

Systematic generalisation with group invariant predictions

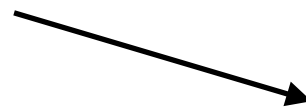
Faruk Ahmed^{1,*}, Yoshua Bengio^{1,2}, Harm van Seijen³, Aaron Courville^{1,2}

¹ Université de Montréal, Mila, ² CIFAR Fellow, ³ Microsoft Research

*Correspondence to faruk.ahmed@umontreal.ca.

easier-to-learn

h_n

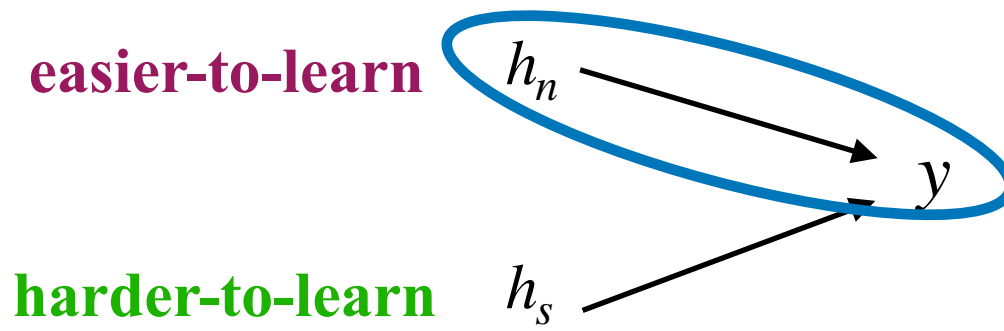


y

harder-to-learn

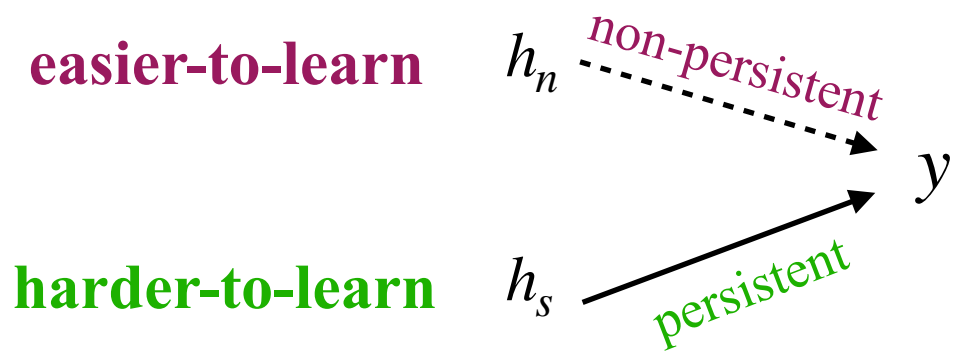
h_s

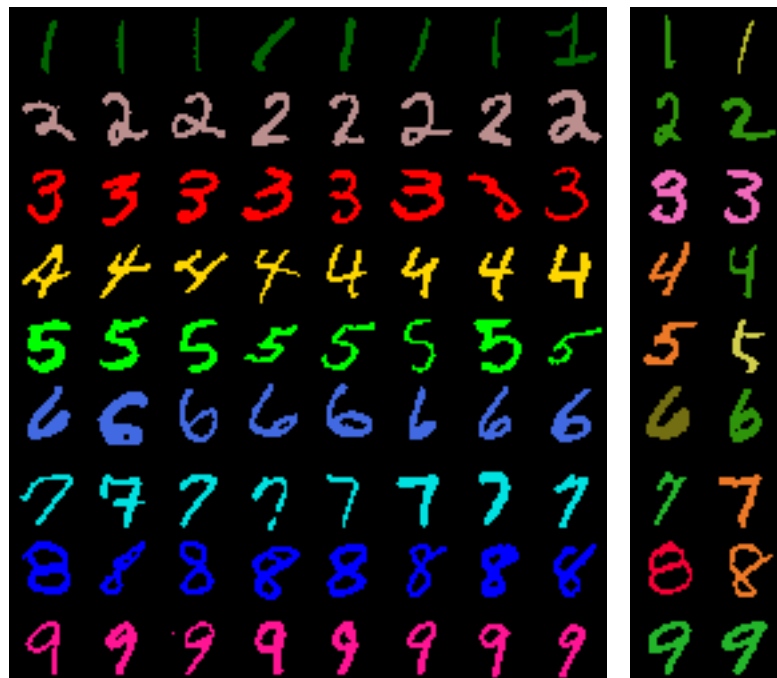




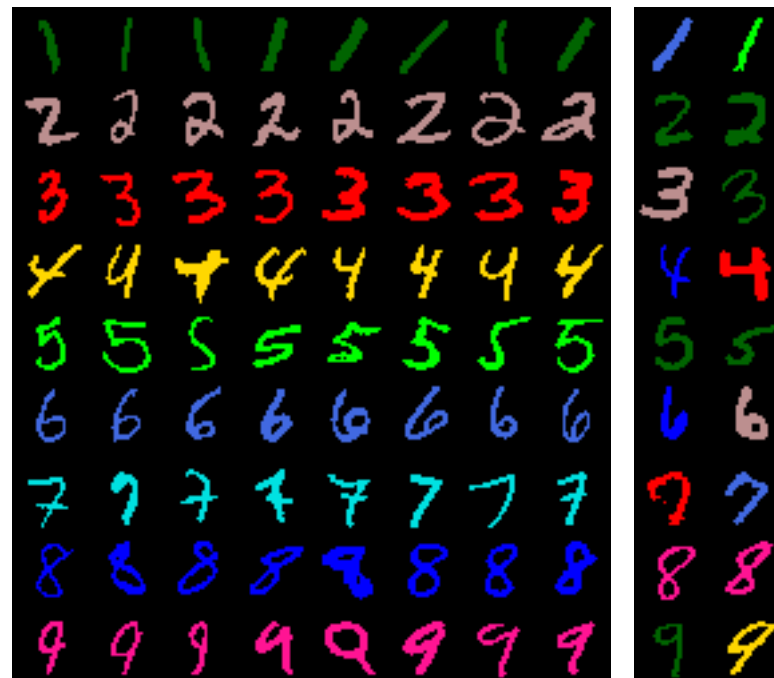
Occam's razor







Minority group colours are different



Minority group colours are recombinations



Systematically shifted test set

Minority colours	In-distribution	Systematic shift
Different	99.60 ± 0.02	38.72 ± 2.27
Recombinations	98.67 ± 0.39	97.56 ± 0.05

Assume data x is synthesised from *non-semantic factors* h_n and *semantic factors* h_s

$$x = \mathcal{C}(h_n, h_s) .$$

If $\hat{y}(x)$ is our prediction, we can compute average accuracy, for synthesised x

$$\mathbb{E} \left[\mathbf{1} \{ \hat{y}(\mathcal{C}(h_s, h_n)) = y \} \right]$$

with different sampling choices of h_s and h_n .

- *In-distribution generalisation:*

$$h_s \sim p(h_s | y), h_n \sim p(h_n | y)$$

- *Generalisation under systematic-shift:*

$$h_s \sim p(h_s | y), h_n \sim p(h_n | y') \quad \text{where } y' \sim p(y) \text{ s.t. } y' \neq y$$

