# Learning with Few Labeled Data

Pratik Chaudhari

Electrical and Systems Engineering & Computer and Information Science



University of Pennsylvania

## Menu

#### - An introduction

- Few-shot image classification
- Fundamental limits of representation learning





$$\theta^* = \operatorname*{argmin}_{ heta} \frac{1}{N} \sum_{i=1}^{N} f_i( heta)$$



$$heta^* = \operatorname*{argmin}_{ heta} rac{1}{N} \sum_{i=1}^N f_i( heta)$$

In practice, Stochastic Gradient Descent (SGD) works extremely well

$$egin{aligned} & heta_{t+1} = heta_t - rac{\eta}{ heta} \sum_{k=1}^{ heta} 
abla f_{\omega_k}( heta_t) \ & \omega_k \in \{1, \dots, N\} \end{aligned}$$



$$heta^* = \operatorname*{argmin}_{ heta} rac{1}{N} \sum_{i=1}^N f_i( heta)$$

In practice, Stochastic Gradient Descent (SGD) works extremely well

$$egin{aligned} & heta_{t+1} = heta_t - rac{\eta}{ heta} \sum_{k=1}^{ heta} 
abla f_{\omega_k}( heta_t) \ &\omega_k \in \{1,\ldots,N\} \end{aligned}$$



# Why is SGD so special?

#### SGD finds "wide minima"



#### SGD finds "wide minima"



Local Entropy [Chaudhari et al., ICLR17] is a modified loss function for training deep networks

$$f_{\gamma}(\theta) = -\log(G_{\gamma} * e^{-f(\theta)})$$

#### Parle: parallelizing stochastic gradient descent

Couple MCMC and distributed updates to get state-of-the-art performance [Chaudhari et al., SysML18]



#### Parle: parallelizing stochastic gradient descent

Couple MCMC and distributed updates to get state-of-the-art performance [Chaudhari et al., SysML18]



# Diving deeper with statistical physics

A continuous-time analysis shows that SGD has a non-equilibrium steady-state distribution [Chaudhari & Soatto, ICLR 2018]

# Diving deeper with statistical physics

A continuous-time analysis shows that SGD has a non-equilibrium steady-state distribution [Chaudhari & Soatto, ICLR 2018]

The posterior distribution of SGD maintains a balance between energy and entropy

$$\rho^{\rm ss} = \underset{\rho}{\operatorname{argmin}} \quad \underset{\theta \sim \rho}{\mathbb{E}} [\Phi(\theta)] - \beta^{-1} H(\rho)$$

# Diving deeper with statistical physics

A continuous-time analysis shows that SGD has a non-equilibrium steady-state distribution [Chaudhari & Soatto, ICLR 2018]

The posterior distribution of SGD maintains a balance between energy and entropy

$$\rho^{\rm ss} = \underset{\rho}{\operatorname{argmin}} \quad \underset{\theta \sim \rho}{\mathbb{E}} [\Phi(\theta)] - \beta^{-1} H(\rho)$$

Noise in SGD is highly non-isotropic



#### Knots in our understanding



#### Menu

- An introduction
- Few-shot image classification
- Fundamental limits of representation learning

# Three regimes of image classification

High-shot regime 100–1000 samples/class



# Three regimes of image classification

High-shot regime 100–1000 samples/class

Low-shot regime 10 samples/class





# Three regimes of image classification

High-shot regime 100–1000 samples/class

Low-shot regime 10 samples/class

Extreme low-shot regime (anomaly detection) 1 sample/class







#### **Problem formulation**

Training set consists of labeled samples from lots of "tasks", e.g., classifying cars, cats, dogs, planes ...

#### **Problem formulation**

Training set consists of labeled samples from lots of "tasks", e.g., classifying cars, cats, dogs, planes ...

Data from the new task, e.g., classifying strawberries has

- *w* "ways": number of classes,
- *s* "shots": number of labeled samples per class.

#### **Problem formulation**

Training set consists of labeled samples from lots of "tasks", e.g., classifying cars, cats, dogs, planes ...

Data from the new task, e.g., classifying strawberries has

- w "ways": number of classes,
- *s* "shots": number of labeled samples per class.

Few-shot setting considers the case when s is small.

