

A Dualistic View of Activations in Deep Neural Networks

Ehsan Amid

Joint work with: Rohan Anil, Manfred Warmuth, Abel L. Peirson, Tomer Koren, Yatong Chen,
Christopher Fifty, Vlad Fineberg



Overview

- **Convexity and Duality**
- **A Dualistic View of Activations in DNNs**
- **Application**
 - LocoProp: Local Loss Optimization [[AISTATS22](#)]
- **Conclusions and Future Directions**

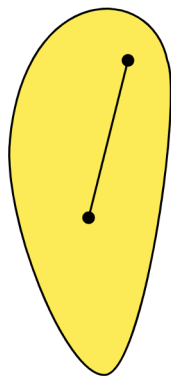


Convexity and Duality

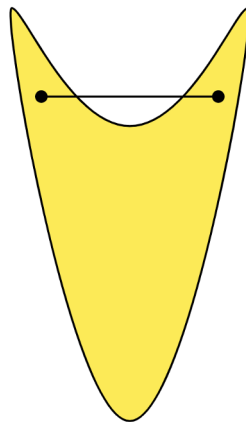
Convex Sets, Convex Functions, Legendre Dual, Bregman Divergence

Convex Sets and Convex Functions

A **set** is **convex** if all the points along the line connecting two points inside the set are also belong to the set



convex



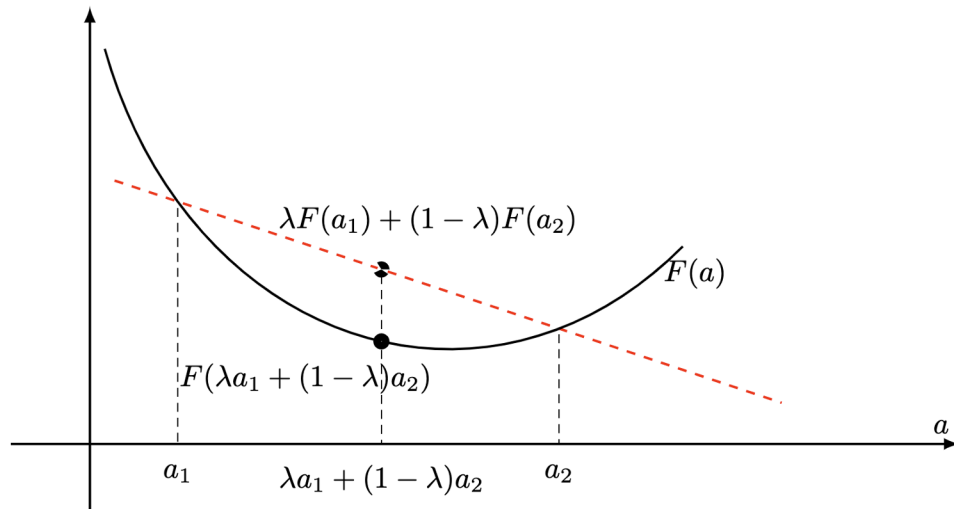
non-convex

Source code: <https://tex.stackexchange.com/questions/639115/how-to-draw-convex-strictly-convex-and-nonconvex-sets>

Convex Sets and Convex Functions

A function is **convex** if it has a convex domain and

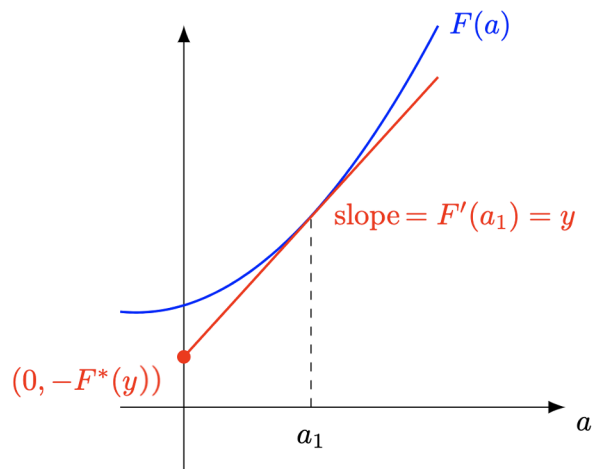
$$F(\lambda a_1 + (1 - \lambda) a_2) \leq \lambda F(a_1) + (1 - \lambda) F(a_2), \quad a_1, a_2 \in \text{Dom}(F), \quad \lambda \in [0, 1]$$



