

# Robust Bi-Tempered Logistic Loss Based on Bregman Divergences

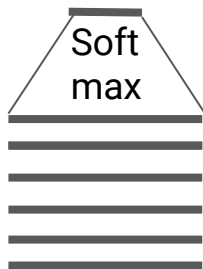
Ehsan Amid  
Google Brain & UC Santa Cruz

{eamid, manfred, rohananil, tkoren}@google.com

Nov 12, 2019 - Google NYC

# Logistic Loss

**Logistic Loss = relative entropy (KL) divergence + softmax probabilities**



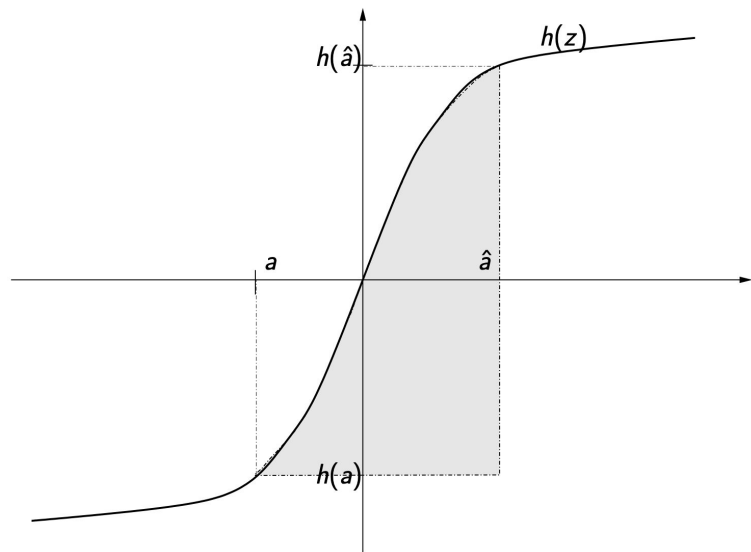
- The most commonly used loss for classification in NN
  - Always convex in the activations/weights of the last layer
  - However, anyway non-convex in weights of lower layers
  - **We replace last layer by non-convex loss that makes NN robust to outliers**

# Matching Loss

[HKW95]

Softmax  $\hat{y}_i = \frac{\exp(\hat{a}_i)}{\sum_{j=1}^k \exp(\hat{a}_j)}$

Matching loss



Logistic loss = area under sigmoid

$$\Delta_{\log}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_i y_i \log \frac{y_i}{\hat{y}_i}$$

In general

$$\begin{aligned} \Delta_h(\hat{\mathbf{a}}, \mathbf{a}) &= \int_{\mathbf{a}}^{\hat{\mathbf{a}}} (h(\mathbf{z}) - h(\mathbf{a})) d\mathbf{z} \\ &= \Delta_{h^{-1}}(\underbrace{h(\mathbf{a})}_{\mathbf{y}}, \underbrace{h(\hat{\mathbf{a}})}_{\hat{\mathbf{y}}}) \end{aligned}$$

# The simplicity of the matching loss

- $$\Delta_h(\hat{\mathbf{a}}, \mathbf{a}) = \int_{\mathbf{a}}^{\hat{\mathbf{a}}} (h(\mathbf{z}) - h(\mathbf{a})) d\mathbf{z}$$

- **Convex for any increasing transfer function**

$$\frac{\partial}{\partial \hat{\mathbf{a}}} \Delta_h(\hat{\mathbf{a}}, \mathbf{a}) = h(\hat{\mathbf{a}}) - h(\mathbf{a}) = \underbrace{\hat{\mathbf{y}} - \mathbf{y}}_{\text{delta rule}}$$

- Examples

$h = \text{id}$

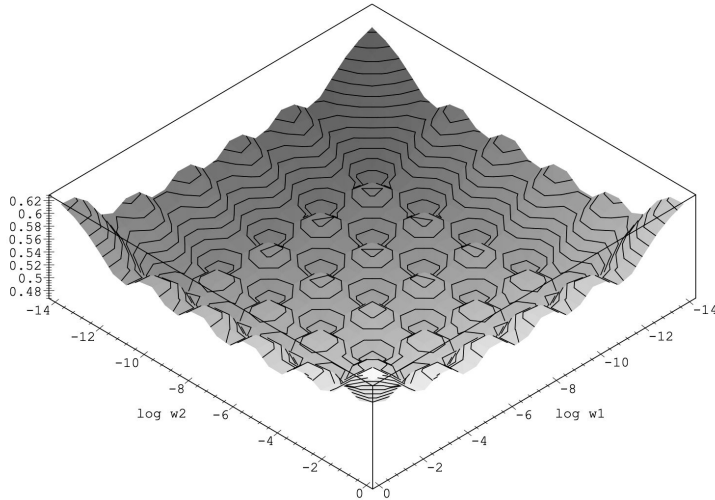
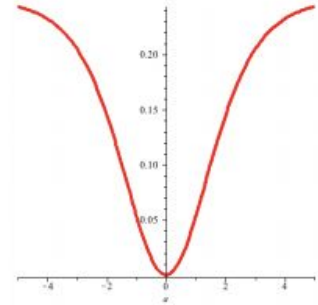
$$\Delta_{\text{id}}(\mathbf{y}, \hat{\mathbf{y}}) = 1/2 \sum_i (y_i - \hat{y}_i)^2$$

$h = \text{softmax}$

$$\Delta_{\log}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_i y_i \log \frac{y_i}{\hat{y}_i} \quad (\text{logistic loss})$$

# Canonical mismatched case: sigmoid with square loss

$$\frac{\partial}{\partial \hat{a}} (\sigma(\hat{a}) - y)^2 = (\sigma(\hat{a}) - y) \underbrace{\sigma'(\hat{a})}_{\sigma(\hat{a})(1-\sigma(\hat{a}))}$$