#### A flavor of current few-shot algorithms

#### Prototypical Networks [Snell et al., 2017]

- Collect a meta-training set, this consists of a large number of related tasks
- Train one model on all these tasks to ensure that the clustering of features of this model correctly classifies the task
- If the test task comes from the same distribution as the meta-training tasks, we can use the clustering on the new task to classify new classes



A classifier trained on a dataset  $D_s$  is a function F that classifies data x using

$$\widehat{y} = F(x; D_s).$$

A classifier trained on a dataset  $D_s$  is a function F that classifies data x using

$$\widehat{y} = F(x; D_s).$$

The parameters  $\theta^* = \theta(D_s)$  of the classifier are a statistic of the dataset  $D_s$  obtained after training. Maintaining this statistic avoids having to search over functions F at inference time.

A classifier trained on a dataset  $D_s$  is a function F that classifies data x using

$$\widehat{y} = F(x; D_s).$$

The parameters  $\theta^* = \theta(D_s)$  of the classifier are a statistic of the dataset  $D_s$  obtained after training. Maintaining this statistic avoids having to search over functions F at inference time.

We cannot learn a good (sufficient) statistic using few samples. So we will search over functions at test-time more explicitly

$$\widehat{y} = \underset{y_{N_s+1}}{\operatorname{argmin}} \min_{\theta} \frac{1}{N_s+1} \sum_{i=1}^{N_s+1} -\log p_{\theta}(y_i \mid x_i) + \frac{1}{2\lambda} \|\theta - \theta^*(D_s)\|^2.$$

#### **Transductive Learning**



# A very simple baseline

- 1. Train a large deep network on the meta-training dataset with the standard classification loss
- 2. Initialize a new "classifier head" on top of the logits to handle new classes
- 3. Fine-tune with the few labeled data from the new task
- 4. Perform transductive learning using the unlabeled test data

# A very simple baseline

- 1. Train a large deep network on the meta-training dataset with the standard classification loss
- 2. Initialize a new "classifier head" on top of the logits to handle new classes
- 3. Fine-tune with the few labeled data from the new task
- 4. Perform transductive learning using the unlabeled test data

with a few practical tricks like cosine annealing of step-sizes, mixup regularization, 16-bit training, very heavy data augmentation, and label smoothing cross-entropy