Legendre Dual

Given a (convex) function, the **Legendre dual** convex function is given by

$$F^*(y) = \sup_{\tilde{a} \in \text{Dom}(F)} \{ \langle \tilde{a}, y \rangle - F(\tilde{a}) \}$$



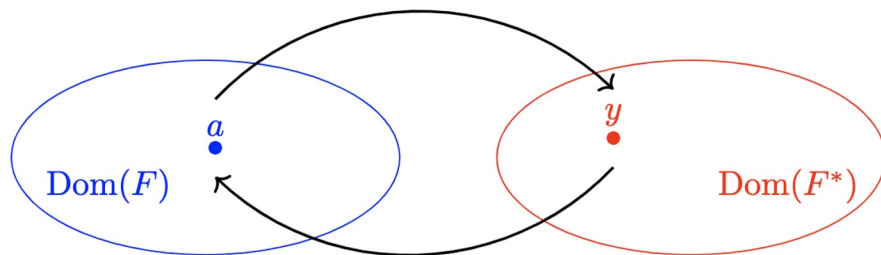
Legendre Dual

When F is strictly convex, we have the **dually coupled variables**:

$$a = \arg \sup_{\tilde{a} \in \text{Dom}(F)} \{ \langle \tilde{a}, y \rangle - F(\tilde{a}) \}$$

(link function)

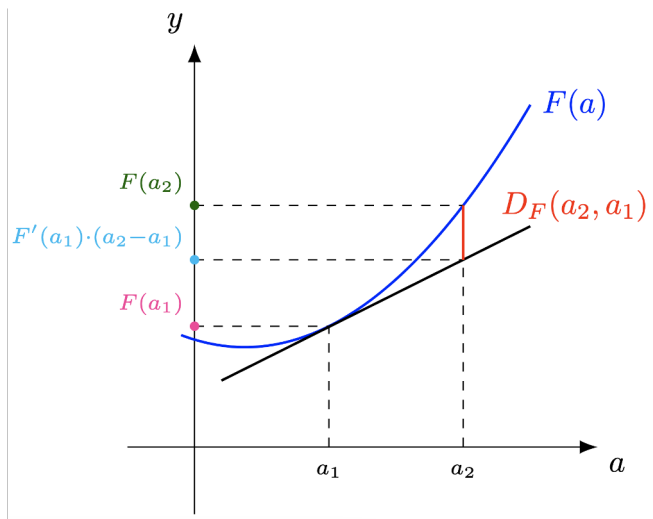
$$f \doteq \nabla F$$



$$f^* \doteq \nabla F^* = f^{-1}$$

(inverse link function)

Bregman Divergence: Definition



$$D_F(a_2, a_1) = \underbrace{F(a_2) - F(a_1)}_{\text{growth of the convex function}} - \underbrace{\nabla F(a_1)^\top (a_2 - a_1)}_{\text{growth of the linear function}}$$

Bregman Divergence: Properties

Convexity: always in the left argument (not necessarily the right)

Non-negativity: $D_F(\hat{\mathbf{a}}, \mathbf{a}) \geq 0$ and $D_F(\hat{\mathbf{a}}, \mathbf{a}) = 0$ iff $\hat{\mathbf{a}} = \mathbf{a}$

Gradient: $\nabla_{\hat{\mathbf{a}}} D_F(\hat{\mathbf{a}}, \mathbf{a}) = \nabla F(\hat{\mathbf{a}}) - \nabla F(\mathbf{a})$

Many well-known cases:

Squared Euclidean: $D_F(\hat{\mathbf{a}}, \mathbf{a}) = \frac{1}{2} \|\hat{\mathbf{a}} - \mathbf{a}\|^2$ (with $F(\mathbf{a}) = \frac{1}{2} \|\mathbf{a}\|^2$)

Relative Entropy: $D_F(\hat{\mathbf{a}}, \mathbf{a}) = \sum_i \left\{ \hat{a}_i \log \frac{\hat{a}_i}{a_i} - \hat{a}_i + a_i \right\}$ (with $F(\mathbf{a}) = \sum_i \{ a_i \log a_i - a_i \}$)