- *Generalisation under non-systematic-shift:*

$$h_s \sim p(h_s | y), h_n \approx p(h_n)$$

- *Semantic anomaly detection:*

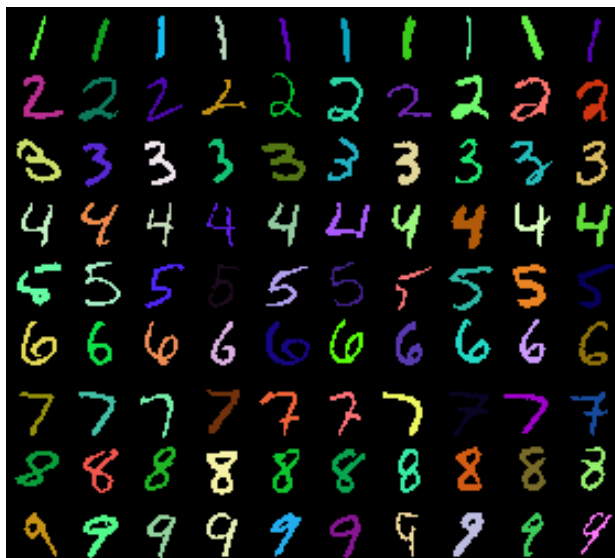
$$h_s \approx p(h_s), h_n \sim p(h_n)$$



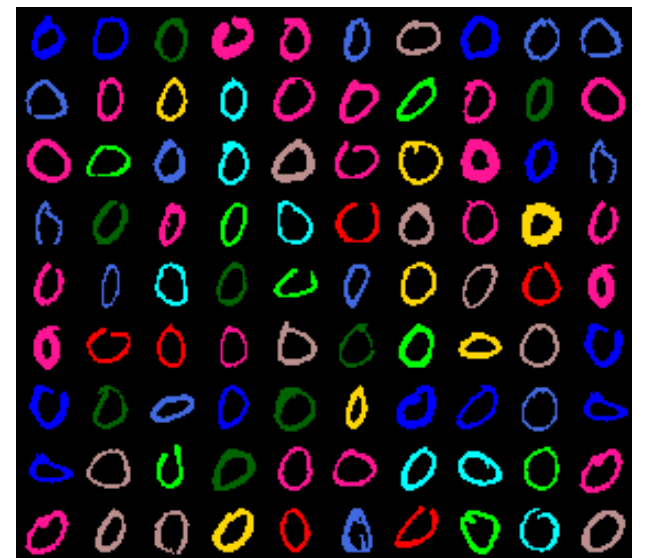
In-distribution



Systematic
shift



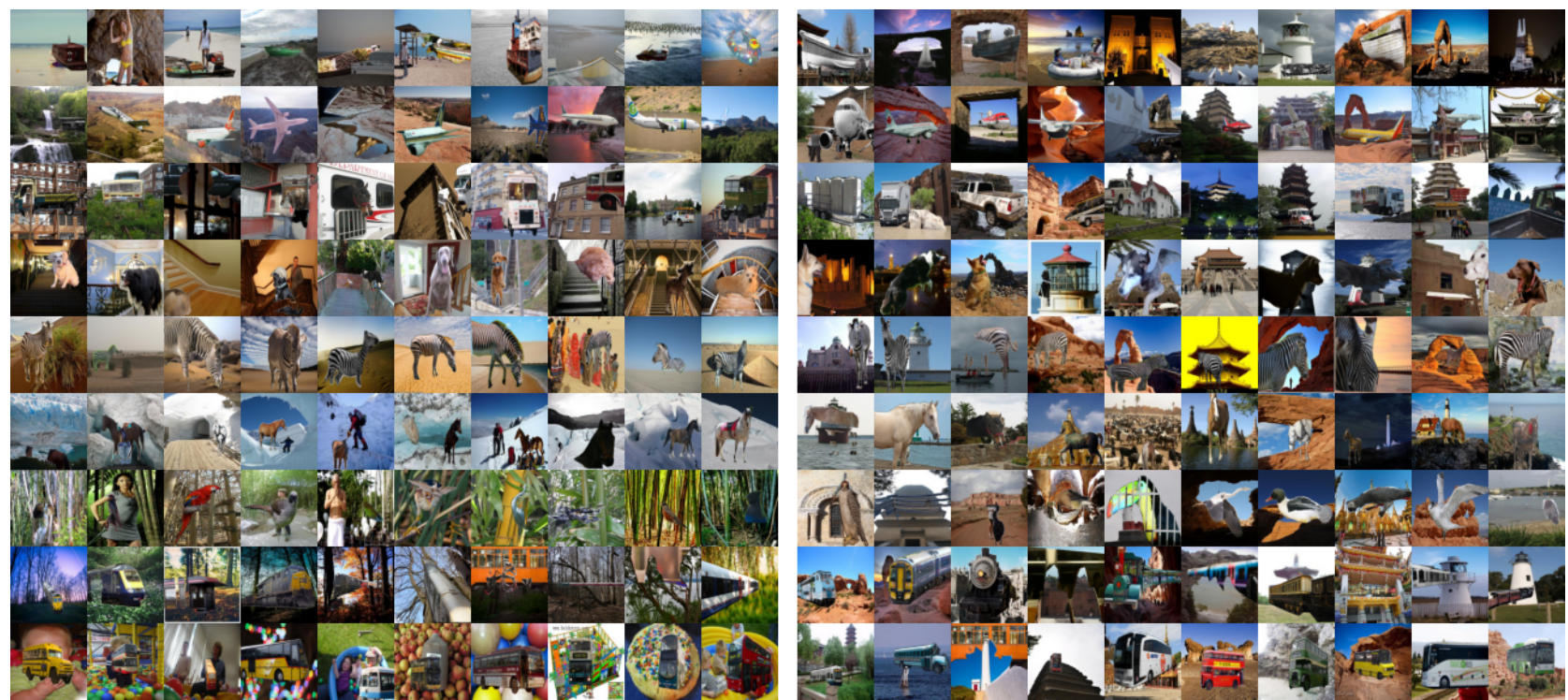
Non-systematic
shift



Semantic anomalies



COCO-on-Colours



COCO-on-Places

Do invariance methods/penalties for multi-group data help at such shifts?

Feature-distribution matching
Group-distributionally robust optimisation
Invariant Risk Minimisation
Others

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *CoRR*, 2019.

David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). *CoRR*, 2020.

H. Li, S. J. Pan, S. Wang, and A. C. Kot. Domain generalization with adversarial feature learning. pp. 5400–5409, 2018.

Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. Domain generalization via conditional invariant representations. 2018.

Shiori Sagawa, Pang Wei Koh, Tatsunori Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *ICLR*, 2020.

Predictive Group Invariance (PGI)

Feature extractor $f_\theta(x)$, predictive distribution:

$$p_w(y|x) = \sigma(w^\top f_\theta(x)).$$