# An example



#### **Results on benchmark datasets**

|                                                              |                                  | Mini-ImageNet                      |                                    | Tiered-ImageNet                    |                                    | CIFAR-FS                           |                                    | FC-100                             |                                    |
|--------------------------------------------------------------|----------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| Algorithm                                                    | Architecture                     | 1-shot (%)                         | 5-shot (%)                         |
| Matching networks (Vinyals et al., 2016)                     | $conv (64)_{\times 4}$           | 46.6                               | 60                                 |                                    |                                    |                                    |                                    |                                    |                                    |
| LSTM meta-learner (Ravi & Larochelle, 2016)                  | $\mathbf{conv}\;(64)_{\times 4}$ | $43.44\pm0.77$                     | $60.60\pm0.71$                     |                                    |                                    |                                    |                                    |                                    |                                    |
| Prototypical Networks (Snell et al., 2017)                   | conv (64)×4                      | $49.42\pm0.78$                     | $68.20\pm0.66$                     |                                    |                                    |                                    |                                    |                                    |                                    |
| MAML (Finn et al., 2017)                                     | $conv (32)_{\times 4}$           | $48.70 \pm 1.84$                   | $63.11\pm0.92$                     |                                    |                                    |                                    |                                    |                                    |                                    |
| R2D2 (Bertinetto et al., 2018)                               | $conv (96^k)_{\times 4}$         | $51.8\pm0.2$                       | $68.4\pm0.2$                       |                                    |                                    | $65.4\pm0.2$                       | $79.4\pm0.2$                       |                                    |                                    |
| TADAM (Oreshkin et al., 2018)                                | ResNet-12                        | $58.5\pm0.3$                       | $76.7\pm0.3$                       |                                    |                                    |                                    |                                    | $40.1\pm0.4$                       | $56.1\pm0.4$                       |
| Transductive Propagation (Liu et al., 2018b)                 | $\mathbf{conv}\;(64)_{\times 4}$ | $55.51\pm0.86$                     | $69.86\pm0.65$                     | $59.91\pm0.94$                     | $73.30 \pm 0.75$                   |                                    |                                    |                                    |                                    |
| Transductive Propagation (Liu et al., 2018b)                 | ResNet-12                        | 59.46                              | 75.64                              |                                    |                                    |                                    |                                    |                                    |                                    |
| MetaOpt SVM (Lee et al., 2019)                               | ResNet-12 *                      | $62.64 \pm 0.61$                   | $\textbf{78.63} \pm \textbf{0.46}$ | $65.99 \pm 0.72$                   | $81.56\pm0.53$                     | $72.0\pm0.7$                       | $84.2\pm0.5$                       | $41.1\pm0.6$                       | $55.5\pm0.6$                       |
| Support-based initialization (train)                         | WRN-28-10                        | $56.17\pm0.64$                     | $73.31\pm0.53$                     | $67.45\pm0.70^\dagger$             | $82.88\pm0.53^\dagger$             | $70.26\pm0.70$                     | $83.82\pm0.49^\dagger$             | $36.82\pm0.51$                     | $49.72\pm0.55$                     |
| Fine-tuning (train)                                          | WRN-28-10                        | $57.73 \pm 0.62$                   | $\textbf{78.17} \pm \textbf{0.49}$ | $66.58 \pm 0.70$                   | $\textbf{85.55} \pm \textbf{0.48}$ | $68.72 \pm 0.67$                   | $\textbf{86.11} \pm \textbf{0.47}$ | $38.25\pm0.52$                     | $\textbf{57.19} \pm \textbf{0.57}$ |
| Transductive fine-tuning (train)                             | WRN-28-10                        | $\textbf{65.73} \pm \textbf{0.68}$ | $\textbf{78.40} \pm \textbf{0.52}$ | $\textbf{73.34} \pm \textbf{0.71}$ | $\textbf{85.50} \pm \textbf{0.50}$ | $\textbf{76.58} \pm \textbf{0.68}$ | $\textbf{85.79} \pm \textbf{0.50}$ | $\textbf{43.16} \pm \textbf{0.59}$ | $\textbf{57.57} \pm \textbf{0.55}$ |
| Activation to Parameter (Qiao et al., 2018)<br>(train + val) | WRN-28-10                        | $59.60\pm0.41$                     | 73.74 ± 0.19                       |                                    |                                    |                                    |                                    |                                    |                                    |
| LEO (Rusu et al., 2018) (train + val)                        | WRN-28-10                        | $61.76\pm0.08$                     | $77.59 \pm 0.12$                   | $66.33 \pm 0.05$                   | $81.44\pm0.09$                     |                                    |                                    |                                    |                                    |
| MetaOpt SVM (Lee et al., 2019) (train + val)                 | ResNet-12 *                      | $64.09\pm0.62$                     | $\textbf{80.00} \pm \textbf{0.45}$ | $65.81\pm0.74$                     | $81.75\pm0.53$                     | $72.8\pm0.7$                       | $85.0\pm0.5$                       | $47.2\pm0.6$                       | $62.5\pm0.6$                       |
| Support-based initialization (train + val)                   | WRN-28-10                        | $58.47 \pm 0.66$                   | $75.56\pm0.52$                     | $67.34 \pm 0.69^\dagger$           | $83.32\pm0.51^\dagger$             | $72.14 \pm 0.69^{\dagger}$         | $85.21\pm0.49^\dagger$             | $45.08 \pm 0.61$                   | $60.05\pm0.60$                     |
| Fine-tuning (train + val)                                    | WRN-28-10                        | $59.62\pm0.66$                     | $\textbf{79.93} \pm \textbf{0.47}$ | $66.23 \pm 0.68$                   | $\textbf{86.08} \pm \textbf{0.47}$ | $70.07\pm0.67$                     | $\textbf{87.26} \pm \textbf{0.45}$ | $43.80 \pm 0.58$                   | $64.40\pm0.58$                     |
| Transductive fine-tuning (train + val)                       | WRN-28-10                        | $\textbf{68.11} \pm \textbf{0.69}$ | $\textbf{80.36} \pm \textbf{0.50}$ | $\textbf{72.87} \pm \textbf{0.71}$ | $\textbf{86.15} \pm \textbf{0.50}$ | $\textbf{78.36} \pm \textbf{0.70}$ | $\textbf{87.54} \pm \textbf{0.49}$ | $50.44 \pm 0.68$                   | $\textbf{65.74} \pm \textbf{0.60}$ |

#### The ImageNet-21k dataset



1-shot, 5-way accuracies are as high as 89%, 1-shot 20-way accuracies are about 70%.