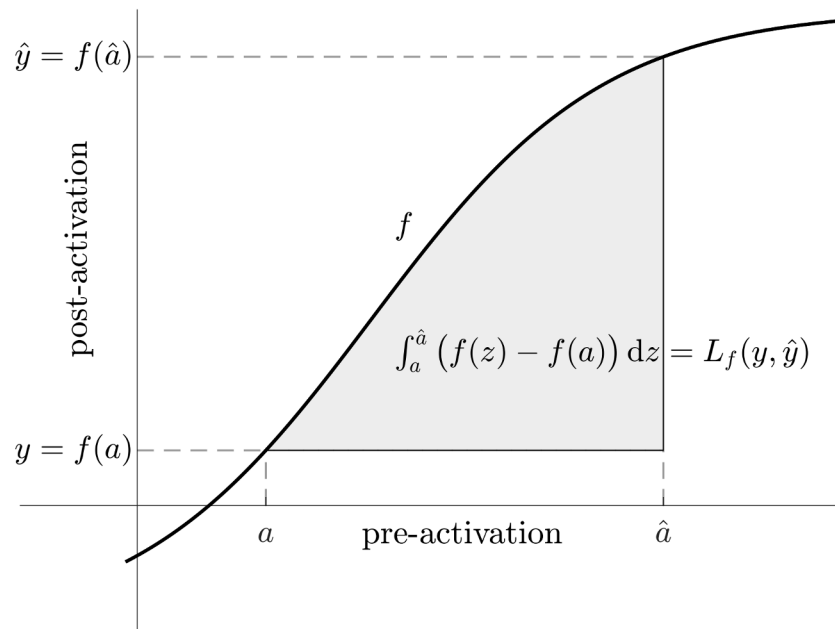
Matching Loss: Definition

Definition: integral under the curve

$$L_f(\mathbf{y}, \hat{\mathbf{y}}(\hat{\mathbf{a}})) \doteq \int_{\mathbf{a}}^{\hat{\mathbf{a}}} (f(\mathbf{z}) - f(\mathbf{a}))^\top d\mathbf{z}$$

... is a Bregman divergence!

$$\begin{aligned} L_f(\mathbf{y}, \hat{\mathbf{y}}(\hat{\mathbf{a}})) &= \left(F(\mathbf{z}) - f(\mathbf{a})^\top \mathbf{z} \right) \Big|_{\mathbf{a}}^{\hat{\mathbf{a}}} \\ &= F(\hat{\mathbf{a}}) - F(\mathbf{a}) - f(\mathbf{a})^\top (\hat{\mathbf{a}} - \mathbf{a}) \doteq D_F(\hat{\mathbf{a}}, \mathbf{a}) \end{aligned}$$



Matching Loss Example: Softmax \rightarrow KL Divergence

Link function:

$$\mathbf{y} = f(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\sum_i \exp(a_i)}$$

Integral function:

$$F(\mathbf{a}) = \log \sum_i \exp(a_i)$$

Matching loss (in terms of **post-activations**):

$$D_F(\hat{\mathbf{a}}, \mathbf{a}) = D_{F^*}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_i y_i \log \frac{y_i}{\hat{y}_i}$$



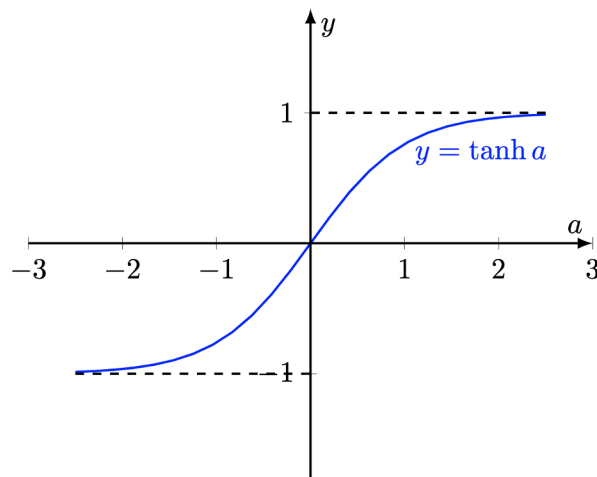
A Dualistic View of Activations in DNNs

The Role of Activation Functions

Which link functions are valid?

