## Minimal Achievable Sufficient Statistic Learning

Milan Cvitkovic California Institute of Technology/Amazon Web Services

> Günther Koliander Austrian Academy of Sciences



Goal: Find representation Z of X that is useful for producing Y

$$X \rightarrow Z \rightarrow Y$$

#### Goal: Find representation Z of X that is useful for producing Y

Hypothesis: Z should be a *minimal* sufficient statistic of X for Y [1, 2]

$$X \rightarrow Z \rightarrow Y$$

# Hypothesis: Z should be a *minimal* sufficient statistic of X for Y [1, 2]

$$X \rightarrow Z \rightarrow Y$$

Statistic: Z = f(X)

Sufficient: p(X | Y, Z) = p(X | Z)

Minimal: g(Z) isn't sufficient for any non-invertible g(Z)

# Hypothesis: Z should be a *minimal* sufficient statistic of X for Y [1, 2]

$$X \rightarrow Z \rightarrow Y$$

Statistic: Z = f(X)

Sufficient: p(X | Y, Z) = p(X | Z)

A little strict for ML

Minimal: g(Z) isn't sufficient for any non-invertible g

# Hypothesis: Z should be a *minimal* achievable sufficient statistic of X for Y

$$X \rightarrow Z \rightarrow Y$$

Statistic: Z = f(X) for f in F

Sufficient: p(X | Y, Z) = p(X | Z)

#### Minimal Achievable:

g(Z) isn't sufficient for any Lipschitz, non-invertible g where gof in F

### How do we find a minimal (achievable) sufficient statistic?

X discrete:

$$f \in \arg\min_{S \in \mathcal{F}} I(X, S(X))$$
 
$$s.t. \ I(S(X), Y) = \max_{S' \in \mathcal{F}} I(S'(X), Y)$$

X, S continuous: X

### How do we find a minimal (achievable) sufficient statistic?

X discrete:

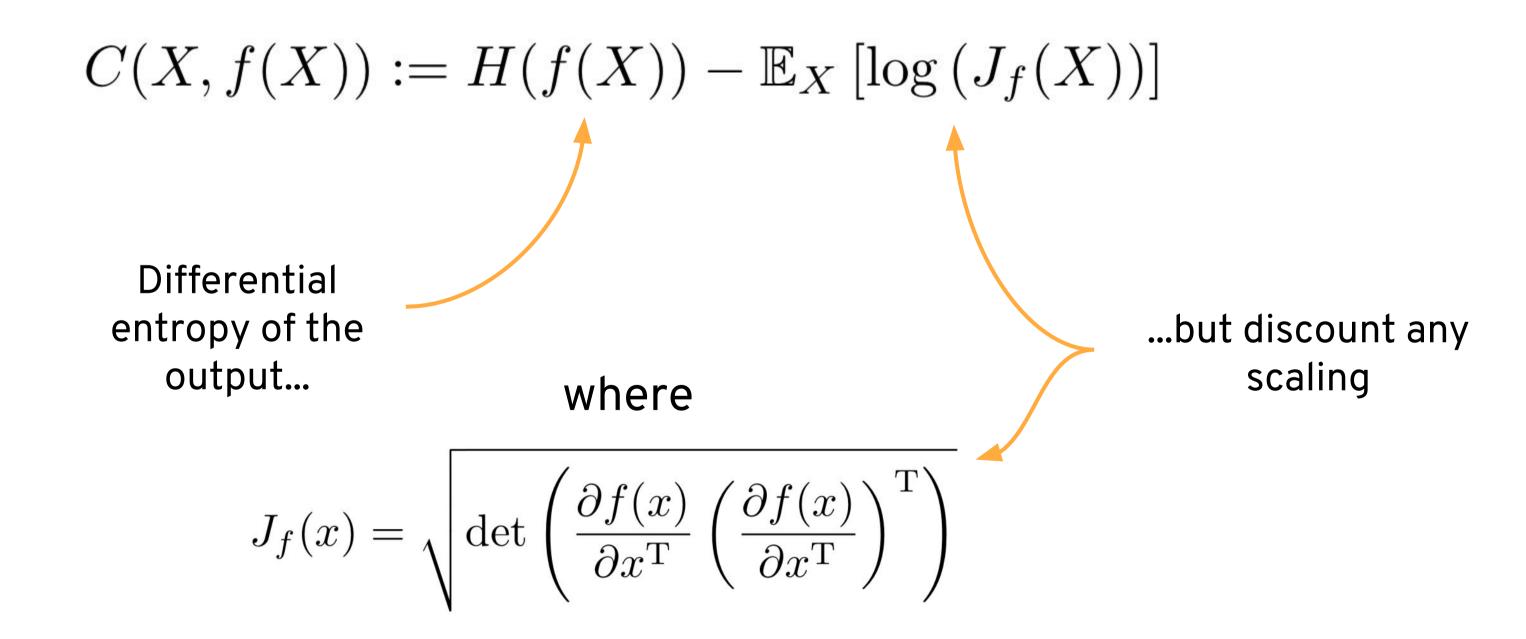
$$f \in \arg\min_{S \in \mathcal{F}} I(X, S(X))$$
 