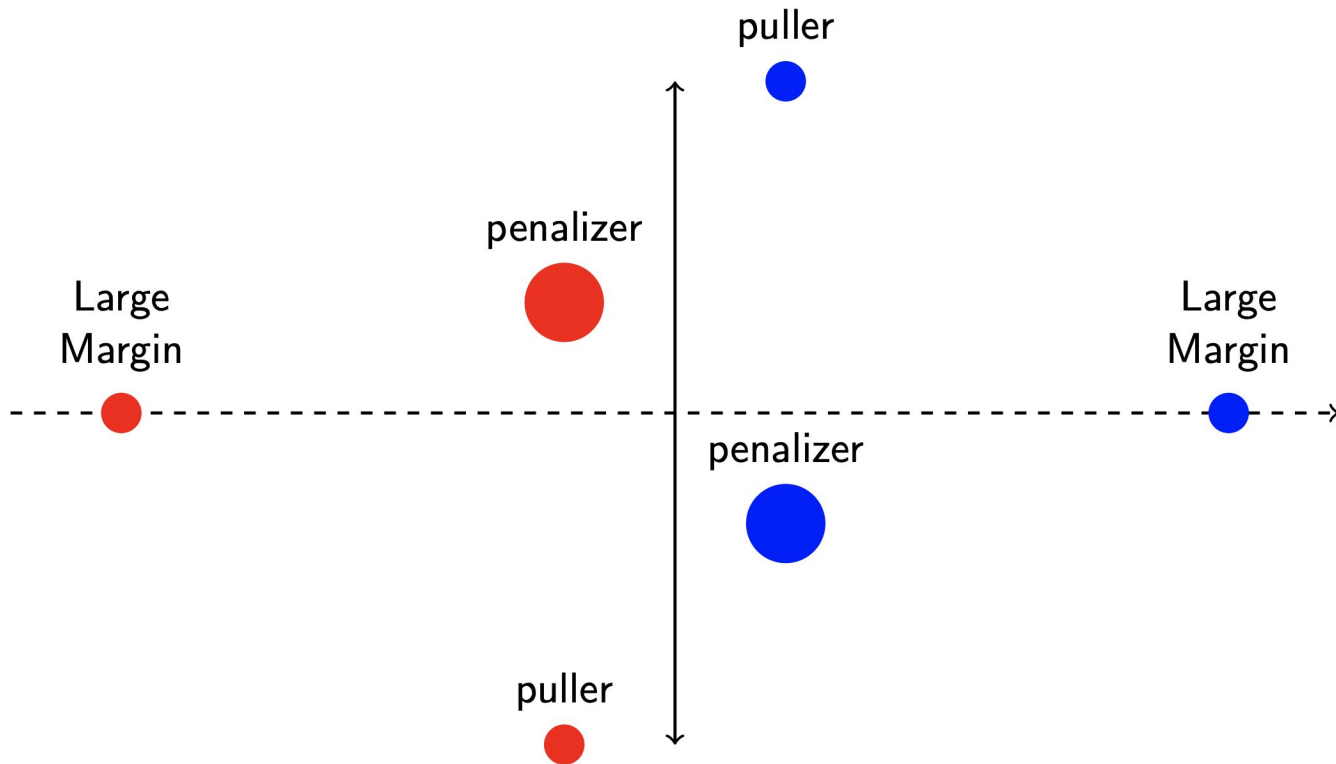


Can lead to  
exponent. many  
minimas

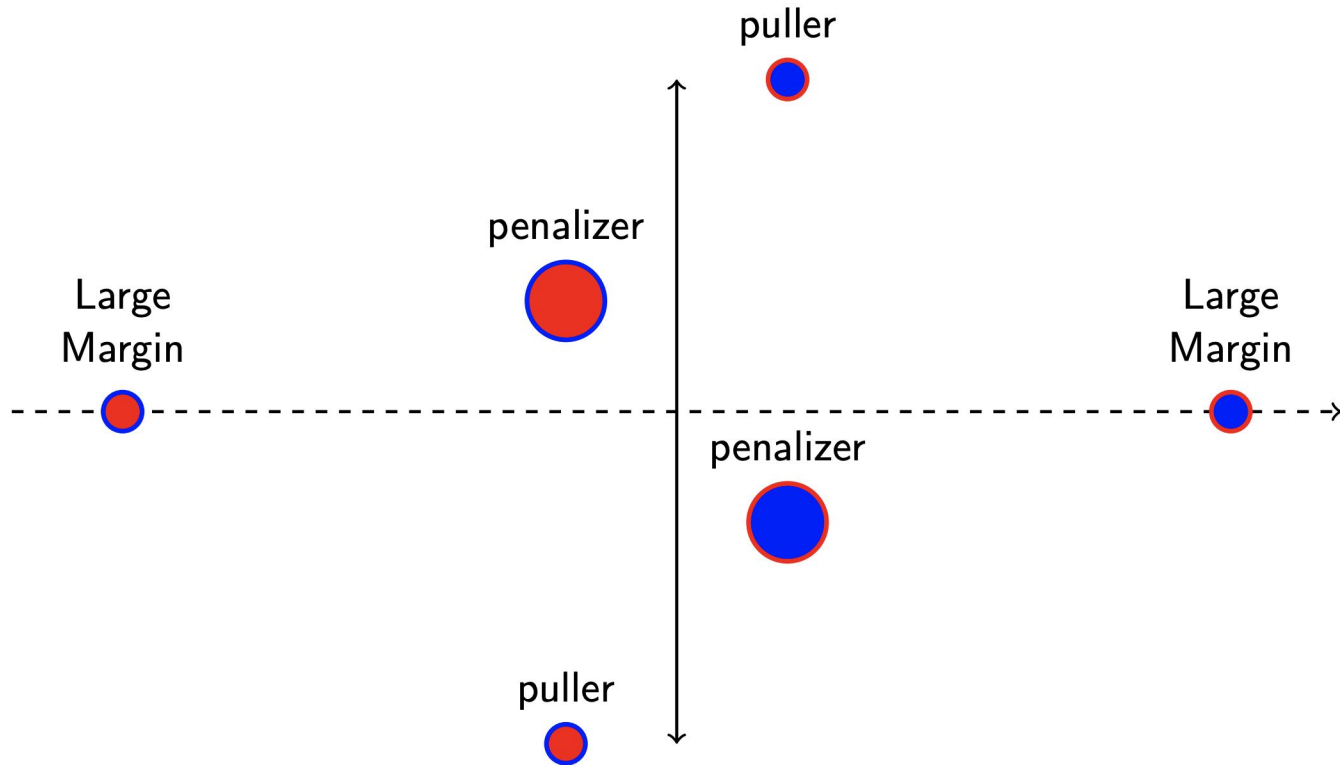
[AHW95]

