





# Implicit differentiation

Cheng Soon Ong  
Marc Peter Deisenroth

December 2020



# Motivation

- ▶ In machine learning, we use gradients to train predictors
- ▶ For functions  $f(x)$  we can directly obtain its gradient  $\nabla f(x)$
- ▶ **How to represent a constraint?**

$$G(x, y) = 0$$

E.g. conservation of mass

# High school gradients

This is an equation

$$y = x^3 + 2x^2 + x + 4$$

What is the gradient  $\frac{dy}{dx}$ ?

# High school gradients

This is an equation

$$y = x^3 + 2x^2 + x + 4$$

What is the gradient  $\frac{dy}{dx}$ ?

$$\frac{dy}{dx} = 3x^2 + 4x + 1$$

# High school gradients

This is an equation

$$y = x^3 + 2x^2 + x + 4$$

What is the gradient  $\frac{dy}{dx}$ ?

$$\frac{dy}{dx} = 3x^2 + 4x + 1$$

Observe that we can write the equation as

$$x^3 + 2x^2 + x + 4 - y = 0$$

which is of the form  $G(x, y) = 0$ .

# Optimization with constraints

Given a constrained continuous optimization problem

$$\begin{aligned} & \min_{x,y} F(x, y) \\ & \text{subject to } G(x, y) = 0 \end{aligned}$$

We can solve the equality constraint to get

$$y = g(x)$$

and substitute into the objective  $F(x, g(x))$ , and calculate the gradient

$$\frac{d}{dx} F(x, g(x))$$

# Advanced high school gradients

This is an equation

$$y = x^3 + 2x^2 + \textcolor{blue}{xy} + 4$$

What is the gradient  $\frac{dy}{dx}$ ?

# Advanced high school gradients

This is an equation

$$y = x^3 + 2x^2 + \textcolor{blue}{xy} + 4$$

What is the gradient  $\frac{dy}{dx}$ ?

## Implicit function theorem

- ▶ Solve for  $y = g(x)$  and use quotient rule
- ▶ Directly differentiate, and use product rule

## Function of two variables

Consider a function

$$G(x, y) = 0$$

where  $x, y \in \mathbb{R}$ .

Assume that near a particular point  $x_0$ , we can write a closed form expression for  $y$  in terms of  $x$ , that is

$$y = g(x).$$

## Solve for one variable

Substituting  $y = g(x)$  into  $G(x, y)$ , near  $x_0$ , we get

$$G(x, g(x)) = 0.$$

We calculate the derivative of  $G$  with respect to  $x$  using the **chain rule**,

$$G_x(x, g(x)) + G_y(x, g(x)) \cdot g'(x) = 0.$$

If  $G_y \neq 0$ , then

$$g'(x) = -\frac{G_x}{G_y}.$$

# Implicit function theorem

The implicit function theorem **provides conditions** under which we can write  $G(x, y) = 0$  as  $y = g(x)$ , and also conditions when  $g'(x) = -\frac{G_x}{G_y}$ .

## Hand wavy intuition

Given three variables (could be any topology)  $x, y, z$  and

$$G(x, y) = z,$$

What are the conditions for the following to be well behaved?

$$\begin{aligned} g(x) &= y \\ G(x, g(x)) &= z. \end{aligned}$$

# Implicit function theorem as ansatz

## Ansatz

1. Explicitly solve one variable in terms of another
  2. Chain the gradients
- ▶ Ansatz = a way to look at problems
  - ▶ The implicit function theorem is a way to solve equations

Generalizations depending on where equations live (Krantz and Parks, 2013):

- ▶ Inverse function theorem
- ▶ Constant rank theorem
- ▶ Banach fixed point theorem
- ▶ Nash-Moser theorem

## Pointers to literature

- ▶ From linear algebra to calculus (Hubbard and Hubbard, 2015; Spivak, 2008)
- ▶ Implicit function theorem from variational analysis (Dontchev and Rockafellar, 2014)
- ▶ Many versions of implicit function theorem (Krantz and Parks, 2013)
- ▶ Deep Declarative Networks

<https://anucvml.github.io/ddn-cvprw2020/>

Augustin-Louis Cauchy, 1916



## Gradient of loss w.r.t. optimal value

- ▶ Consider the problem of structured prediction (Nowozin et al., 2014)
- ▶ Let a sample be  $(\mathbf{x}_n, \mathbf{y}_n)$
- ▶ Energy given the parameter of the learner  $\mathbf{w}$  is  $E(\mathbf{y}_n, \mathbf{x}_n, \mathbf{w})$
- ▶ We assume that the best predictor is found by an optimization algorithm over  $\mathbf{y}$ ,

$$\mathbf{y}^*(\mathbf{x}_n, \mathbf{w}) = \text{opt}_{\mathbf{y}} E(\mathbf{y}, \mathbf{x}_n, \mathbf{w}).$$

- ▶ Measure the error with  $\ell(\mathbf{y}, \mathbf{y}_n)$
- ▶ Loss per sample  $L(\mathbf{x}_n, \mathbf{y}_n, \mathbf{w})$  (assume predict with  $\mathbf{y}^*$ )
- ▶ **Want to take the gradient of loss  $L$  w.r.t. parameters  $\mathbf{w}$**
- ▶ but have an optimization problem inside (to find  $\mathbf{y}^*$ )

## Gradient of $L$ w.r.t. $w$

By implicit differentiation (Do et al., 2007; Samuel and Tappen, 2009), the gradient of the loss  $L(\mathbf{x}_n, \mathbf{y}_n, \mathbf{w})$  with respect to the parameters  $\frac{dL}{d\mathbf{w}}$  has a closed form.

