt1, vHAt2 = sum(dJdM[:, 1:m, 0:n], dim=2), sum(dJdM, dim=1)
12
13    # compute  $[v1, v2] = -v^T H^{-1} A^T (A H^{-1} A^T)^{-1}$ 
14    P_over_c = P[:, 1:m, 0:n] / c.view(batches, 1, n)
15    lmd_11 = cholesky(diag_embed(r[:, 1:m])) - bmm(P[:, 1:m, 0:n], P_over_c.transpose(1, 2))
16    lmd_12 = cholesky_solve(P_over_c, lmd_11)
17    lmd_22 = diag_embed(1.0 / c) + bmm(lmd_12.transpose(1, 2), P_over_c)
18
19    v1 = torch.cholesky_solve(vHAt1, lmd_11) - torch.bmm(lmd_12, vHAt2)
20    v2 = torch.bmm(lmd_22, vHAt2) - torch.bmm(lmd_12.transpose(1, 2), vHAt1)
21
22    # compute  $v^T H^{-1} A^T (A H^{-1} A^T)^{-1} A H^{-1} B - v^T H^{-1} B$ 
23    dLdM[:, 1:m, 0:n] -= v1.view(batches, m-1, 1) * P[:, 1:m, 0:n]
24    dJdM -= v2.view(batches, 1, n) * P
25
26    # return gradients
27    return dJdM
```

Running Time



Memory Usage



Application to Blind Perspective-n-Point

(Campbell et al., ECCV 2020)



find the location where the photograph was taken

Coupled Problem

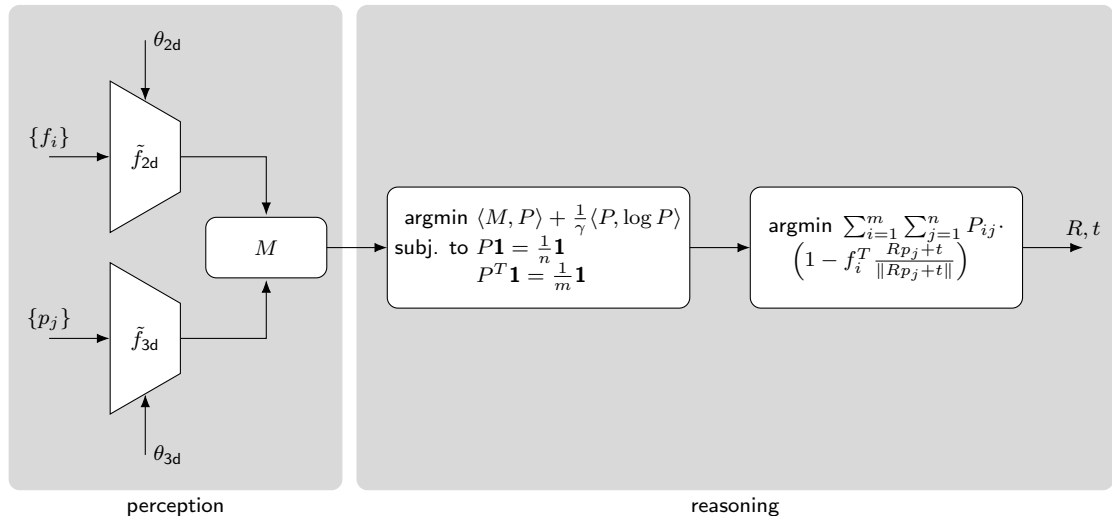


- ▶ if we knew **correspondences** then determining **camera pose** would be easy

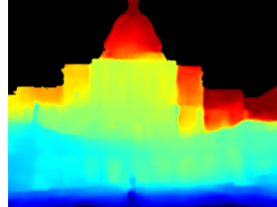
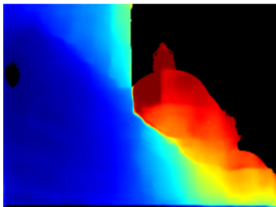
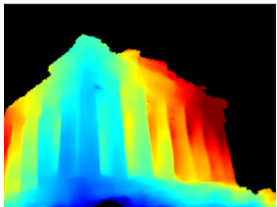


- ▶ if we knew **camera pose** then determining **correspondences** would be easy

Blind Perspective-n-Point Network Architecture



Blind Perspective-n-Point Results



Further Resources

Where to from here?

► background reading



- Deep declarative networks (<http://deepdeclarativenetworks.com>)
 - lots of small code examples and tutorials
- CVXPylayers (<https://github.com/cvxgrp/cvxpylayers>)
- Theseus (<https://sites.google.com/view/theseus-ai>)
- JAXopt (<https://github.com/google/jaxopt>)

lecture notes available at <https://users.cecs.anu.edu.au/~sgould>

break-out slides

automatic differentiation

Toy Example: Babylonian Algorithm [▶ back](#)

Consider the following implementation for a forward operation:

```
1: procedure FWDFCN( $x$ )  
2:    $y_0 \leftarrow \frac{1}{2}x$   
3:   for  $t = 1, \dots, T$  do  
4:      $y_t \leftarrow \frac{1}{2} \left( y_{t-1} + \frac{x}{y_{t-1}} \right)$   
5:   end for  
6:   return  $y_T$   
7: end procedure
```

Toy Example: Babylonian Algorithm [▶▶ back](#)

Consider the following implementation for a forward operation:

```
1: procedure FWDFCN( $x$ )  
2:    $y_0 \leftarrow \frac{1}{2}x$   
3:   for  $t = 1, \dots, T$  do  
4:      $y_t \leftarrow \frac{1}{2} \left( y_{t-1} + \frac{x}{y_{t-1}} \right)$   
5:   end for  
6:   return  $y_T$   
7: end procedure
```

Automatic differentiation algorithmically generates the backward code:

```
1: procedure BCKFCN( $x, y_T, \frac{dL}{dy_T}$ )  
2:    $\frac{dL}{dx} \leftarrow 0$   
3:   for  $t = T, \dots, 1$  do  
4:      $\frac{dL}{dx} \leftarrow \frac{dL}{dx} + \overbrace{\frac{dL}{dy_t} \left( \frac{1}{2y_{t-1}} \right)}^{\partial y_t / \partial x}$   
5:      $\frac{dL}{dy_{t-1}} \leftarrow \frac{dL}{dy_t} \underbrace{\left( \frac{1}{2} - \frac{x}{2y_{t-1}^2} \right)}_{\partial y_t / \partial y_{t-1}}$   
6:   end for  
7:    $\frac{dL}{dx} \leftarrow \frac{dL}{dx} + \frac{dL}{dy_0} \frac{1}{2}$   
8:   return  $\frac{dL}{dx}$   
9: end procedure
```


Toy Example: Babylonian Algorithm [▶▶ back](#)

Consider the following implementation for a forward operation:

```
1: procedure FWDFCN( $x$ )
2:    $y_0 \leftarrow \frac{1}{2}x$ 
3:   for  $t = 1, \dots, T$  do
4:      $y_t \leftarrow \frac{1}{2} \left( y_{t-1} + \frac{x}{y_{t-1}} \right)$ 
5:   end for
6:   return  $y_T$ 
7: end procedure
```

► computes $y = \sqrt{x}$

► derivative computed directly is
$$\frac{dy}{dx} = \frac{1}{2\sqrt{x}} = \frac{1}{2y}$$

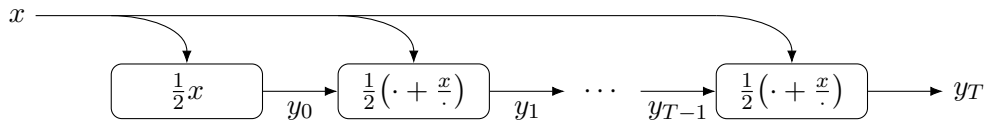
Automatic differentiation algorithmically generates the backward code:

```
1: procedure BCKFCN( $x, y_T, \frac{dL}{dy_T}$ )
2:    $\frac{dL}{dx} \leftarrow 0$ 
3:   for  $t = T, \dots, 1$  do
4:     
$$\frac{dL}{dx} \leftarrow \frac{dL}{dx} + \overbrace{\frac{dL}{dy_t} \left( \frac{1}{2y_{t-1}} \right)}^{\partial y_t / \partial x}$$

5:     
$$\frac{dL}{dy_{t-1}} \leftarrow \overbrace{\frac{dL}{dy_t} \left( \frac{1}{2} - \frac{x}{2y_{t-1}^2} \right)}^{\partial y_t / \partial y_{t-1}}$$