# Key example justifying non-convexity

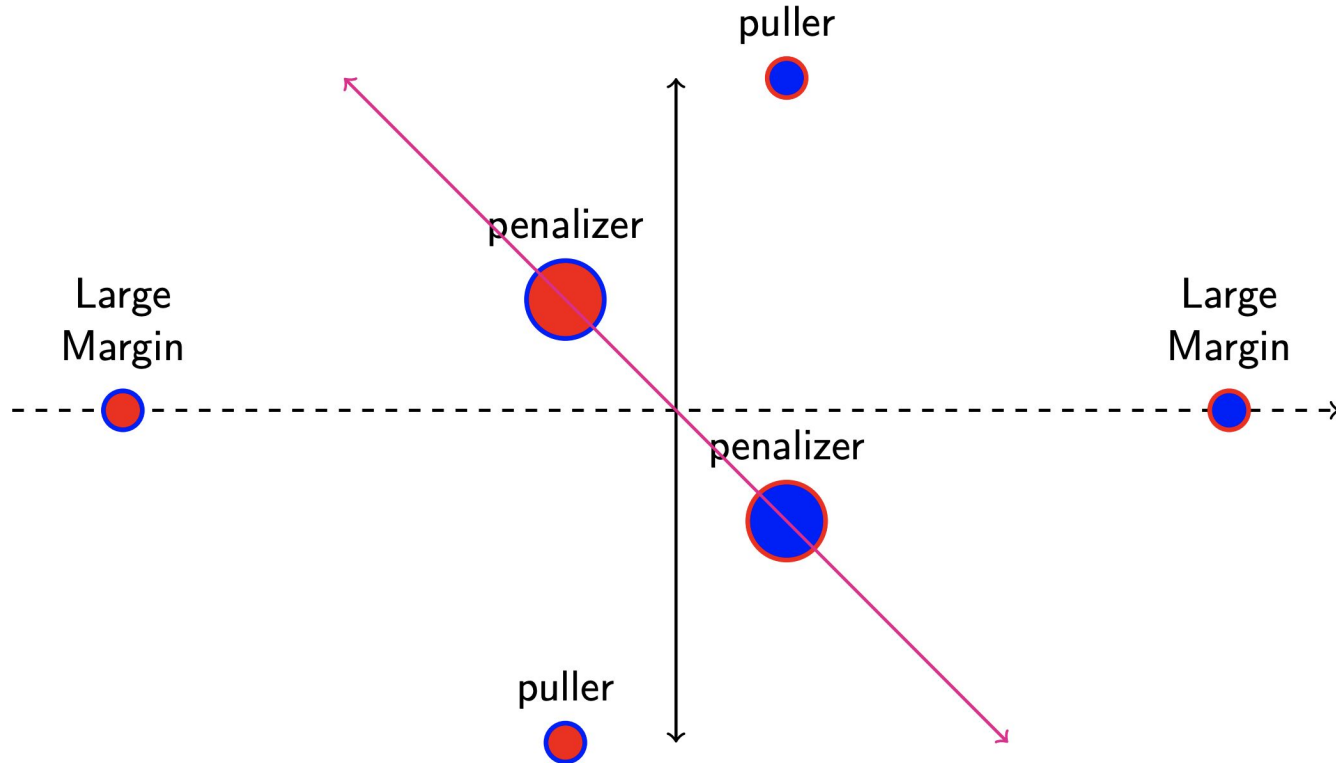
[LS08]



# 10% label noise



# Solution determined by the outliers





# Logistic Loss as a Matching Loss

Start: Logistic Loss = relative entropy divergence + softmax probabilities

- Introduce temperatures into links, i.e.  $\log_{t_1}$  and  $\exp_{t_2}$
- When the 2 temperatures are **equal**, we again obtain a **convex** loss
- By increasing the temperature in the exponential, loss becomes non-convex
- Tuning the two temperatures will be crucial

# Tempered Logarithm and Exponential

[N02]

Generalization of log and exp functions endowed with a **temperature**  $t \geq 0$

$$\log_t(x) := \frac{1}{1-t} (x^{1-t} - 1)$$

Bounded by  $-1/(1-t)$  at 0 for  $0 \leq t < 1$

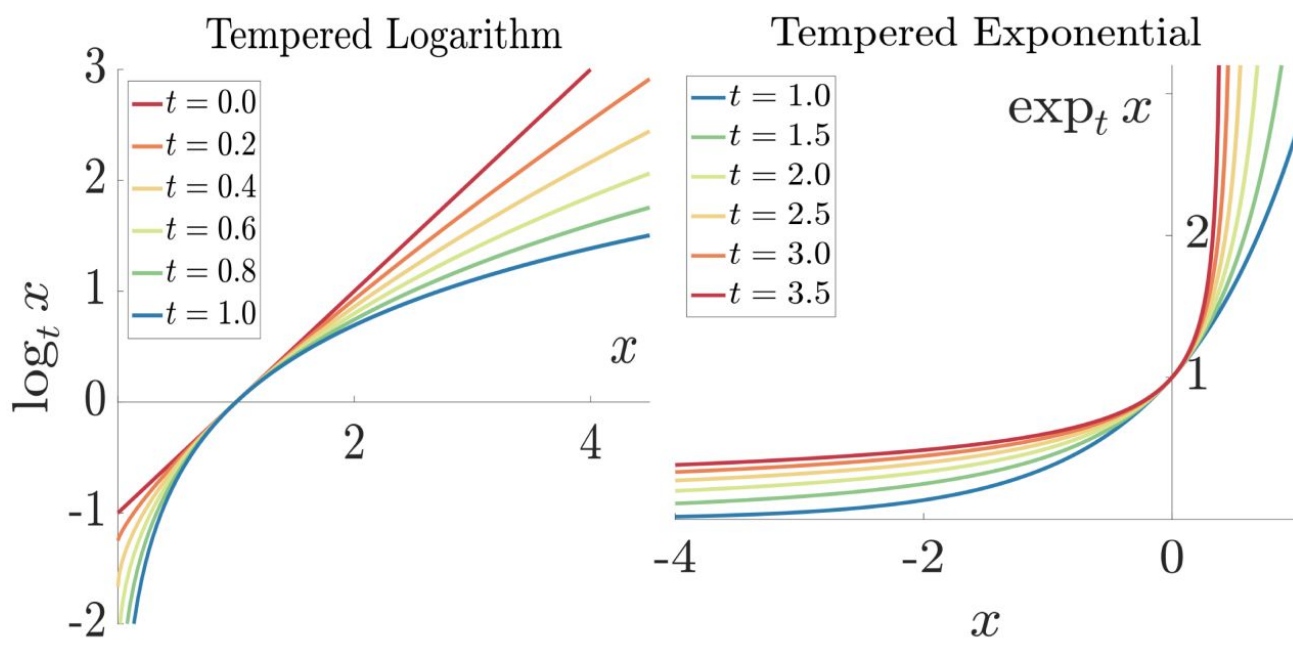
$$\exp_t(x) := [1 + (1-t)x]_+^{1/(1-t)}$$

Heavier tail for  $x < 0$  for  $t > 1$

Standard log and exp are recovered at the limit  $t \rightarrow 1$

# Tempered Logarithm and Exponential

Generalization of log and exp functions endowed with a **temperature**  $t \geq 0$



# Two Drawbacks of Logistic Loss

## 1. Convex losses are not robust to noise

[LS08]

- Convex losses increase unboundedly
- A single bad example can dominate the cumulative loss
- Extreme examples can cause large gradients
- Non-convex (bending down) losses have been shown to perform significantly better

## 2. Softmax probabilities have exponentially decaying tail

- The margin becomes small for mislabeled examples near the boundary
- Heavy-tailed alternatives yield better margins and improved results [DV,10]

# 1. Replacing the Relative Entropy Divergence

Tempered relative entropy divergence ( $0 \leq t_1 < 1$ ):

$$\sum_{i=1}^k \left( y_i (\log_{t_1} y_i - \log_{t_1} \hat{y}_i) - \frac{1}{2-t_1} (y_i^{2-t_1} - \hat{y}_i^{2-t_1}) \right) \stackrel{\text{if } \mathbf{y} \text{ one-hot}}{=} -\log_{t_1} \hat{y}_c - \frac{1}{2-t_1} \left( 1 - \sum_{i=1}^k \hat{y}_i^{2-t_1} \right)$$

