

# Loss Functions and Backpropagation

Yao Zhang

The guy is a populace

Mostly based on Thomas Hofmann's lecture in ETH

<https://zhims.github.io/>

Dec 1, 2019

# Reminder: Notation

- Neural networks implements map  $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Compositional structure **layers**:

$$F = F^L \circ F^{L-1} \circ \dots \circ F^1 \quad (1)$$

- **Linear + activation** function

$$F^l = \sigma^l \circ \vec{F}^l, \quad \vec{F}^l(x) = W^l x + b^l, \quad l = 1, \dots, L \quad (2)$$

- F minus output layer non-linearity

$$\bar{F} = \bar{F}^L \circ F^{L-1} \circ \dots \circ F^1 \quad (3)$$

# Loss Function

For learning, we need to assess the goodness-of-fit of network.

## Definition 1 (Loss function)

A loss (or cost) function is non-negative function

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}, \quad (y, \nu) \mapsto \ell(y, \nu) \quad (4)$$

such that  $\ell(y, y) = 0$  ( $\forall y \in \mathcal{Y}$ ) and  $\ell(y, \nu) > 0$  ( $\forall \nu \neq y$ ).

- 1 here:  $\mathcal{Y}$ : output space
- 2 general convention:  $y$  is the truth and  $\nu$  predicted

# Loss Function: Examples

## Example 1 (Squared-error)

$$\mathcal{Y} = \mathbb{R}^m, \ell(y, \nu) = \frac{1}{2} \|y - \nu\|_2^2 = \frac{1}{2} \sum_{i=1}^m (y_i - \nu_i)^2 \quad (5)$$

## Example 2 (Classification error)

$$\mathcal{Y} = [1 : m], \ell(y, \nu) = 1 - \delta_{y\nu} \quad (6)$$

with Kronecker delta:

$$\delta_{ab} = \begin{cases} 1, & \text{if } a = b \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

## Definition 2 (Expected Risk)

Assume inputs and outputs are governed by a distribution  $p(x, y)$  over  $\mathcal{X} \times \mathcal{Y}$ ,  $\mathcal{X} \subseteq \mathbb{R}^n$ ,  $\mathcal{Y} \subseteq \mathbb{R}^m$ . The expected risk of  $F$  is given by

$$\mathcal{R}^*(F) = E_{x,y}[\ell(y, F(x))] \quad (8)$$

- 1 as  $p$  is generally unknown, we cannot evaluate  $\mathcal{R}^*$  directly, but it serves as a point of reference in learning theory
- 2  $\mathcal{R}^*$  is a **functional** (mapping functions to scalars)
- 3 parameterized functions  $\{F_\theta : \theta \in \Theta\} \Rightarrow \mathcal{R}^*(\theta) \triangleq \mathcal{R}^*(F_\theta)$

## Definition 3 (Empirical Risk)

Assume we have a random sample of  $N$  input-output pairs,

$$\mathcal{S}_N \triangleq \left\{ (x_i, y_i) \stackrel{i.i.d.}{\sim} p : 1, \dots, N \right\}. \quad (9)$$

The empirical risk of  $F$  is defined as

$$\mathbb{R}(F, \mathcal{S}_N) = \frac{1}{N} \sum_{i=1}^N \ell(y_i, F(x_i)) \quad (10)$$

- 1 a.k.a. training risk = expected risk under the empirical distribution induced by the sample  $\mathcal{S}_N$ .

# Empirical Risk Minimization

For a family  $\mathcal{F} = \{F_\theta : \theta \in \Theta\}$  (e.f. neural network) and training data  $\mathcal{S}_N$ : find function with lowest empirical risk.

## Definition 4 (Empirical risk minimization)

The empirical risk minimizer is defined as

$$\hat{F}(\mathcal{S}_N) \in \arg \min_{F \in \mathcal{F}} \mathcal{R}(F, \mathcal{S}_N) \quad (11)$$

with the corresponding parameters  $\hat{\theta}(\mathcal{S}_N)$ .

- 1 one may also add a regularizer  $\Omega(F)$  or  $\Omega(\theta)$  to the risk (more on that later)
- 2 finding  $\hat{F} \in \mathcal{F}$  amounts to solving an **optimization** problem

# Probability Distributions as Outputs

It is often constructive to think of functions  $F$  as mappings from inputs to distribution  $\mathcal{P}(\mathcal{Y})$  over outputs  $y \in \mathcal{Y}$ .