For f to be a “valid” link function (i.e., gradient of a strictly convex function):

- Needs to be **strictly increasing** in 1-D
- And cyclically strictly monotone in higher dimension

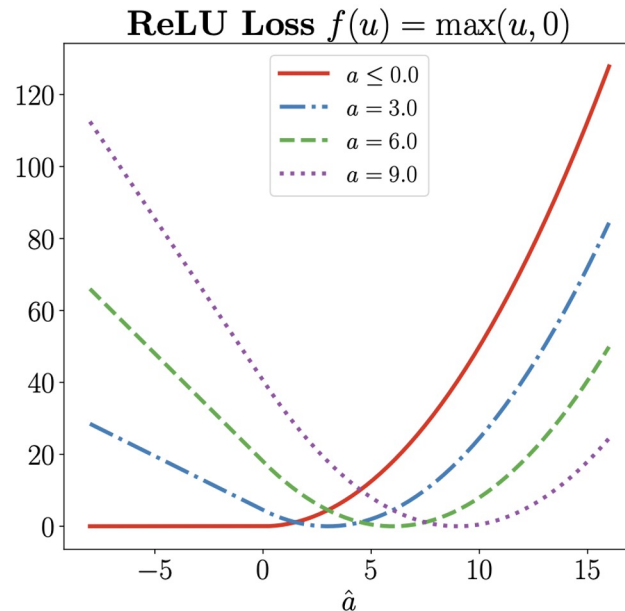
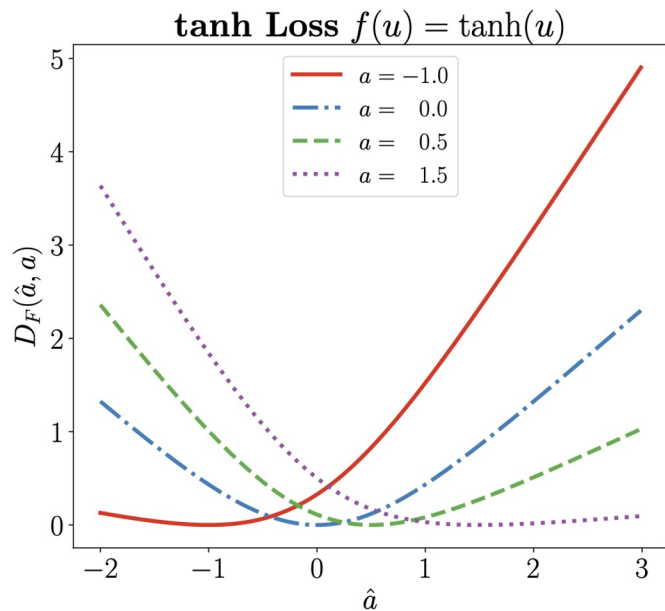


Many **activation functions** in modern neural networks (Linear, Leaky ReLU, tanh, sigmoid, softmax*, etc.) satisfy this property!

Matching Losses of Common Activation Functions

NAME	TRANSFER FUNCTION $f(\mathbf{a})$	CONVEX INTEGRAL FUNCTION $F(\mathbf{a})$	NOTE
STEP FUNCTION	$1/2 (1 + \text{sign}(\mathbf{a}))$	$\sum_i \max(a_i, 0)$	–
LINEAR	\mathbf{a}	$1/2 \ \mathbf{a}\ ^2$	–
(LEAKY) RELU	$\max(\mathbf{a}, 0) - \beta \max(-\mathbf{a}, 0)$	$1/2 \sum_i a_i (\max(a_i, 0) - \beta \max(-a_i, 0))$	$\beta \geq 0$
SIGMOID	$(1 + \exp(-\mathbf{a}))^{-1}$	$\sum_i (a_i + \log(1 + \exp(-a_i)))$	–
SOFTMAX	$\exp(\mathbf{a}) / \sum_i \exp(a_i)$	$\log \sum_i \exp(a_i)$	–
HYPERBOLIC TAN	$\tanh(\mathbf{a})$	$\sum_i \log \cosh(a_i)$	–
ARC TAN	$\arctan(\mathbf{a})$	$\sum_i (a_i \arctan(a_i) - \log \sqrt{1 + a_i^2})$	–
SOFTPLUS	$\log(1 + \exp(\mathbf{a}))$	$-\sum_i \text{Li}_2(-\exp(a_i))$	$\text{Li}_2 := \text{SPENCE'S FUNC.}$
ELU	$[f(\mathbf{a})]_i = \begin{cases} a_i & a_i \geq 0 \\ \beta(\exp a_i - 1) & \text{OTHERWISE} \end{cases}$	$\sum_i (a_i^2/2 \mathbb{I}(a_i \geq 0) + \beta(\exp a_i - a_i - 1) \mathbb{I}(a_i < 0))$	$\beta \geq 0$