Split data into the biased majority ($\alpha = 0$) and unbiased minority group ($\alpha = 1$), for every class c , such that

$$\begin{aligned} x^{(i)} &\sim \mathbb{P}^c \text{ if } \alpha^{(i)} = 0, y^{(i)} = c, \\ x^{(j)} &\sim \mathbb{Q}^c \text{ if } \alpha^{(j)} = 1, y^{(j)} = c. \end{aligned}$$

Loss function:

$$\ell(\theta, w | \mathcal{D}) + \lambda \left[\sum_c d\left(\mathbb{E}_{x \sim \mathbb{Q}^c}[p_{\tilde{w}}(y|x)], \mathbb{E}_{x \sim \mathbb{P}^c}[p_{\tilde{w}}(y|x)]\right) \right]_{\tilde{w}=w \text{ (fixed)}}$$

ℓ = standard (regularised) ERM term

d = KL-divergence

Predictive Group Invariance (PGI)

Feature extractor $f_\theta(x)$, predictive distribution:

$$p_w(y|x) = \sigma(w^\top f_\theta(x)).$$

Split data :

“Environment Inference for Invariant Learning”

COLOURED MNIST	COCO-ON-COLOURS	COCO-ON-PLACES
97.26 ± 0.71	98.22 ± 1.05	80.43 ± 1.41

Loss function:

Partitioning accuracy

$$\ell(\theta, w | \mathcal{D}) + \lambda \left[\sum_c d\left(\mathbb{E}_{x \sim \mathbb{Q}^c}[p_{\tilde{w}}(y|x)], \mathbb{E}_{x \sim \mathbb{P}^c}[p_{\tilde{w}}(y|x)]\right) \right]_{\tilde{w}=w \text{ (fixed)}}$$

ℓ = standard (regularised) ERM term

d = KL-divergence

Table 2: Generalisation results on COLOURED MNIST.