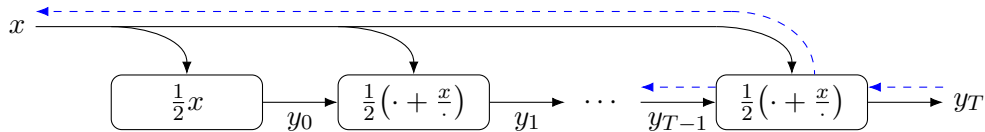
6:   end for
7:    $\frac{dL}{dx} \leftarrow \frac{dL}{dx} + \frac{dL}{dy_0} \frac{1}{2}$ 
8:   return  $\frac{dL}{dx}$ 
9: end procedure
```

Computation Graph for Babylonian Algorithm [▶ back](#)



$$y_T = f(x, f(x, f(x, \dots f(x, \frac{1}{2}x)))) \text{ with } f(x, y) = \frac{1}{2}\left(y + \frac{x}{y}\right)$$

Computation Graph for Babylonian Algorithm [▶ back](#)



$$y_T = f(x, f(x, f(x, \dots f(x, \frac{1}{2}x)))) \text{ with } f(x, y) = \frac{1}{2}\left(y + \frac{x}{y}\right)$$

$$\frac{dL}{dx} = \frac{dL}{dy_T} \left(\frac{\partial y_T}{\partial x} + \frac{\partial y_T}{\partial y_{T-1}} \left(\frac{\partial y_{T-1}}{\partial x} + \frac{\partial y_{T-1}}{\partial y_{T-2}} \left(\dots + \frac{\partial y_0}{\partial x} \right) \right) \right)$$

least squares

Least Squares Backward Pass Derivation [» back](#)

Differentiating x^\star with respect to single element A_{ij} , we have

$$\begin{aligned}\frac{d}{dA_{ij}} x^\star &= \frac{d}{dA_{ij}} (A^T A)^{-1} A^T b \\ &= \left(\frac{d}{dA_{ij}} (A^T A)^{-1} \right) A^T b + (A^T A)^{-1} \left(\frac{d}{dA_{ij}} A^T b \right)\end{aligned}$$

Using the identity $\frac{d}{dz} Z^{-1} = -Z^{-1} \left(\frac{d}{dz} Z \right) Z^{-1}$ we get, for the first term,

$$\begin{aligned}\frac{d}{dA_{ij}} (A^T A)^{-1} &= - (A^T A)^{-1} \left(\frac{d}{dA_{ij}} (A^T A) \right) (A^T A)^{-1} \\ &= - (A^T A)^{-1} (E_{ij}^T A + A^T E_{ij}) (A^T A)^{-1}\end{aligned}$$

where E_{ij} is a matrix with one in the (i, j) -th element and zeros elsewhere. Furthermore, for the second term,

$$\frac{d}{dA_{ij}} A^T b = E_{ij}^T b$$

Plugging these back into parent equation we have

$$\begin{aligned}\frac{d}{dA_{ij}}x^* &= -(A^TA)^{-1}(E_{ij}^TA + A^TE_{ij}) (A^TA)^{-1}A^Tb + (A^TA)^{-1}E_{ij}^Tb \\ &= -(A^TA)^{-1}(E_{ij}^TA + A^TE_{ij}) x^* + (A^TA)^{-1}E_{ij}^Tb \\ &= -(A^TA)^{-1}(E_{ij}^T(Ax^* - b) + A^TE_{ij}x^*) \\ &= -(A^TA)^{-1}((a_i^Tx^* - b_i)e_j + x_j^*a_i)\end{aligned}$$

where $e_j = (0, 0, \dots, 1, 0, \dots) \in \mathbb{R}^n$ is the j -th canonical vector, i.e., vector with a one in the j -th component and zeros everywhere else, and $a_i^T \in \mathbb{R}^{1 \times n}$ is the i -th row of matrix A .

Least Squares Backward Pass Derivation (cont.) [▶▶ back](#)

Let $r = b - Ax^*$ and let v^T denote the backward coming gradient $\frac{d}{dx^*}L$. Then

$$\begin{aligned}\frac{dL}{dA_{ij}} &= v^T \frac{dx^*}{dA_{ij}} \\ &= v^T (A^T A)^{-1} (r_i e_j - x_j^* a_i) \\ &= w^T (r_i e_j - x_j^* a_i) \\ &= r_i w_j - w^T a_i x_j^*\end{aligned}$$

where $w = (A^T A)^{-1} v$. We can compute the entire matrix of $m \times n$ derivatives efficiently as the sum of outer products

$$\left(\frac{dL}{dA} \right)^T = \left[\frac{dL}{dA_{ij}} \right]_{\substack{i=1,\dots,m \\ j=1,\dots,n}} = w r^T - x^* (A w)^T$$

differentiating equality constrained problems

Differentiating Equality Constrained Optimisation Problems [▶▶ back](#)

Consider functions $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ and $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^q$. Let