Matching Losses of Common Activation Functions





Application

LocoProp: Local Loss Optimization

Local Loss Optimization

LocoProp conceives neural networks as a **modular composition** of layers

Local Loss Optimization

Pre $\hat{\mathbf{a}}_m = \mathbf{W}_m \hat{\mathbf{y}}_{m-1}$ and **post** $\hat{\mathbf{y}}_m = f_m(\hat{\mathbf{a}}_m)$ activations in layer m

$$\mathbf{W}^{\text{new}} = \arg \min_{\tilde{\mathbf{W}}} \left\{ \underbrace{D(\tilde{\mathbf{W}} \hat{\mathbf{y}}_{m-1}, \mathbf{a}_m)}_{\text{loss}} + \underbrace{R(\tilde{\mathbf{W}}, \mathbf{W})}_{\text{regularizer}} \right\}$$

target \mathbf{a}_m (or $\mathbf{y}_m = f_m(\mathbf{a}_m)$)

The Case of Squared Loss

Local regularized squared loss:

$$\underbrace{1/2 \|\widetilde{\mathbf{W}} \hat{\mathbf{y}}_{m-1} - \mathbf{a}_m\|^2}_{\text{loss to the target}} + \underbrace{1/2\eta \|\widetilde{\mathbf{W}} - \mathbf{W}_m\|^2}_{\text{keep the weight close}} \quad \text{where} \quad \underbrace{\mathbf{a}_m = \hat{\mathbf{a}}_m - \gamma \nabla_{\hat{\mathbf{a}}_m} L(\mathbf{y}, \hat{\mathbf{y}})}_{\text{GD w.r.t. the final loss}}$$

Solution: fixed point equation

$$\mathbf{W}_m^{\text{new}} = \mathbf{W}_m - \eta \left(\mathbf{W}_m^{\text{new}} \hat{\mathbf{y}}_{m-1} - \mathbf{a}_m \right) \hat{\mathbf{y}}_{m-1}^{\top}$$

The Case of Squared Loss

How about a single iteration?

$$\begin{aligned} \mathbf{W}_m^{\text{new}} &\approx \mathbf{W}_m - \eta (\mathbf{W}_m \hat{\mathbf{y}}_{m-1} - \mathbf{a}_m) \hat{\mathbf{y}}_{m-1}^\top \\ &= \mathbf{W}_m - \eta (\mathbf{W}_m \hat{\mathbf{y}}_{m-1} \\ &\quad - (\mathbf{W}_m \hat{\mathbf{y}}_{m-1} - \gamma \nabla_{\hat{\mathbf{a}}_m} L(\mathbf{y}, \hat{\mathbf{y}}))) \hat{\mathbf{y}}_{m-1}^\top \\ &= \mathbf{W}_m - \eta \gamma \nabla_{\hat{\mathbf{a}}_m} L(\mathbf{y}, \hat{\mathbf{y}}) \hat{\mathbf{y}}_{m-1}^\top \\ &= \mathbf{W}_m - \eta_e \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial \hat{\mathbf{a}}_m} \frac{\partial \hat{\mathbf{a}}_m}{\partial \mathbf{W}_m}. \quad (\text{BackProp}) \end{aligned}$$

This is just a single step of **BackProp**!