#### Menu

- An introduction
- Few-shot image classification
- Fundamental limits of representation learning

#### Information Bottleneck Principle

A generalization of rate-distortion theory for learning representations of data [Tishby et al., 2000]

$$X \to Z \to Y$$

Z is a representation of the data X. We want

- -Z to be sufficient to predict the target Y, and
- -Z to be small in size, e.g., few number of bits.

#### Information Bottleneck Principle

A generalization of rate-distortion theory for learning representations of data [Tishby et al., 2000]

$$X \to Z \to Y$$

Z is a representation of the data X. We want

- -Z to be sufficient to predict the target Y, and
- -Z to be small in size, e.g., few number of bits.

$$\min_{Z|X, Y|Z} \{I(X; Z) - I(Z; Y)\}.$$

Doing well on one task requires throwing away nuisance information [Achille & Soatto, 2017].

The IB Lagrangian simply minimizes I(X; Z), it does not let us measure what was thrown away.

The IB Lagrangian simply minimizes I(X; Z), it does not let us measure what was thrown away.

Choose a canonical task to measure discarded information. Setting

Y := X,

i.e., reconstruction of data, gives a special task. It is the superset of all tasks and forces the model to learn lossless representations.

The IB Lagrangian simply minimizes I(X; Z), it does not let us measure what was thrown away.

Choose a canonical task to measure discarded information. Setting

Y := X,

i.e., reconstruction of data, gives a special task. It is the superset of all tasks and forces the model to learn lossless representations.

The architecture we will focus on is

$$\begin{array}{cccc} X & \stackrel{\mathsf{Encoder} \ e_{\theta}(z|x)}{\longrightarrow} & Z & \stackrel{\mathsf{Classifier} \ c_{\theta}(y|z)}{\longrightarrow} & Y \\ & & & \downarrow \\ & & & \downarrow \\ & & & & \widehat{\chi} \end{array}$$

#### An auto-encoder

Distortion D measures the quality of reconstruction

$$D = \mathop{\mathbb{E}}_{x \sim p(x)} \left[ -\int \mathrm{d}z \ e(z|x) \log d(x|z) \right].$$

Rate R measures the average excess bits used to encode the representation

$$R = \mathop{\mathbb{E}}_{x \sim p(x)} \left[ \int \mathrm{d}z \ e(z|x) \log \frac{e(z|x)}{m(z)} \right].$$

#### **Rate-Distortion curve**

The Lagrangian

$$F(\lambda) = \min_{e_{\theta}(z|x), m_{\theta}(z), d_{\theta}(x|z)} \{R + \lambda D\}.$$

is a relaxation of the fact that given a variational family and data there is an optimal value R = func(D) (that best sandwiches the inequality  $H - D \le I(X; Z) \le R$ ).

#### **Rate-Distortion curve**

The Lagrangian

$${\sf F}(\lambda) = \min_{e_{ heta}(z|x), m_{ heta}(z), d_{ heta}(x|z)} \ \{R + \lambda D\}.$$

is a relaxation of the fact that given a variational family and data there is an optimal value R = func(D) (that best sandwiches the inequality  $H - D \le I(X; Z) \le R$ ).



#### Rate-Distortion-Classification (RDC) surface

Let us extend the Lagrangian to

$$F(\lambda,\gamma) = \min_{e_{\theta}(z|x), m_{\theta}(z), d_{\theta}(x|z)} \{R + \lambda D + \gamma C\}$$

where the classification loss is

$$C = \mathbb{E}_{x \sim p(x), y \sim p(y|x)} \left[ -\int dz \ e(z|x) \log c(y|z) \right]$$

#### Rate-Distortion-Classification (RDC) surface

Let us extend the Lagrangian to

$$F(\lambda, \gamma) = \min_{e_{\theta}(z|x), m_{\theta}(z), d_{\theta}(x|z)} \{R + \lambda D + \gamma C\}$$

where the classification loss is

$$C = \mathbb{E}_{x \sim p(x), y \sim p(y|x)} \left[ -\int dz \ e(z|x) \log c(y|z) \right]$$

Can also include other quantities like the entropy  ${\cal S}$  of the model parameters

$$S = \mathbb{E}_{x \sim p(x), y \sim p(y|x)} \left[ \log \frac{p(\theta|\{x, y\})}{m(\theta)} \right]$$

#### Rate-Distortion-Classification (RDC) surface



The existence of a convex surface func(R, D, C, S) = 0 tying together these functionals allows a formal connection to thermodynamics [Alemi & Fischer 2018]

$$\mathrm{d}R = -\lambda \; \mathrm{d}D - \gamma \; \mathrm{d}C - \sigma \; \mathrm{d}S.$$

Just like energy is conserved in physical processes, information is conserved in the model, either it is in the encoder-classifier pair or it is in the decoder.