**bounded by  $1/(1-t_1)$**

where  $c = \operatorname{argmax}_i y_i$  is the index of the one-hot class

## 2. Replacing the Softmax Probabilities

$\mathbf{z}$  : vector of inputs to softmax layer

$\mathbf{w}_i$  : trainable weight vector for class  $i$

$\mathbf{y}$  : target vector

Softmax:

$$\hat{y}_i = \frac{\exp(\hat{a}_i)}{\sum_{j=1}^k \exp(\hat{a}_j)} = \exp\left(\hat{a}_i - \log \sum_{j=1}^k \exp(\hat{a}_j)\right), \text{ for linear activation } \hat{a}_i = \mathbf{w}_i \cdot \mathbf{z} \text{ for class } i.$$

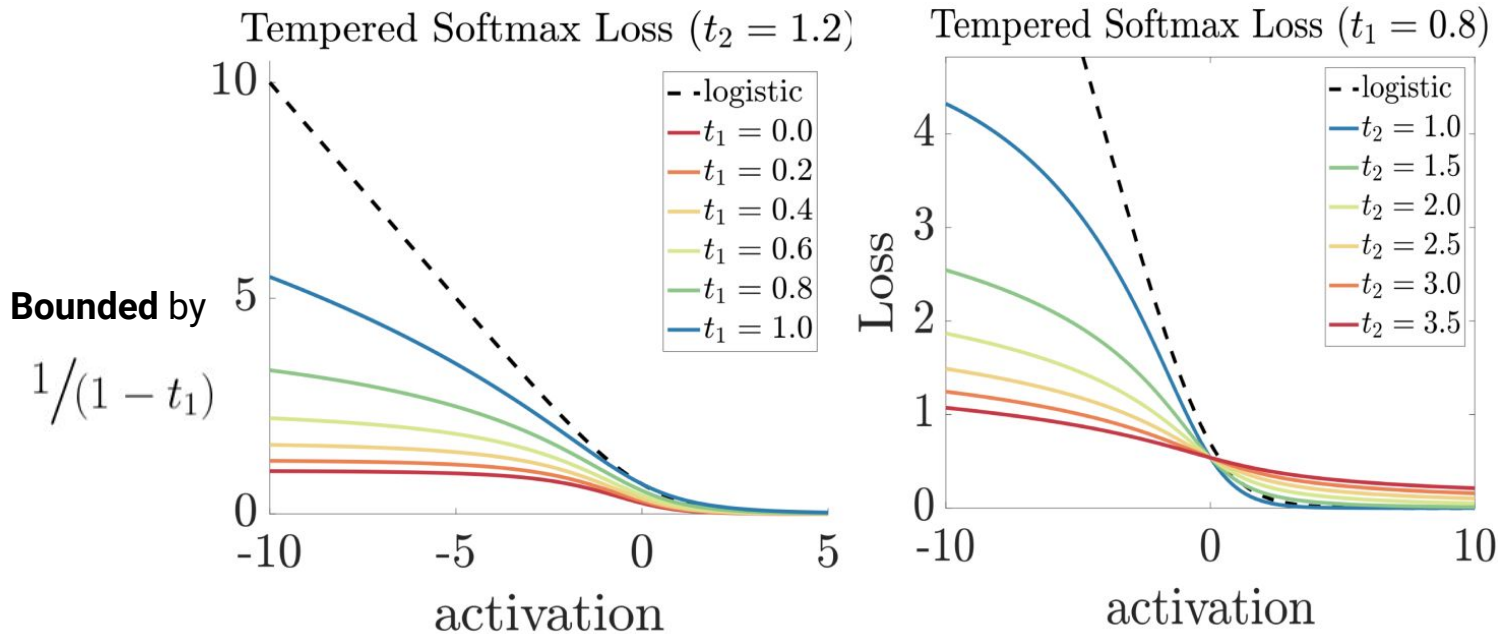
Tempered Softmax ( $t_2 > 1$ ):

$$\hat{y}_i = \exp_{t_2}\left(\hat{a}_i - \lambda_{t_2}(\hat{\mathbf{a}})\right), \text{ where } \lambda_{t_2}(\hat{\mathbf{a}}) \in \mathbb{R} \text{ is s.t. } \sum_{j=1}^k \exp_{t_2}\left(\hat{a}_j - \lambda_{t_2}(\hat{\mathbf{a}})\right) = 1$$

**Tail-heavy for  $t_2 > 1$ !**

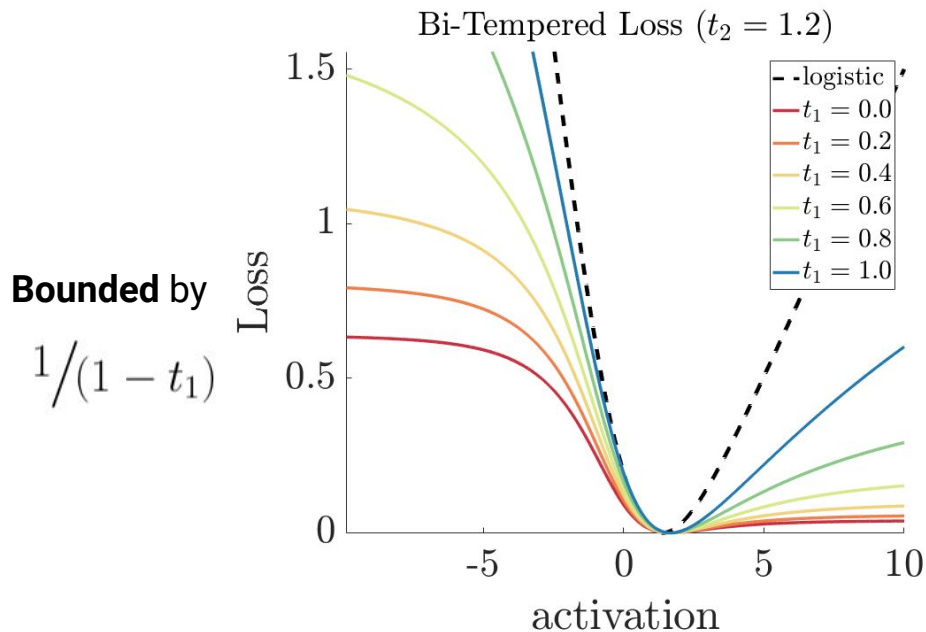
# Examples of the Bi-Tempered Logistic Loss ( $y = 1.0$ )