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad x \mapsto \nu, \quad \nu \xrightarrow{\text{fixed}} p(y, \nu) \in \mathcal{P}(\mathcal{Y}), \quad y \sim p(\cdot, \nu) \quad (12)$$

Each  $F$  effectively defines a conditional probability distribution (or conditional probability density function) via

$$p(y|x, F) = p(y, \nu = F(x)) \quad (13)$$



# Example: Multivariate Normal Distribution

## Example 3 (mean of a normal distribution)

$$p(y|x, F) = \left[ \frac{1}{\sqrt{2\pi\gamma}} \right]^m e^{\left[ -\frac{1}{2\gamma^2} \|y - F(x)\|^2 \right]} \quad (14)$$

so that

$$-\log p(y|x, F) = mC(\gamma) + \frac{1}{2\gamma^2} \|y - F(x)\|^2 \quad (15)$$

which is equivalent to the squared error loss.

- ①  $F(x) = \nu$  and  $y$  live in same space ( $\mathbb{R}^m$ )

## Definition 5 (Generalized linear model (simplified))

A generalized linear model over  $y \in \mathcal{Y} \subseteq \mathbb{R}$  takes the form

$$E[y|x] = \sigma(w^T x). \quad (16)$$

where  $\sigma$  is invertible and  $\sigma^{-1}$  is called the **link function**.

- ① can be extended to also predict variances or dispersions
- ② can be extended to multidimensional outputs

# Example: Logistic Regression

## Example 4 (Logistic regression)

$\mathcal{Y} = \{0, 1\}$ ,  $\mathcal{P} = [0, 1]$ ,  $\sigma = \frac{1}{1+e^{-x}}$ , then:

$$E[y|x] = p(1|x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} \quad (17)$$

*Link function: logit*

$$\sigma^{-1}(t) = \log\left(\frac{t}{1-t}\right), \quad t \in (0, 1) \quad (18)$$

# Example: Multinomial Logistic Regression

## Example 5

$\mathcal{Y} = [1 : m]$ ,  $\mathcal{P}(\mathcal{Y})$  can be represented via soft-max

$$p(y|x) = \frac{e^{z_y}}{\sum_{i=1}^m e^{z_i}}, \quad z \triangleq w_i^T x, \quad i = 1, \dots, m \quad (19)$$

- 1 over-parametrized model: set  $w_1 = 0$ , s.t.  $z_1 = 0$  (w.l.o.g)
- 2 generalizes (binary) logistic regression

# Generalized Linear Units

In neural networks:

- **non-linear** functions replace linear functions
- output layer units implement **inverse link function**

## Example 6 (Normal model)

*Linear output layer*

$$E[y|x] = \bar{F}(x) = W^L \left( F^{L-1} \circ \dots \circ F^1 \right) (x) + b^L \quad (20)$$

## Example 7 (Logistic model)

*Sigmoid output layer*

$$E[y|x] = \sigma(\bar{F}(x)) \quad (21)$$

Use conditional probability distribution to define generalized loss between target value  $y \in \mathcal{Y}$  and a distribution over  $\mathcal{Y}$ .

## Definition 6 (Negative log-loss)

Canonical way of defining a generalized loss functions: negative of a log-likelihood function

$$\ell(y, \theta, x) = -\log p(y|x, \theta) \quad (22)$$

- 1 non-linearity of output layer is "absorbed" in loss function
- 2 i.e.  $\ell$  depends on  $\bar{F}$
- 3 provides a "template" for generalized loss/risk functions

# Cross-Entropy Loss

Let us look at the (implied) risk function for the logistic function

## Definition 7 (Cross-entropy Loss)

Use shorthand  $z \triangleq \bar{F}(x) \in \mathbb{R}$  then the cross entropy loss over a binary response variable  $y \in \{0, 1\}$  is defined as

$$\begin{aligned} -\log p(y|z) &= -\log \sigma((2y - 1)z) \\ &= \zeta((1 - 2y)z) \end{aligned} \quad (23)$$

where  $\zeta = \log(1 + e^{(\cdot)})$  is the **soft-plus** function.

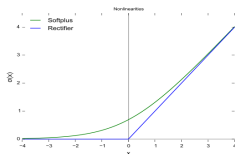


Figure 1: rectifier and softplus functions

# Multinomial Log-Likelihood

## Definition 8 (Multinomial cross-entropy loss)

Assume multinomial response variable  $y \in [1 : m]$ . Use shorthand:

$$z \triangleq \bar{F}(x) \in \mathbb{R}^m \quad (24)$$

then with the soft-max activation function

$$\begin{aligned} \ell(y, \bar{F}(x)) &= -\log p(y|\bar{F}(x)) = -\log \left[ \frac{e^{z_y}}{\sum_{i=1}^m e^{z_i}} \right] \\ &= -z_y + \underbrace{\log \sum_{i=1}^m e^{z_i}}_{\text{log-partition}} = \log \left[ 1 + \sum_{i \neq y} e^{(z_i - z_y)} \right] \end{aligned} \quad (25)$$



# Gradient Descent

Learning in neural networks = gradient-based optimization (with very few exceptions).