### Theorem

Let  $\mathbf{y}^*(\mathbf{w}) = \operatorname{argmin}_{\mathbf{y}} E(\mathbf{y}, \mathbf{w})$ , and  $L(\mathbf{w}) = \ell(\mathbf{y}^*(\mathbf{w}))$ . Then

$$\frac{dL}{d\mathbf{w}} = -\frac{\partial^2 E}{\partial \mathbf{w} \partial \mathbf{y}^\top} \left( \frac{\partial^2 E}{\partial \mathbf{y} \partial \mathbf{y}^\top} \right)^{-1} \frac{d\ell}{d\mathbf{y}}.$$

## Sketch: chain rule

$$L(\mathbf{w}) = \ell(\text{opt}_{\mathbf{y}} E(\mathbf{y}, \mathbf{w}))$$

By the chain rule,

$$\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial \ell}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{w}}.$$

## Sketch: gradient with respect to energy $E$

Denote by

$$g(\mathbf{y}, \mathbf{w}) = \frac{\partial E(\mathbf{y}, \mathbf{w})}{\partial \mathbf{y}}$$

The gradient of  $g$  with respect to  $\mathbf{w}$  is given by the chain rule again. Note that  $\mathbf{y}$  is actually a function of  $\mathbf{w}$ , i.e.  $g(\mathbf{y}(\mathbf{w}), \mathbf{w})$ .

$$\frac{\partial}{\partial \mathbf{w}} g(\mathbf{y}(\mathbf{w}), \mathbf{w}) = \frac{\partial g}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{w}} + \frac{\partial g}{\partial \mathbf{w}}.$$

## Sketch: Stationarity conditions

At optimality of  $E(\mathbf{y}, \mathbf{w})$ , its gradient is zero, i.e.  $g(\mathbf{y}, \mathbf{w}) = 0$ . Solving

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}} = -\frac{\partial g}{\partial \mathbf{w}} \left( \frac{\partial g}{\partial \mathbf{y}} \right)^{-1}$$

Substituting  $\frac{\partial \mathbf{y}}{\partial \mathbf{w}}$  into  $\frac{\partial L}{\partial \mathbf{w}}$ ,

$$\frac{\partial L}{\partial \mathbf{w}} = -\frac{\partial g}{\partial \mathbf{w}} \left( \frac{\partial g}{\partial \mathbf{y}} \right)^{-1} \frac{\partial \ell}{\partial \mathbf{y}}.$$

Sketch: Substitute back to obtain  $\frac{\partial L}{\partial \mathbf{w}}$

---

Recall the definition of  $g$  as gradient of energy  $E$ ,

$$g(\mathbf{y}, \mathbf{w}) = \frac{\partial E(\mathbf{y}, \mathbf{w})}{\partial \mathbf{y}}$$

And hence we get the Hessian with respect to the energy

$$\frac{\partial L}{\partial \mathbf{w}} = -\frac{\partial^2 E}{\partial \mathbf{w} \partial \mathbf{y}^\top} \left( \frac{\partial^2 E}{\partial \mathbf{y} \partial \mathbf{y}^\top} \right)^{-1} \frac{d\ell}{d\mathbf{y}}.$$

## Result: Gradient of $L$ w.r.t. $w$

The gradient of  $L(\mathbf{w}) = \ell(\operatorname{argmin}_{\mathbf{y}} E(\mathbf{y}(\mathbf{w}), \mathbf{w}))$  with respect to  $w$  has a closed form.

### Theorem

Let  $\mathbf{y}^*(\mathbf{w}) = \operatorname{argmin}_{\mathbf{y}} E(\mathbf{y}, \mathbf{w})$ , and  $L(\mathbf{w}) = \ell(\mathbf{y}^*(\mathbf{w}))$ . Then

$$\frac{dL}{d\mathbf{w}} = -\frac{\partial^2 E}{\partial \mathbf{w} \partial \mathbf{y}^\top} \left( \frac{\partial^2 E}{\partial \mathbf{y} \partial \mathbf{y}^\top} \right)^{-1} \frac{d\ell}{d\mathbf{y}}.$$

Domke (2012)

# Summary

- ▶ We want to take the gradient with respect to an equality constraint
- ▶ Implicit function theorem gives conditions where we can "invert" a derivative
- ▶ Implicit function theorem as a way to solve equations
- ▶ Useful to take a gradient over an optimum

Backpropagation is just...

the implicit function theorem

## References

- Do, C. B., Foo, C.-S., and Ng, A. (2007). Efficient multiple hyperparameter learning for log-linear models. In *Neural Information Processing Systems*.
- Domke, J. (2012). Generic methods for optimization-based modeling. In *AISTATS*.
- Dontchev, A. L. and Rockafellar, R. T. (2014). *Implicit Functions and Solution Mappings: A View from Variational Analysis*. Springer, 2nd edition.
- Hubbard, J. H. and Hubbard, B. B. (2015). *Vector Calculus, Linear Algebra, and Differential Forms: A Unified Approach*. Matrix Editions, 5th edition.
- Krantz, S. G. and Parks, H. R. (2013). *The Implicit Function Theorem: History, Theory, and Applications*. Springer.
- Nowozin, S., Gehler, P. V., Jancsary, J., and Lampert, C. H. (2014). *Advanced Structured Prediction*. MIT Press.

## References (cont.)

- Samuel, K. and Tappen, M. (2009). Learning optimized map estimates in continuously-valued MRF models. In *CVPR*.
- Spivak, M. (2008). *Calculus*. Publish or Perish, Inc., 4th edition.