A binary classification task with true label = 1.0

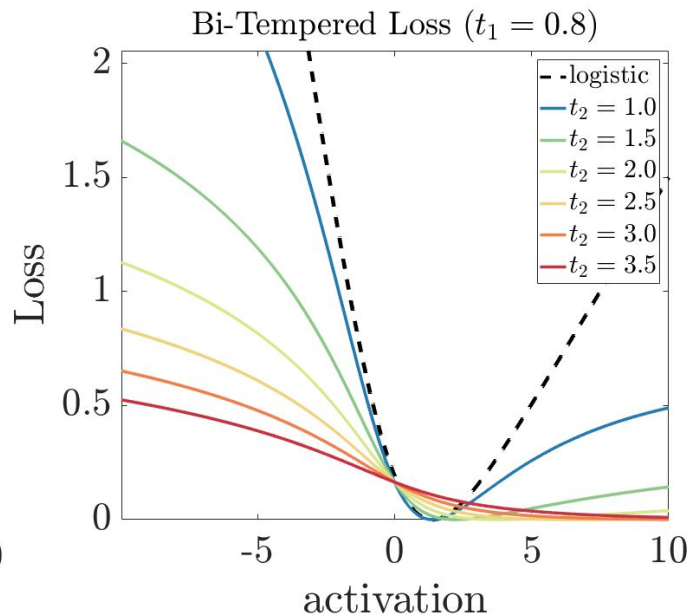


# Examples of the Bi-Tempered Logistic Loss ( $y = 0.8$ )

A binary classification task with true label = 0.8



Heavier-tail requires higher activations to produce the same probability



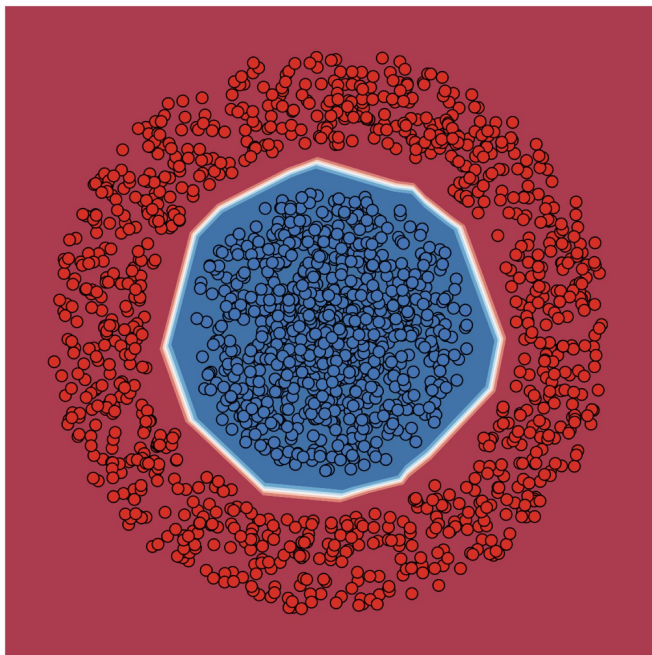


# An Illustration

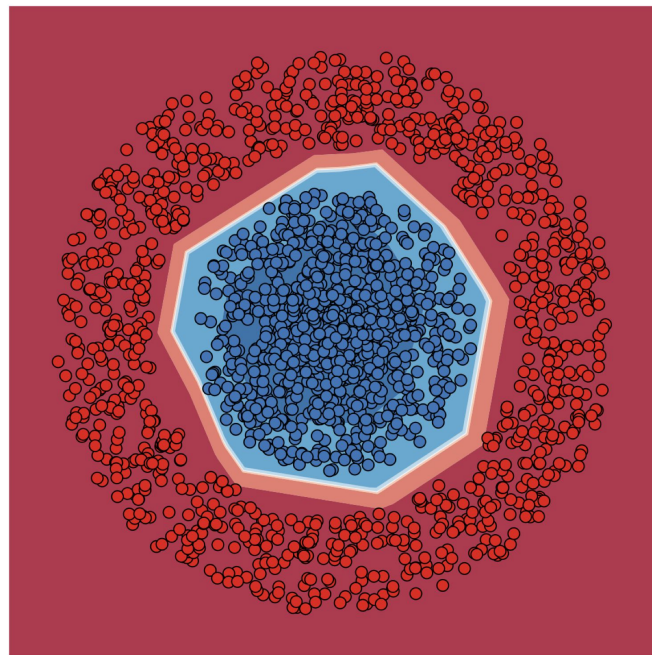
- A small two layer feed-forward neural net on a synthetic binary classification problem in two dimension
  - 10 and 5 units in the first and second layer, respectively
  - Trained using logistic and our bi-tempered logistic loss
  - We add synthetic label noise by flipping the labels
  - Four cases:
    - i. Noise-free
    - ii. Small-margin noise (targeting point near the boundary)
    - iii. Large-margin noise (targeting points far away from the boundary)
    - iv. Random noise (points are selected uniformly at random)

# Noise-free Case

Logistic



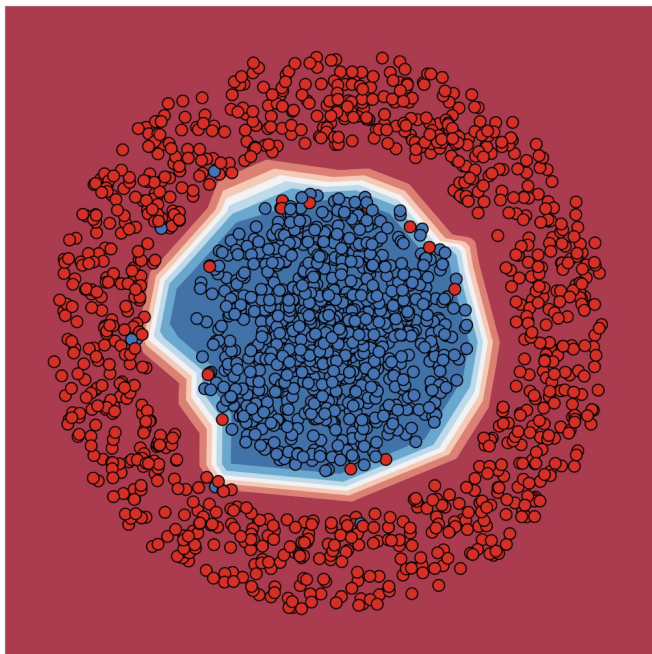
Bi-Tempered (0.2, 4.0)