## Definition 9 (Gradient)

**Gradient** of objective with regard to parameters  $\theta$

$$\nabla_{\theta} = \left( \frac{\partial \mathcal{R}}{\partial \theta_1}, \dots, \frac{\partial \mathcal{R}}{\partial \theta_d} \right)^T \quad (26)$$

## Definition 10 (Steepest descent and stochastic gradient decent)

**Steepest descent** and **stochastic gradient decent**

$$\theta(t+1) \leftarrow \theta(t) - \eta \nabla_{\theta} \mathcal{R}(\mathcal{S}) \quad (27)$$

- 1 here  $t = 0, 1, 2, \dots$  is an iteration index
- 2  $\mathcal{S}$  = all training data  $\Rightarrow$  steepest descent
- 3  $\mathcal{S}$  = mini batch of data  $\Rightarrow$  SGD

# Gradient Computation via Backpropagation

Computational challenge: how to compute  $\nabla_{\theta} \mathcal{R}$ ?

Exploit compositional structure of network = **backpropagation**

Basic steps:

- 1 perform a forward pass (for given training input  $x$ ) to compute activations for all units
- 2 compute gradient of  $\mathcal{R}$  w.r.t. output layer activations (for given target  $y$ )
- 3 iteratively propagate activation gradient information from outputs to inputs
- 4 compute local gradients of activations w.r.t weights

# Backpropagation in Plain English

- 1 How do changes in the output layer activities change the objective?
  - depends on choice of objective
- 2 How does the activity of a parent unit influence the activity of each of its child units (in DAG)?
  - layer structure  $\Rightarrow$  concurrently between subsequent layers
- 3 Propagate influence information through reverse DAG
  - details are implied by chain rule of differentiation
- 4 What is the effect of a change of an incoming weight on the activity of a unit?
  - can only change activities (given  $x$ ) by modifying weights

Compositional of functions  $\Rightarrow$  use of **chain rule**

## Proposition 1 (Chain Rule)

$$(f \circ g)' = (f' \circ g) \cdot g' \quad (28)$$

*or equivalently with formal variables*

$$\frac{d(f \circ g)}{dx} \Big|_{x=x_0} = \frac{df}{dz} \Big|_{z=g(x_0)} \cdot \frac{dg}{dx} \Big|_{x=x_0} \quad (29)$$

# Jacobi Matrix

Vector-valued function (map)  $F : \mathcal{R}^n \rightarrow \mathcal{R}^m$ : each component function has gradient  $\nabla F_i \in \mathcal{R}^n$ ,  $i \in [1 : m]$

## Definition 11 (Jacobi matrix)

$$J_F \triangleq \begin{pmatrix} \nabla^T F_1 \\ \nabla^T F_2 \\ \vdots \\ \nabla^T F_m \end{pmatrix} = \begin{pmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} & \cdots & \frac{\partial F_1}{\partial x_n} \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} & \cdots & \frac{\partial F_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_m}{\partial x_1} & \frac{\partial F_m}{\partial x_2} & \cdots & \frac{\partial F_m}{\partial x_n} \end{pmatrix} \in \mathbb{R}^{m \times n} \quad (30)$$

derivative of outputs with regard to inputs, i.e.  $(J_F)_{ij} = \frac{\partial F_i}{\partial x_j}$ .

# Jacobian Matrix Chain Rule

Vector-valued functions  $G : \mathbb{R}^n \rightarrow \mathbb{R}^q$ ,  $H : \mathbb{R}^q \rightarrow \mathbb{R}^m$ ,  $F \triangleq H \circ G$ .  
Componentwise rule

$$\frac{\partial F_i}{\partial x_j} \Big|_{x=x_0} = \frac{\partial (H \circ G)_i}{\partial x_j} \Big|_{x=x_0} = \sum_{k=1}^q \frac{\partial H_i}{\partial z_k} \Big|_{z=G(x_0)} \cdot \frac{\partial G_k}{\partial x_j} \Big|_{x=x_0} \quad (31)$$

Lemma 1 (Jacobi matrix chain rule)

$$J_{H \circ G} \Big|_{x=x_0} = J_H \Big|_{z=G(x_0)} \cdot J_G \Big|_{x=x_0} \quad (32)$$

# Function Composition

Special case: composition of a map with a function

$$G : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad h : \mathbb{R}^m \rightarrow \mathbb{R}, \quad h \circ G : \mathbb{R}^n \rightarrow \mathbb{R} \quad (33)$$

$x \in \mathbb{R}^n$ , use more intuitive variable notation

$$x \xrightarrow{G} y \xrightarrow{h} z \in \mathbb{R} \quad (34)$$

Then

$$\nabla_x^T z = \nabla_y^T z \cdot J_G, \quad \frac{\partial z}{\partial x_i} = \sum_j \frac{\partial y_j}{\partial x_i} \frac{\partial z}{\partial y_j} \quad (35)$$

# Warning: Notation!

We have a lot of indices!