Methods	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
Base (ERM)	99.60 \pm 0.02	53.26 \pm 1.89	38.72 \pm 2.27	7.70 \pm 0.23
IRMv1	99.47 \pm 0.05	63.24 \pm 3.04	55.19 \pm 1.07	11.54 \pm 1.18
REx	98.95 \pm 0.11	72.12 \pm 1.90	71.18 \pm 3.27	15.54 \pm 2.05
GroupDRO	89.47 \pm 4.52	70.53 \pm 1.79	79.17 \pm 1.64	35.15 \pm 10.83
Reweight	98.51 \pm 0.12	75.01 \pm 1.28	84.85 \pm 0.61	28.60 \pm 1.11
cIRMv1	99.36 \pm 0.25	65.78 \pm 3.53	61.09 \pm 5.30	14.16 \pm 2.12
cREx	98.56 \pm 0.12	74.35 \pm 2.09	80.01 \pm 2.11	22.02 \pm 2.52
cGroupDRO	95.65 \pm 3.23	75.41 \pm 3.45	81.14 \pm 2.41	26.61 \pm 6.61
cMMD	99.40 \pm 0.03	97.17 \pm 0.59	97.86 \pm 0.16	78.32 \pm 4.15
PGI	99.05 \pm 0.08	98.58 \pm 0.06	98.48 \pm 0.05	89.42 \pm 1.95

Table 3: Generalisation performance on COCO-ON-COLOURS.

Methods	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
Base (ERM)	90.57 \pm 1.28	26.81 \pm 4.93	1.10 \pm 0.36	5.47 \pm 0.08
IRMv1	91.61 \pm 0.38	32.30 \pm 4.52	2.11 \pm 0.30	5.81 \pm 0.17
REx	91.69 \pm 0.50	36.57 \pm 4.03	2.69 \pm 0.81	5.73 \pm 0.14
GroupDRO	43.06 \pm 2.26	41.32 \pm 4.39	43.24 \pm 2.89	20.05 \pm 3.08
Reweight	42.42 \pm 3.47	47.56 \pm 2.27	49.12 \pm 1.63	18.15 \pm 3.81
cIRMv1	91.53 \pm 0.31	31.11 \pm 4.51	1.74 \pm 0.40	5.87 \pm 0.16
cREx	74.75 \pm 14.14	32.29 \pm 7.71	29.75 \pm 5.16	19.77 \pm 14.98
cGroupDRO	41.10 \pm 2.37	41.83 \pm 2.96	42.10 \pm 2.15	21.81 \pm 5.40
cMMD	89.87 \pm 1.13	55.02 \pm 2.29	27.36 \pm 1.57	8.82 \pm 0.70
PGI	78.23 \pm 2.01	55.57 \pm 4.60	51.62 \pm 3.09	18.84 \pm 2.11

Table 4: Generalisation performance on COCO-ON-PLACES.

Methods	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
Base (ERM)	81.06 \pm 1.01	45.25 \pm 0.96	29.18 \pm 1.24	9.21 \pm 0.21
IRMv1	80.93 \pm 0.71	45.17 \pm 0.92	28.78 \pm 0.73	9.39 \pm 0.60
REx	81.55 \pm 0.70	45.35 \pm 0.92	29.56 \pm 0.77	9.46 \pm 0.51
GroupDRO	76.05 \pm 0.87	43.72 \pm 0.43	31.83 \pm 0.54	9.61 \pm 0.55
Reweight	81.14 \pm 0.80	45.84 \pm 0.70	30.37 \pm 1.16	9.75 \pm 0.69
cIRMv1	80.08 \pm 1.90	44.96 \pm 2.88	30.06 \pm 2.07	9.64 \pm 0.94
cREx	81.50 \pm 0.76	45.44 \pm 0.96	29.12 \pm 0.97	9.17 \pm 0.59
cGroupDRO	78.25 \pm 0.31	41.69 \pm 0.08	28.16 \pm 0.91	9.45 \pm 0.22
cMMD	79.64 \pm 0.73	49.44 \pm 0.99	35.86 \pm 0.66	9.80 \pm 0.45
PGI	75.00 \pm 0.85	46.10 \pm 0.79	36.25 \pm 0.42	11.12 \pm 0.85