$$s.t. \ I(S(X), Y) = \max_{S' \in \mathcal{F}} I(S'(X), Y)$$

X, S continuous: X

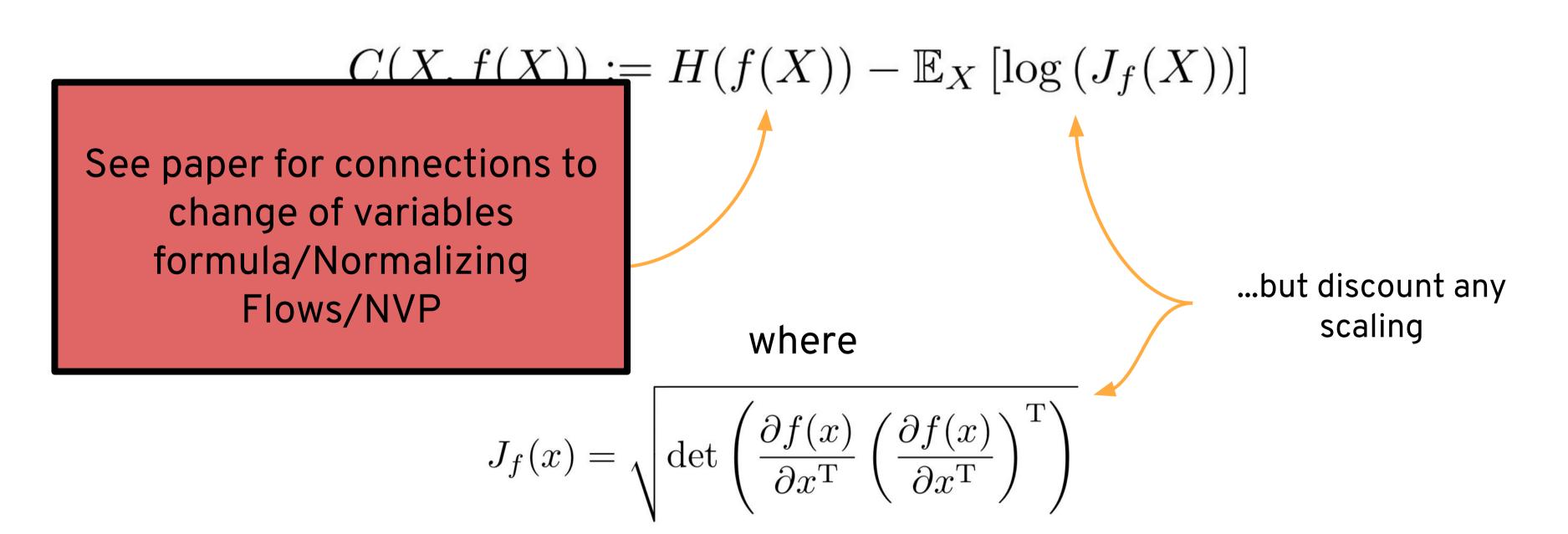
$$I(X, S(X)) = \infty$$

This is a generally problematic issue in machine learning [3, 4].