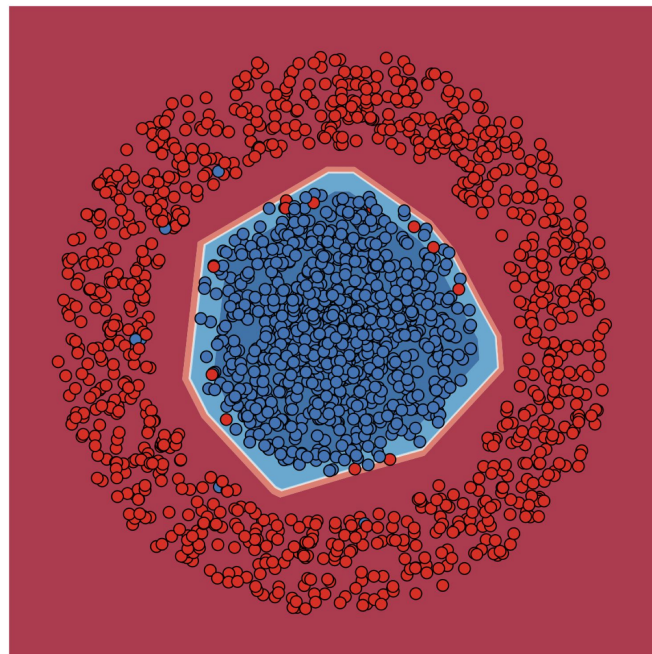
tuned for bounded & tail-heavy

# Small-margin Noise

Logistic



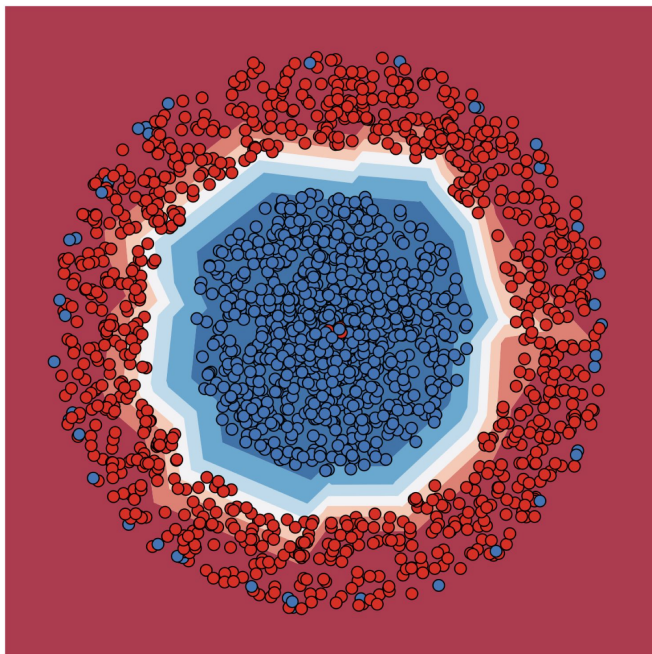
Bi-Tempered (1.0, 4.0)



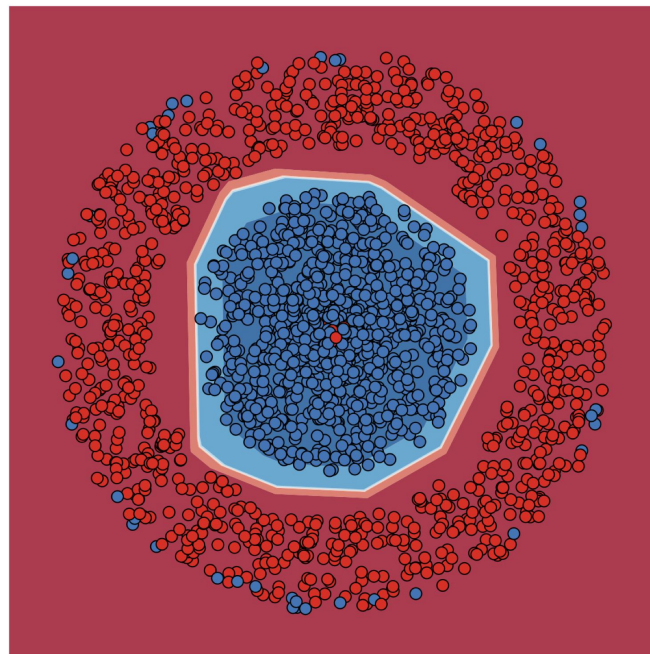
tuned for tail-heavy

# Large-margin Noise

Logistic



Bi-Tempered (0.2, 1.0)



tuned for bounded



# Experiments

# Synthetic Label Noise

**For MNIST:** 2 convolution layers: 32, 64. Followed by 2 FFN of size 1024 and 10, trained for 500 epochs

**For CIFAR-100:** a Resnet56 architecture with SGD + momentum optimizer trained for 50k steps with batch size of 128

We search over range  $[0.5, 1)$  and  $(1, 4]$  for  $t_1$  and  $t_2$ , respectively

# Synthetic Label Noise

Dataset	Loss	Label Noise Level					
		0.0	0.1	0.2	0.3	0.4	0.5
MNIST	Logistic	<b>99.40</b>	98.96	98.70	98.50	97.64	96.13
	Bi-Tempered (0.5, 4.0)	99.24	<b>99.13</b>	<b>99.02</b>	<b>98.62</b>	<b>98.56</b>	<b>97.69</b>
CIFAR-100	Logistic	74.03	69.94	66.39	63.00	53.17	52.96
	Bi-Tempered (0.8, 1.2)	<b>75.30</b>	<b>73.30</b>	<b>70.69</b>	<b>67.45</b>	<b>62.55</b>	<b>57.80</b>

Table 1: Top-1 accuracy on a clean test set for MNIST and CIFAR-100 datasets where a fraction of the training labels are corrupted.



# Large-scale Experiments

**On the Imagenet-2012 dataset** with state-of-the-art Resnet-18 and Resnet-50 models

Trained on a 4x4 CloudTPU-v2 device with a batch size of 4096

180 epochs, SGD + momentum optimizer with staircase learning rate decay schedule

# Imagenet Results

Model	Logistic	Bi-tempered (0.9,1.05)
Resnet18	71.333 $\pm$ 0.069	<b>71.618</b> $\pm$ 0.163
Resnet50	76.332 $\pm$ 0.105	<b>76.748</b> $\pm$ 0.164

Top-1 accuracy

# Theoretical Preliminaries

# Convex Duality and Bregman Divergences

For a continuously-differentiable strictly convex function  $F : \mathcal{D} \rightarrow \mathbb{R}$

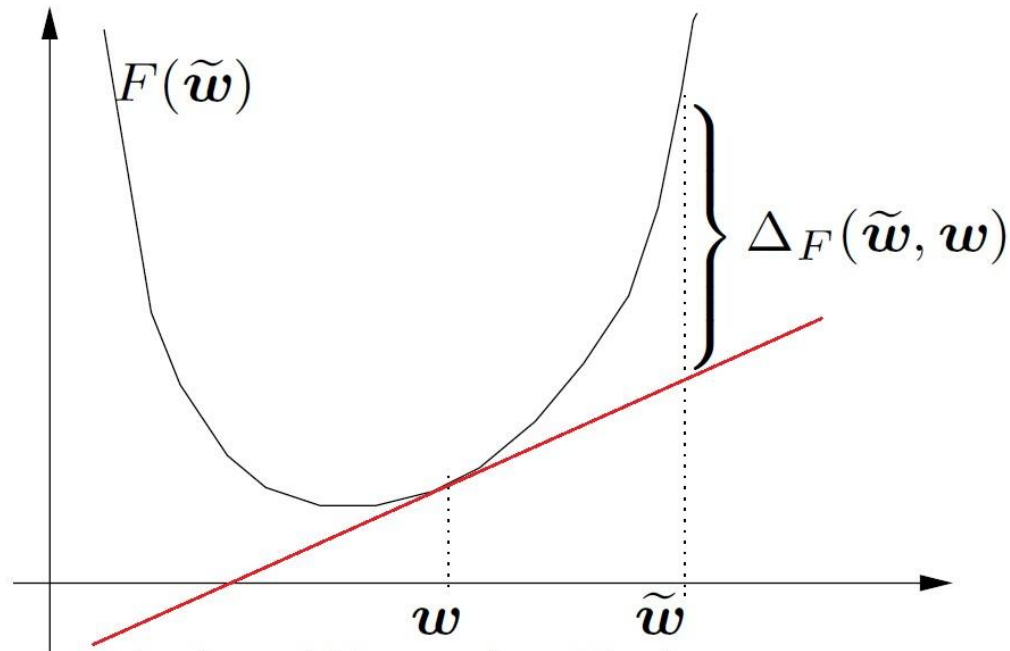
Bregman divergence between  $\mathbf{y}, \hat{\mathbf{y}} \in \mathcal{D}$