- index of a layer: put as a **superscript**
- index of a dimension of a vector: put as a **subscript**
- shorthand for layer activations

$$x^l \triangleq (F^l \circ \dots \circ F^1)(x) \in \mathbb{R}^{m_l} \quad (36)$$

$x_i^l \in \mathbb{R}$  : *activation of  $i$ -th unit in layer  $l$*

- index of a data point, omitted where possible, rectangular brackets  
( $x[i], y[i]$ )



# Deep Function Compositions

Composition of multiple maps with a final cost function

$$\begin{aligned} F &= F^L \circ \dots \circ F^1 : \mathbb{R}^n \rightarrow \mathbb{R}^m \\ x &= x^0 \xrightarrow{F^1} x^1 \xrightarrow{F^2} x^2 \mapsto \dots \xrightarrow{F^L} x^L = \nu \mapsto \ell(y, \nu) \end{aligned} \quad (37)$$

## Proposition 2 (Activity Backpropagation)

$$e^L \triangleq \nabla_{\nu}^T \mathcal{R}, \quad e^l \triangleq \nabla_{x^l}^T \mathcal{R} = e^L \cdot J_{F^L} \cdots J_{F^{l+1}} = e^{l+1} \cdot J_{F^{l+1}} \quad (38)$$

Compute **activity gradients** is backward order via successive multiplication with Jacobians. Backpropagation of error terms  $e^l$ .

Linear network in reversed direction with "activities"  $e^l$ .

# Jacobian Matrix: Ridge Functions

How does a Jacobian matrix for a ridge function look like?

$$x^l = F^l(x^{l-1}) = \sigma(W^l x^{l-1} + b^l) \quad (39)$$

Hence (assuming differentiability of  $\sigma$ ):

$$\frac{\partial x_i^l}{\partial x_j^{l-1}} = \sigma'(\langle w_i^l, x^{l-1} \rangle + b_i^l) w_{ij}^l \triangleq \tilde{w}_{ij}^l \quad (40)$$

and thus simply

$$J_{F^l} = \tilde{W}^l \quad (41)$$

① for ReLU  $\tilde{w}_{ij}^l \in \{0, w_{ij}^l\} \Rightarrow \tilde{W}^l = \text{sparsified matrix}$

# Loss Function (Negative) Gradients

Quadratic loss

$$-\nabla_{\nu} \ell(y, \nu) = -\nabla_{\nu} \frac{1}{2} \|y - \nu\|^2 = y - \nu \quad (42)$$

Multivariate logistic loss

$$\begin{aligned} -\frac{\partial \ell(y, \nu)}{\partial z_y} &= \frac{\partial}{\partial z_y} \left[ z_{y^*} - \log \sum_i e^{z_i} \right] \\ &= \delta_{yy^*} - \frac{e^{z_y}}{\sum_i e^{z_i}} \\ &= \delta_{yy^*} - p(y|x) \end{aligned} \quad (43)$$

# From Activations to Weights

How can we get from gradients w.r.t. activations to gradients w.r.t. weights? Easily!

Need to apply chain rule one more time- locally:

$$\frac{\partial \ell}{\partial w_{ij}^l} = \frac{\partial \ell}{\partial x_i^l} \cdot \frac{\partial x_i^l}{\partial w_{ij}^l} = \underbrace{\frac{\partial \ell}{\partial x_i^l}}_{\text{backprop}} \cdot \underbrace{\sigma' \left( \left\langle w_i^l, x^{l-1} \right\rangle + b_i^l \right)}_{\text{sensitivity of } i\text{-th unit}} \cdot \underbrace{x_j^{l-1}}_{j\text{-th unit activity}} \quad (44)$$

$$\frac{\partial \ell}{\partial b_i^l} = \frac{\partial \ell}{\partial x_i^l} \cdot \frac{\partial x_i^l}{\partial b_i^l} = \frac{\partial \ell}{\partial x_i^l} \cdot \sigma' \left( \left\langle w_i^l, x^{l-1} \right\rangle + b_i^l \right) \cdot 1$$

- each weight/bias influences exactly one unit
- can "reshape" gradient into matrix/tensor form

# Specialized Programming Languages: Theano

Symbolic representation of mathematical expressions.

Access to full computation graph (stability, optimization).

Symbolic differentiation.

[Bergstra 2015]



J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley and Y. Bengio (2010)

Theano: A CPU and GPU Math Compiler in Python

*9th Annual Python In Science Conference (SciPy 2010)*

Thank you all of you! –Yao