$$\begin{aligned} y(x) \in \arg \min_{u \in \mathbb{R}^m} f(x, u) \\ \text{subject to} \quad h(x, u) = 0_q \end{aligned}$$

Assume that $y(x)$ exists, that f and h are twice differentiable in the neighbourhood of $(x, y(x))$, and that $\mathbf{rank}(\frac{\partial h(x, y)}{\partial y}) = q$.

Differentiating Equality Constrained Optimisation Problems [▶▶ back](#)

Consider functions $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ and $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^q$. Let

$$\begin{aligned} y(x) \in \arg \min_{u \in \mathbb{R}^m} & f(x, u) \\ \text{subject to} & h(x, u) = 0_q \end{aligned}$$

Assume that $y(x)$ exists, that f and h are twice differentiable in the neighbourhood of $(x, y(x))$, and that $\mathbf{rank}(\frac{\partial h(x, y)}{\partial y}) = q$. Then for H non-singular

$$\frac{dy(x)}{dx} = H^{-1}A^T(AH^{-1}A^T)^{-1}(AH^{-1}B - C) - H^{-1}B$$

where

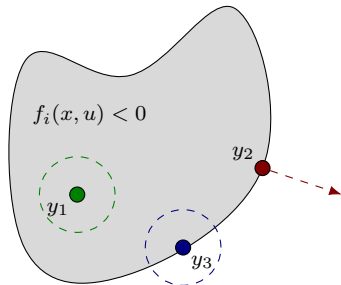
$$\begin{aligned} A &= \frac{\partial h(x, y)}{\partial y} \in \mathbb{R}^{q \times m} & B &= \frac{\partial^2 f(x, y)}{\partial x \partial y} - \sum_{i=1}^q \nu_i \frac{\partial^2 h_i(x, y)}{\partial x \partial y} \in \mathbb{R}^{m \times n} \\ C &= \frac{\partial h(x, y)}{\partial x} \in \mathbb{R}^{q \times n} & H &= \frac{\partial^2 f(x, y)}{\partial y^2} - \sum_{i=1}^q \nu_i \frac{\partial^2 h_i(x, y)}{\partial y^2} \in \mathbb{R}^{m \times m} \end{aligned}$$

and $\nu \in \mathbb{R}^q$ satisfies $\nu^T A = \frac{\partial f(x, y)}{\partial y}$.

Dealing with Inequality Constraints [» back](#)

$$\begin{aligned} y(x) \in \arg \min_{u \in \mathbb{R}^m} & f_0(x, u) \\ \text{subject to} & h_i(x, u) = 0, \quad i = 1, \dots, q \\ & f_i(x, u) \leq 0, \quad i = 1, \dots, p. \end{aligned}$$

- ▶ Replace inequality constraints with log-barrier approximation
- ▶ Treat as equality constraints if active (y_2 or y_3) and ignore otherwise (y_1 or y_3)
 - ▶ may lead to one-sided gradients since $\nu \succeq 0$



eigen decomposition

Deriving the Gradient for Eigen Decomposition [▶ back](#)

Implicit differentiation of the optimality conditions with respect to X_{ij} gives,

$$\frac{d}{dX_{ij}}(Xy - \lambda_{\max}y) = \frac{1}{2}(E_{ij} + E_{ji})y - \frac{d\lambda_{\max}}{dX_{ij}}y + (X - \lambda_{\max}I)\frac{dy}{dX_{ij}} = 0 \quad (1)$$

$$\frac{d}{dX_{ij}}(y^T y - 1) = 2y^T \frac{dy}{dX_{ij}} = 0 \quad (2)$$

Pre-multiplying (1) by y^T , and using (2) and $y^T y = 1$, we get

$$\frac{d\lambda_{\max}}{dX_{ij}} = \frac{1}{2}y^T(E_{ij} + E_{ji})y$$

Pre-multiplying (1) by $(X - \lambda_{\max}I)^\dagger$, we get

$$\begin{aligned} \frac{1}{2}(X - \lambda_{\max}I)^\dagger(E_{ij} + E_{ji})y - (X - \lambda_{\max}I)^\dagger \frac{d\lambda_{\max}}{dX_{ij}}y + \frac{dy}{dX_{ij}} &= 0 \\ \therefore \frac{dy}{dX_{ij}} &= -\frac{1}{2}(X - \lambda_{\max}I)^\dagger(E_{ij} + E_{ji})y \end{aligned}$$

since $(X - \lambda_{\max}I)^\dagger \frac{d\lambda_{\max}}{dX_{ij}}y = \frac{d\lambda_{\max}}{dX_{ij}}(X - \lambda_{\max}I)^\dagger y = 0$ since if $Az = 0$, then $A^\dagger z = 0$.