#### Equilibrium surface of optimal free-energy

The RDC surface determines all possible representations that can be learnt from given data. Can solve the variational problem for  $F(\lambda, \gamma)$  to get

$$Z_{ heta,x} = \int \mathrm{d}z \; m_{ heta}(z) \; d_{ heta}(x|z)^{\lambda} \; c_{ heta}(y_x|z)^{\gamma}$$

and

$$F(\lambda,\gamma) = \min_{\theta \in \Theta} \langle -\log Z_{\theta,x} \rangle_{x \sim \rho(x)} := J(\theta,\lambda,\gamma)$$

#### Equilibrium surface of optimal free-energy

The RDC surface determines all possible representations that can be learnt from given data. Can solve the variational problem for  $F(\lambda, \gamma)$  to get

$$Z_{ heta,x} = \int \mathrm{d}z \; m_{ heta}(z) \; d_{ heta}(x|z)^{\lambda} \; c_{ heta}(y_x|z)^{\gamma}$$

and

$$F(\lambda,\gamma) = \min_{\theta \in \Theta} \langle -\log Z_{\theta,x} \rangle_{x \sim p(x)} := J(\theta,\lambda,\gamma)$$

The surface depends on data p(x, y).

A quasi-static process happens slowly enough for the system to remain in equilibrium with its surroundings, e.g., reversible expansion of an ideal gas.

A quasi-static process happens slowly enough for the system to remain in equilibrium with its surroundings, e.g., reversible expansion of an ideal gas.

We will create a quasi-static process to travel on the RDC surface. This constraint is

$$abla_{ heta} J( heta,\lambda,\gamma) = 0 ext{ for all } heta \in \Theta_{\lambda,\gamma}.$$

A quasi-static process happens slowly enough for the system to remain in equilibrium with its surroundings, e.g., reversible expansion of an ideal gas.

We will create a quasi-static process to travel on the RDC surface. This constraint is

$$abla_ heta \, J( heta,\lambda,\gamma) = \mathsf{0} \, \, \mathsf{for} \, \, \mathsf{all} \, \, heta \in \Theta_{\lambda,\gamma}.$$

e.g., if we want classification loss to be constant in time, we need

$$\frac{d}{dt} \nabla_{\theta} J = 0 \qquad (Quasi-Static Condition)$$
$$\frac{d}{dt} C = 0 \qquad (Iso-classification Condition).$$

A quasi-static process happens slowly enough for the system to remain in equilibrium with its surroundings, e.g., reversible expansion of an ideal gas.

We will create a quasi-static process to travel on the RDC surface. This constraint is

$$abla_ heta \, J( heta,\lambda,\gamma) = \mathsf{0} ext{ for all } heta \in \Theta_{\lambda,\gamma}.$$

e.g., if we want classification loss to be constant in time, we need

$$\begin{aligned} &\frac{\mathrm{d}}{\mathrm{d}t} \, \nabla_\theta \, J = 0 \qquad & (\text{Quasi-Static Condition}) \\ &\frac{\mathrm{d}}{\mathrm{d}t} \, C = 0 \qquad & (\text{Iso-classification Condition}). \end{aligned}$$

Can also impose other constraints, e.g.,

$$\frac{\mathsf{d}}{\mathsf{d}t}\{C+\gamma^{-1}R\}=0$$

which is the objective for learning Bayesian neural networks.

#### Implementing processes on the RDC surface



Could pick particular values of  $(\dot{\lambda},\dot{\gamma})$  to get

$$0 = \frac{\mathsf{d}}{\mathsf{d} t} \, \nabla_\theta \, J = \nabla_\theta^2 \, J \, \dot{\theta} + \dot{\lambda} \frac{\partial}{\partial \lambda} \, \nabla_\theta \, J + \dot{\gamma} \frac{\partial}{\partial \gamma} \, \nabla_\theta \, J$$

this requires inverting  $\nabla_{\theta}^2 J$ .