#### Solution: Conserved Differential Information



#### Solution: Conserved Differential Information



### How do we find a minimal (achievable) sufficient statistic?

X discrete:

$$f \in \arg\min_{S \in \mathcal{F}} I(X, S(X))$$
 
$$s.t. \ I(S(X), Y) = \max_{S' \in \mathcal{F}} I(S'(X), Y)$$

X, S continuous:

$$f \in \arg\min_{S \in \mathcal{F}} C(X, S(X))$$
 
$$s.t. \ I(S(X), Y) = \max_{S'} I(S'(X), Y)$$

## MASS Learning

#### Minimize:

$$\mathcal{L}_{MASS}(f) := H(Y|f(X)) + \beta H(f(X)) - \beta \mathbb{E}_X[\log J_f(X)]$$

Make the output predictable from the representation...

...while keeping the representation as low-entropy as possible...

...but scaling things doesn't count

## Results

#### Same accuracy as standard training and VIB

CIFAR-10, ResNet20, 4 trials

METHOD	TRAINING SET SIZE			
Метнор	2500	10,000	40,000	
SoftmaxCE	$50.0 \pm 0.7$	$\textbf{67.5} \pm \textbf{0.8}$	$\textbf{81.7} \pm \textbf{0.3}$	
VIB, $\beta$ =1e-3	$49.5 \pm 1.1$	$66.9 \pm 1.0$	$81.0 \pm 0.3$	
VIB, $\beta$ =1e-4	$49.4 \pm 1.0$	$66.4 \pm 0.5$	$81.2 \pm 0.4$	
VIB, $\beta$ =1e-5	$50.0 \pm 1.1$	$67.9 \pm 0.8$	$80.9 \pm 0.5$	
VIB, $\beta$ =0	$50.6 \pm 0.8$	$67.1 \pm 1.0$	$81.5 \pm 0.2$	
MASS, $\beta$ =1e-3	$38.2 \pm 0.7$	$59.6 \pm 0.8$	$75.8 \pm 0.5$	
MASS, $\beta$ =1e-4	$49.9 \pm 1.0$	$66.6 \pm 0.4$	$80.6 \pm 0.5$	
MASS, $\beta$ =1e-5	$50.1 \pm 0.5$	$67.4 \pm 1.0$	$81.6 \pm 0.4$	
MASS, $\beta$ =0	$\textbf{50.2} \pm \textbf{1.0}$	$67.4 \pm 0.3$	$81.5 \pm 0.2$	

#### Improved Uncertainty Quantification

CIFAR-10, ResNet20, 4 trials

Method	Test Accuracy	Entropy	NLL	Brier Score
SoftmaxCE	$81.7 \pm 0.3$	$0.087 \pm 0.002$	$1.45 \pm 0.04$	$0.0324 \pm 0.0005$
VIB, $\beta$ =1e-3	$81.0 \pm 0.3$	$0.089 \pm 0.003$	$1.51 \pm 0.04$	$0.0334 \pm 0.0005$
VIB, $\beta$ =1e-4	$81.2 \pm 0.4$	$0.092 \pm 0.002$	$1.46 \pm 0.05$	$0.0331 \pm 0.0007$
VIB, $\beta$ =1e-5	$80.9 \pm 0.5$	$0.087 \pm 0.005$	$1.58 \pm 0.08$	$0.0339 \pm 0.0008$
VIB, $\beta$ =0	$81.5 \pm 0.2$	$0.079 \pm 0.001$	$1.70 \pm 0.06$	$0.0331 \pm 0.0007$
MASS, $\beta$ =1e-3	$75.8 \pm 0.5$	$0.139 \pm 0.003$	$1.66 \pm 0.07$	$0.0417 \pm 0.0011$
MASS, $\beta$ =1e-4	$80.6 \pm 0.5$	$0.109 \pm 0.002$	$\boldsymbol{1.33 \pm 0.02}$	$0.0337 \pm 0.0008$
MASS, $\beta$ =1e-5	$81.6 \pm 0.4$	$0.095 \pm 0.003$	$1.36 \pm 0.03$	$0.0320 \pm 0.0005$
MASS, $\beta$ =0	$81.5 \pm 0.2$	$0.092 \pm 0.000$	$1.43 \pm 0.04$	$0.0325 \pm 0.0004$

Caveat: current implementation expensive, but ample room for improvement

More results at our poster.

Code available at <a href="mailto:github.com/mwcvitkovic/MASS-Learning">github.com/mwcvitkovic/MASS-Learning</a>.