$$\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = F(\mathbf{y}) - F(\hat{\mathbf{y}}) - (\mathbf{y} - \hat{\mathbf{y}}) \cdot f(\hat{\mathbf{y}})$$

where  $f(\hat{\mathbf{y}}) := \nabla F(\hat{\mathbf{y}})$  denotes the gradient

# Convex Duality and Bregman Divergences

$$F : \mathcal{D} \rightarrow \mathbb{R}$$



$$\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = F(\mathbf{y}) - F(\hat{\mathbf{y}}) - (\mathbf{y} - \hat{\mathbf{y}}) \cdot f(\hat{\mathbf{y}}), \text{ where } f(\hat{\mathbf{y}}) := \nabla F(\hat{\mathbf{y}})$$

# Properties of Bregman Divergence

- **Convexity:** always in the first argument (not necessarily the second)
- **Non-negativity:**  $\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) \geq 0$  and  $\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = 0$  iff  $\mathbf{y} = \hat{\mathbf{y}}$
- **Gradient:**  $\nabla_{\mathbf{y}} \Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = f(\mathbf{y}) - f(\hat{\mathbf{y}})$
- **Invariance to adding affine functions:**

$$\Delta_{F+A}(\mathbf{y}, \hat{\mathbf{y}}) = \Delta_F(\mathbf{y}, \hat{\mathbf{y}}), \text{ where } A(\mathbf{y}) = \mathbf{b} + \mathbf{c} \cdot \mathbf{y}$$

- **Many well-known cases:**

- Squared Euclidean:  $\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$  (with  $F(\mathbf{y}) = \frac{1}{2} \|\mathbf{y}\|_2^2$ )
- Relative entropy:  $\Delta_F(\mathbf{y}, \hat{\mathbf{y}}) = \sum_i (y_i \log \frac{y_i}{\hat{y}_i} - y_i + \hat{y}_i)$  (with  $F(\mathbf{y}) = \sum_i (y_i \log y_i - y_i)$ )

# Tempered Entropy Function

convex function  $F_t : \mathbb{R}_+^k \rightarrow \mathbb{R}$  with a temperature parameter  $t \geq 0$

$$F_t(\mathbf{y}) = \sum_{i=1}^k \left( y_i \log_t y_i + \frac{1}{2-t} (1 - y_i^{2-t}) \right)$$

Gradient:  $f_t(\mathbf{y}) := \nabla F_t(\mathbf{y}) = \log_t \mathbf{y}$

**Lemma 1.** *The function  $F_t$ , with  $0 \leq t \leq 1$ , is  $B^{-t}$ -strongly convex over the set  $\{\mathbf{y} \in \mathbb{R}_+^k : \|\mathbf{y}\|_{2-t} \leq B\}$  w.r.t. the  $L_{2-t}$ -norm.*

# Tempered Relative Entropy Divergence

The Bregman divergence induced by  $F_t$

$$\begin{aligned}\Delta_{F_t}(\mathbf{y}, \hat{\mathbf{y}}) &= \sum_{i=1}^k \left( y_i \log_t y_i - y_i \log_t \hat{y}_i - \frac{1}{2-t} y_i^{2-t} + \frac{1}{2-t} \hat{y}_i^{2-t} \right) \\ &= \sum_{i=1}^k \left( \frac{1}{(1-t)(2-t)} y_i^{2-t} - \frac{1}{1-t} y_i \hat{y}_i^{1-t} + \frac{1}{2-t} \hat{y}_i^{2-t} \right).\end{aligned}$$

Also known as  $\beta$ -divergence with  $\beta = 2 - t$ .



# Some Special Cases of the Tempered RE

$t$	$F_t(\mathbf{y})$	$\Delta_{F_t}(\mathbf{y}, \hat{\mathbf{y}})$	Name
0	$\frac{1}{2} \ \mathbf{y}\ _2^2$	$\frac{1}{2} \ \mathbf{y} - \hat{\mathbf{y}}\ _2^2$	Euclidean
$\frac{1}{2}$	$\frac{1}{3} \sum_i (4 y_i^{\frac{4}{3}} - 6 y_i + 2)$	$\sum_i (\frac{4}{3} y_i^{\frac{3}{2}} - 2 y_i \sqrt{\hat{y}_i} + \frac{3}{2} \hat{y}_i^{\frac{3}{2}})$	
1	$\sum_i (y_i \log y_i - y_i + 1)$	$\sum_i (y_i \log \frac{y_i}{\hat{y}_i} - y_i + \hat{y}_i)$	KL-divergence
$\frac{3}{2}$	$\sum_i (-4 y_i^{\frac{3}{2}} + 2 y_i + 2)$	$2 \sum_i \frac{(\sqrt{y_i} - \sqrt{\hat{y}_i})^2}{\sqrt{\hat{y}_i}}$	Squared Xi on roots
2	$\sum_i (-\log y_i + y_i)$	$\sum_i (\frac{y_i}{\hat{y}_i} - \log \frac{y_i}{\hat{y}_i} - 1)$	Itakura-Saito
3	$\frac{1}{2} \sum_i (-\frac{1}{y_i} + y_i - 2)$	$\frac{1}{2} \sum_i (\frac{1}{y_i} - \frac{2}{\hat{y}_i} + \frac{y_i}{\hat{y}_i^2})$	Inverse

# Tempered Transfer Function

Using the duality argument

$$\check{F}_t^*(\mathbf{a}) = \sup_{\mathbf{y}' \in S^k} (\mathbf{y}' \cdot \mathbf{a} - F_t(\mathbf{y}')) = \sup_{\mathbf{y}' \in \mathbb{R}_+^k} \inf_{\lambda_t \in \mathbb{R}} (\mathbf{y}' \cdot \mathbf{a} - F_t(\mathbf{y}') + \lambda_t (1 - \sum_{i=1}^k y'_i))$$