Hyper-parameter selection

NS+S	Both non-systematic* and systematic
NS	Non-systematic* only
NS+ID	Non-systematic* + in-distribution
ID	In-distribution only

*Non-systematic validation sets use different colours than in training and test sets

Table 5: Hyper-parameters with different validation sets for COLOURED MNIST

Validation	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
NS+S (PGI)	99.05 ± 0.08	98.58 ± 0.06	98.48 ± 0.05	89.42 ± 1.95
NS (PGI)	99.31 ± 0.05	98.21 ± 0.26	97.54 ± 0.41	76.00 ± 4.06
NS+ID (PGI)	99.30 ± 0.07	98.31 ± 0.27	97.48 ± 0.45	76.07 ± 5.67
ID only (PGI)	99.69 ± 0.03	63.62 ± 2.05	58.18 ± 2.05	11.81 ± 1.89
Base (ERM)	99.60 ± 0.02	53.26 ± 1.89	38.72 ± 2.27	7.70 ± 0.23

Table 6: Hyper-parameters with different validation sets for COCO-ON-COLOURS

Validation	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
NS+S (PGI)	78.23 ± 2.01	55.57 ± 4.60	51.62 ± 3.09	18.84 ± 2.11
NS (PGI)	85.78 ± 1.45	51.02 ± 2.32	38.85 ± 2.29	15.71 ± 3.25
NS+ID (PGI)	85.78 ± 1.45	51.02 ± 2.32	38.85 ± 2.29	15.71 ± 3.25
ID only (cMMD)	92.51 ± 0.41	44.59 ± 3.28	10.48 ± 0.98	6.05 ± 0.23
Base (ERM)	90.57 ± 1.28	26.81 ± 4.93	1.10 ± 0.36	5.47 ± 0.08

Table 7: Hyper-parameters with different validation sets for COCO-ON-PLACES

Validation	In-distribution	Non-systematic shift	Systematic shift	Anomaly detection
NS+S (cMMD)	79.64 ± 0.73	49.44 ± 0.99	35.86 ± 0.66	9.80 ± 0.45
NS (cMMD)	79.64 ± 0.73	49.44 ± 0.99	35.86 ± 0.66	9.80 ± 0.45
NS+ID (cMMD)	79.64 ± 0.73	49.44 ± 0.99	35.86 ± 0.66	9.80 ± 0.45
ID only (PGI)	80.99 ± 0.52	47.63 ± 0.90	31.91 ± 0.89	9.59 ± 0.89
Base (ERM)	81.06 ± 1.01	45.25 ± 0.96	29.18 ± 1.24	9.21 ± 0.21

Best performing methods for a mix of non-systematic and systematic generalisation performance

Some takeaways

Invariance methods/penalties across inferred splits of a dataset appear to be useful at improving performance under distributional shift.

Such methods cannot be useful when spurious correlations are completely pervasive; ensuring diversity in data through curation could help.

The question of how to trade off in-distribution performance with that in unexpected situations is relevant. Real-world problems might call for more targeted invariance methods.

Some takeaways

Invariance methods/penalties across inferred splits of a dataset appear to be useful at improving performance under distributional shift.

Such methods cannot be useful when spurious correlations are completely pervasive; ensuring diversity in data through curation could help.

The question of how to trade off in-distribution performance with that in unexpected situations is relevant. Real-world problems might call for more targeted invariance methods.

Some takeaways

Invariance methods/penalties across inferred splits of a dataset appear to be useful at improving performance under distributional shift.

Such methods cannot be useful when spurious correlations are completely pervasive; ensuring diversity in data through curation could help.

The question of how to trade off in-distribution performance with that in unexpected situations is relevant. Real-world problems might call for more targeted invariance methods.

Thank you!