The Case of Squared Loss

$$\mathbf{W}_m^{\text{new}} = \mathbf{W}_m - \eta (\mathbf{W}_m^{\text{new}} \hat{\mathbf{y}}_{m-1} - \mathbf{a}_m) \hat{\mathbf{y}}_{m-1}^\top$$

There is a closed-form solution:

$$\mathbf{W}_m^{\text{new}} = \mathbf{W}_m - \eta_e \nabla_{\mathbf{W}_m} L(\mathbf{y}, \hat{\mathbf{y}}) (\mathbf{I} + \eta \hat{\mathbf{y}}_{m-1} \hat{\mathbf{y}}_{m-1}^\top)^{-1}$$

gradient descent preconditioner matrix

A **preconditioned** gradient descent!

Using Matching Losses

Replace the squared loss with the matching loss:

$$\boxed{1/2 \|\widetilde{\mathbf{W}} \hat{\mathbf{y}}_{m-1} - \mathbf{a}_m\|^2} + 1/2\eta \|\widetilde{\mathbf{W}} - \mathbf{W}_m\|^2 \quad \text{where} \quad \mathbf{a}_m = \hat{\mathbf{a}}_m - \gamma \nabla_{\hat{\mathbf{a}}_m} L(\mathbf{y}, \hat{\mathbf{y}})$$



$$D_{F_m^*}(\boxed{\mathbf{y}_m}, \boxed{f_m(\widetilde{\mathbf{W}} \hat{\mathbf{y}}_{m-1})}) + 1/2\eta \|\widetilde{\mathbf{W}} - \mathbf{W}_m\|^2 \quad \text{where} \quad \boxed{\mathbf{y}_m = \hat{\mathbf{y}}_m - \gamma \nabla_{\hat{\mathbf{a}}_m} L(\mathbf{y}, \hat{\mathbf{y}})}$$

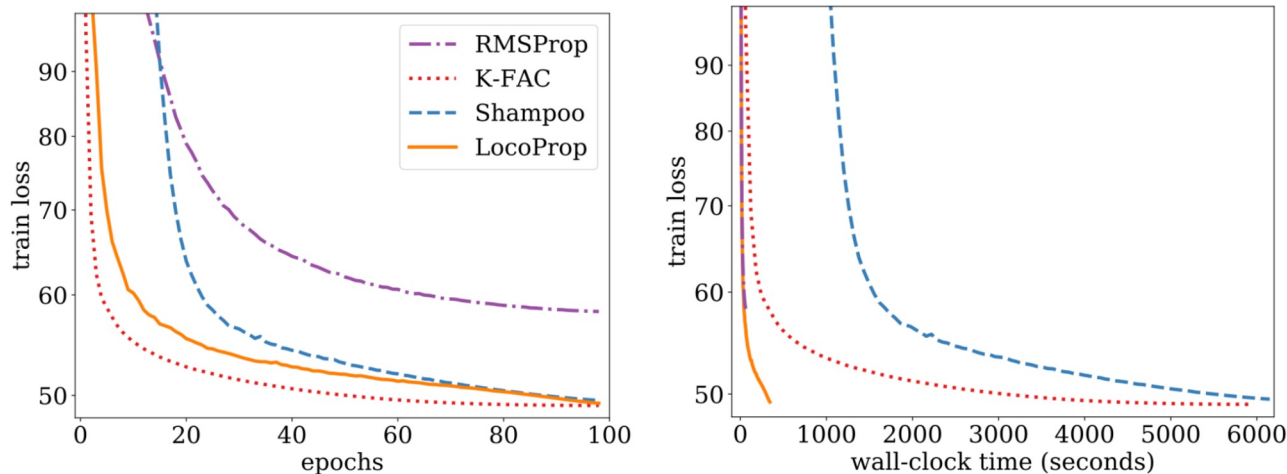
target post-activation

MD w.r.t. the final
loss

Still convex in the **weights**, and yields a preconditioned update

Experiments

Results on the deep autoencoder benchmark



LocoProp performs competitive to **second-order** methods



Conclusions and Future Directions

Conclusions and Future Directions

Further applications:

- Layerwise Fisher approximation via local sampling [[NerIPS22-HITY](#)]
- Bregman knowledge distillation [[TMLR23](#)]
- Robust bi-tempered loss [[NeurIPS19](#)]

Future directions:

- No-sample layerwise Fisher approximation
- Extension to non-monotonic activations
- Low-rank compression of the weights
- Learning losses and activation functions
- ...