#### Implementing processes on the RDC surface



Could pick particular values of  $(\dot{\lambda},\dot{\gamma})$  to get

$$0 = \frac{\mathsf{d}}{\mathsf{d}t} \, \nabla_\theta \, J = \nabla_\theta^2 \, J \; \dot{\theta} + \dot{\lambda} \frac{\partial}{\partial \lambda} \, \nabla_\theta \, J + \dot{\gamma} \frac{\partial}{\partial \gamma} \, \nabla_\theta \, J$$

this requires inverting  $\nabla^2_{\theta} J$ .

We exploit constraints like 0 =  ${\cal C}_\lambda \dot\lambda + {\cal C}_\gamma \dot\gamma$  to get

$$\dot{\lambda} = -\alpha \frac{\partial}{\partial \gamma} C = -\alpha \frac{\partial^2}{\partial \gamma^2} F$$
$$\dot{\gamma} = \alpha \frac{\partial}{\partial \lambda} C = \alpha \frac{\partial^2}{\partial \lambda \partial \gamma} F$$

#### Iso-C process for different initial $(\lambda, \gamma)$ : CIFAR-10



The RDC surface depends on data p(x, y). We now move the data distribution from the source task to the target task, e.g., interpolate it as

$$p(x, y, t) = (1 - t) p^{s}(x, y) + t p^{t}(x, y).$$

The RDC surface depends on data p(x, y). We now move the data distribution from the source task to the target task, e.g., interpolate it as

$$p(x, y, t) = (1 - t) p^{s}(x, y) + t p^{t}(x, y).$$

Can also interpolate data using

 $p(x,t) = \operatorname{argmin}_{p} (1-t)W_{2}^{2}(p^{s},p) + tW_{2}^{2}(p,p^{t})$  and solve the discrete optimal transport problem.

The RDC surface depends on data p(x, y). We now move the data distribution from the source task to the target task, e.g., interpolate it as

$$p(x, y, t) = (1 - t) p^{s}(x, y) + t p^{t}(x, y).$$

Can also interpolate data using

 $p(x,t) = \operatorname{argmin}_{p} (1-t)W_{2}^{2}(p^{s},p) + tW_{2}^{2}(p,p^{t})$  and solve the discrete optimal transport problem.

The quasi-static iso-classification process

$$0 = \frac{\mathsf{d}}{\mathsf{d}t} \, \nabla_\theta \, J = \frac{\mathsf{d}}{\mathsf{d}t} \, C$$

can be executed on this changing data distribution.

The RDC surface depends on data p(x, y). We now move the data distribution from the source task to the target task, e.g., interpolate it as

$$p(x, y, t) = (1 - t) p^{s}(x, y) + t p^{t}(x, y).$$

Can also interpolate data using

 $p(x,t) = \operatorname{argmin}_{p} (1-t)W_{2}^{2}(p^{s},p) + tW_{2}^{2}(p,p^{t})$  and solve the discrete optimal transport problem.

The quasi-static iso-classification process

$$0 = \frac{\mathsf{d}}{\mathsf{d}t} \, \nabla_\theta \, J = \frac{\mathsf{d}}{\mathsf{d}t} \, C$$

can be executed on this changing data distribution.

This is a completely controlled mechanism to transfer representations.

#### Iso-C process: MNIST 0-4 to 5-9



#### Iso-C process: CIFAR-10 Vehicles to Animals



# Summary

Simple methods such as transductive fine-tuning work extremely well for few-shot learning. This is really because of powerful function approximators such as neural networks.

The RDC surface is a fundamental quantity and enables principled methods for transfer learning. Also unlocks new paths to understanding regularization and properties of neural architecture for classical supervised learning.

We did well in the era of big data without understanding much about data; this is unlikely to work in the age of little data.

#### Email questions to pratikac@seas.upenn.edu

#### Read more at

- Chaudhari, P. and Soatto, S.. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. arXiv:1710.11029. ICLR 2018.
- Dhillon, G., Chaudhari, P., Ravichandran, A., and Soatto, S. A baseline for few-shot image classification. arXiv:1909.02729. ICLR 2020.
- Gao, Y., and Chaudhari, P. A free-energy principle for representation learning. arXiv:2002.12406. ICML 2020