The tempered (softmax) transfer function becomes

$$\mathbf{y} = \exp_t(\mathbf{a} - \lambda_t(\mathbf{a}) \mathbf{1}), \quad \text{with } \lambda_t(\mathbf{a}) \text{ s.t. } \sum_{i=1}^k \exp_t(a_i - \lambda_t(\mathbf{a})) = 1$$

# Robust Bi-Tempered Logistic Loss

**Bi-tempered logistic = tempered relative entropy divergence + tempered softmax**

$$\forall 0 \leq t_1 < 1 < t_2: L_{t_1}^{t_2}(\hat{\mathbf{a}} \mid \mathbf{y}) := \Delta_{F_{t_1}}(\mathbf{y}, \exp_{t_2}(\hat{\mathbf{a}} - \lambda_{t_2}(\hat{\mathbf{a}}))), \text{ with } \lambda_t(\hat{\mathbf{a}}) \text{ s.t. } \sum_{i=1}^k \exp_t(a_i - \lambda_t(\mathbf{a})) = 1$$

$0 \leq t_1 < 1$  : controls boundedness of relative entropy

$1 < t_2$  : controls tail-heaviness of softmax

# Two Important Properties

## 1. Properness

- Ensures that we have an unbiased estimator of the expected loss

## 2. Bayes-risk Consistency

- Implies inference can be done via the argmax operation over the activations

# Properness

## Model fitting:

Unknown data distribution:  $P_{\text{UK}}(y \mid \mathbf{x})$

Model distribution:  $P(y \mid \mathbf{x}; \Theta)$  parameterized by  $\Theta$

## Minimize:

$$\mathbb{E}_{P_{\text{UK}}(\mathbf{x})} \left[ \Delta \left( P_{\text{UK}}(y \mid \mathbf{x}), P(y \mid \mathbf{x}; \Theta) \right) \right]$$

We would like to get an unbiased estimator

# Properness of Bi-Tempered Loss

Ignoring the constant terms w.r.t.  $\Theta$

$$\begin{aligned} & \mathbb{E}_{P_{\text{UK}}(\mathbf{x})} \left[ \sum_i \left( -P_{\text{UK}}(i | \mathbf{x}) \log_t P(i | \mathbf{x}; \Theta) + \frac{1}{2-t} P(i | \mathbf{x}; \Theta)^{2-t} \right) \right] \\ & \approx \frac{1}{N} \sum_n \sum_i \left( -P_{\text{UK}}(i | \mathbf{x}_n) \log_t P(i | \mathbf{x}_n; \Theta) + \frac{1}{2-t} P(i | \mathbf{x}_n; \Theta)^{2-t} \right) \\ & \approx \frac{1}{N} \sum_n \left( -\log_t P(y_n | \mathbf{x}_n; \Theta) + \sum_i \frac{1}{2-t} P(i | \mathbf{x}_n; \Theta)^{2-t} \right), \end{aligned}$$

Thus, it is a proper loss

# Tsallis Divergence is not Proper

The approximation using the Tsallis divergence is not proper

[AWS19]

$$\mathbb{E}_{P_{\text{UK}}(\mathbf{x})} \left[ \underbrace{- \sum_i P_{\text{UK}}(i | \mathbf{x}) \log_t \frac{P(i | \mathbf{x}; \Theta)}{P_{\text{UK}}(i | \mathbf{x})}}_{\Delta_t^{\text{Tsallis}}(P_{\text{UK}}(y|\mathbf{x}), P(y|\mathbf{x}; \Theta))} \right] \approx -\frac{1}{N} \sum_n \log_t \frac{P(y_n | \mathbf{x}_n; \Theta)}{P_{\text{UK}}(y_n | \mathbf{x}_n)}$$

But in practice,  $P_{\text{UK}}(y_n | \mathbf{x}_n)$  is unknown and the loss is approximated by

$$-\frac{1}{N} \sum_n \log_t P(y_n | \mathbf{x}_n; \Theta)$$

# Bayes-risk Consistency

The conditional risk of the multiclass loss  $\mathbf{l}(\hat{\mathbf{a}})$  with  $l_i := \ell(\hat{\mathbf{a}} | y = i)$ ,  $i \in [k]$  is defined as

$$R(\boldsymbol{\eta}, \mathbf{l}(\hat{\mathbf{a}})) = \sum_i \eta_i l_i,$$

where  $\eta_i := P_{\text{UK}}(y = i | \mathbf{x})$ .

**Definition 4** (Bayes-risk Consistency). A Bayes-risk consistent loss for multiclass classification is the class of loss functions  $\ell$  for which  $\hat{\mathbf{a}}^*$ , the minimizer of  $R(\boldsymbol{\eta}, \mathbf{l}(\hat{\mathbf{a}}))$ , satisfies

$$\arg \min_i \ell(\hat{\mathbf{a}}^* | y = i) \subseteq \operatorname{argmax}_i \eta_i.$$

**Proposition 2.** *The multiclass bi-tempered logistic loss  $L_{t_1}^{t_2}(\hat{\mathbf{a}} | y)$  is Bayes-risk consistent.*



# Implementation and Future Work

# Implementation

Current version of the paper accepted to NeurIPS 2019:

<https://arxiv.org/pdf/1906.03361.pdf>

Open source TF implementation available at **Google research Github**:

[https://github.com/google-research/google-research/tree/master/bitempered\\_loss](https://github.com/google-research/google-research/tree/master/bitempered_loss)

Replace one line:

~~`softmax_cross_entropy_with_logits(activations, labels)`~~

`bi_tempered_logistic_loss(activation, labels, t1, t2)`

# Future Work

- Better tuning of the temps (possibly dynamic tuning during training)
- Reoptimize all other variables:  
regularization, batch normalization, structure, dropout, ...
- Use Bi-Tempered loss for language models, ad placement, ...
- Can we avoid model blow-ups w. new loss
- Design asymmetric losses
- Long-term: generalization of the matching loss for deep NNs

# References

- [[HKW95](#)] David P. Helmbold, Jyrki Kivinen and Manfred K. Warmuth. Worst-case loss bounds for single neurons. NIPS '95, pp. 309–315.
- [[AHW95](#)] Peter Auer, Mark Herbster, and Manfred K. Warmuth. Exponentially Many Local Minima for Single Neurons. NIPS'95. pp. 315–322.
- [[KW01](#)] Jyrki Kivinen and Manfred K. Warmuth. Relative loss bounds for multidimensional regression problems. Journal of Machine Learning, Vol. 45(3), pp. 301-329.
- [[LS08](#)] Philip M Long and Rocco A Servedio. Random classification noise defeats all convex potential boosters. ICML'08. pp. 608–615.
- [[DV10](#)] Nan Ding and S. V. N. Vishwanathan. t-logistic regression. NIPS'10, pp. 514–522.
- [[AWS19](#)] Ehsan Amid, Manfred K. Warmuth, and Sriram Srinivasan. Two-temperature logistic regression based on the Tsallis divergence. AISTATS'19.