

The Value of Out-of-Distribution Data

Ashwin De Silva^{1,†}, Rahul Ramesh^{2,†},
Carey E. Priebe¹, Pratik Chaudhari^{2,*}, Joshua T. Vogelstein^{1,*}

¹Johns Hopkins University, ²University of Pennsylvania

†Equal Contribution, *Equal Contribution



JOHNS HOPKINS
UNIVERSITY



Penn
UNIVERSITY of PENNSYLVANIA

Does more data always help?

Generalization error for a desired target task always improves with more in-distribution data.

But real data is often $\mathcal{D} \neq \mathcal{D}'$; $m \neq m'$
Even a curated dataset can contain out-of-distribution (OOD) samples.

Does more data always help?

Generalization error for a desired target task always improves with more in-distribution data.

But real data is often ϵ ; δ ; u o] ; m ; o \dagger v
Even a curated dataset can contain out-of-distribution (OOD) samples.

For a model trained on such data, we expect the generalization error on the target task to be ϵ o m o | o i m t e the number of OOD samples.

Modeling heterogeneity in a dataset

Suppose dataset D consists of n target samples and m OOD samples without the knowledge of sample identities.

We seek a hypothesis h that minimizes the generalization error on the target task $\epsilon_t(h)$.

Modeling heterogeneity in a dataset

Suppose dataset D consists of n target samples and m OOD samples without the knowledge of sample identities.

We seek a hypothesis h that minimizes the generalization error on the target task $e_t(h)$.

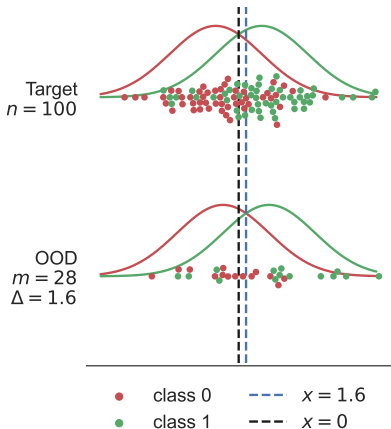
The hypothesis is selected by minimizing the empirical loss,

$$e(h) = \sum_{i=1}^{n+m} \ell(h(x_i); y_i)$$

An example using Fisher's Linear Discriminant

The target and OOD tasks are both gaussian mixture models.

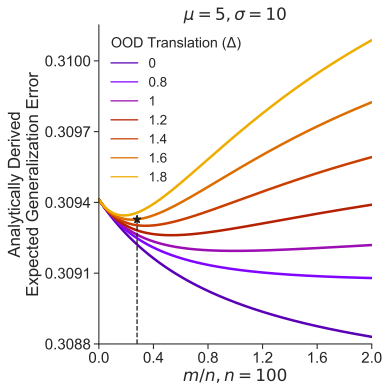
We consider a family of OOD task distributions which are translations of the target distribution.



An example using Fisher's Linear Discriminant

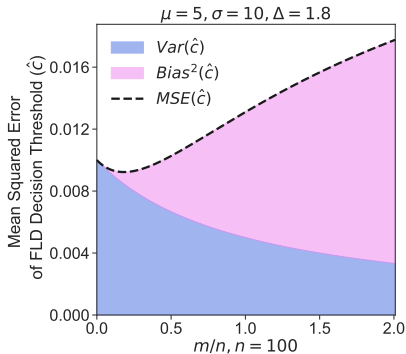
OOD data from the same distribution can both improve or deteriorate the target generalization depending on the number of OOD samples.

Generalization error on the target task can be minimized by choosing the number of OOD samples.



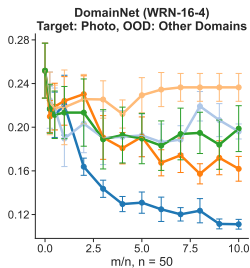
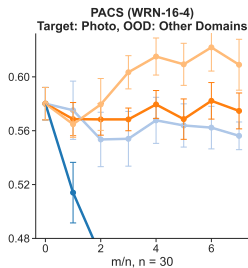
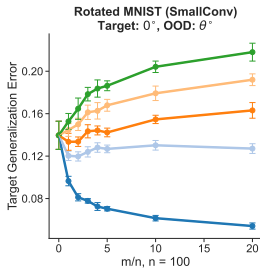
Why does non-monotonicity occur?

More OOD samples decrease the variance but increase the bias. The trade-off depends on the distance between the two distributions.



Non-monotonic trends also occur in popular benchmark datasets

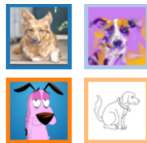
Non-monotonic trends occur due to geometric, semantic nuisances and distribution shifts.



- 15
- 30
- 35
- 40
- 45



- Photo
- Art
- Cartoon
- Sketch



- Photo
- Sketch
- Clipart
- Art
- Quickdraw



Exploiting the non-monotonic trends in generalization error

Assuming that the target and OOD samples are separable, we consider the objective

$$\hat{e}(h) = \hat{e}_t(h) + (1 - \lambda) \hat{e}_o(h):$$

We can compute the optimal λ using an upper bound of the generalization error¹

$$= KBM; \frac{n}{n+m} \left(1 + \frac{s}{4} \frac{m^2}{(n+m)nm} \right):$$

is the ratio between task distance and model capacity.

¹Ben-David et al., "A theory of learning from different domains".

Exploiting the non-monotonic trends in generalization error

Assuming that the target and OOD samples are separable, we consider the objective

$$\hat{e}(h) = \hat{e}_t(h) + (1 - \lambda) \hat{e}_o(h):$$

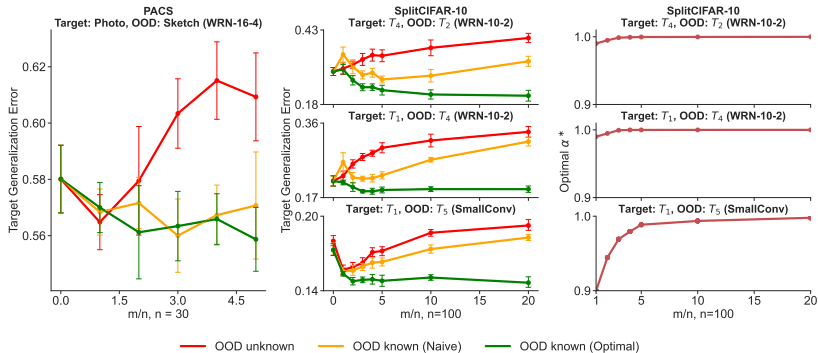
We can compute the optimal λ using an upper bound of the generalization error¹

$$= KBM; \frac{n}{n+m} \left(1 + \frac{s}{4} \frac{m^2}{(n+m)nm} \right) :$$

s is the ratio between task distance and model capacity. In practice we consider, s to be a hyper-parameter.

¹Ben-David et al., "A theory of learning from different domains".

Exploiting the non-monotonic trends in generalization error



Concluding Thoughts

Generalization error can be a non-monotonic function of the number of OOD samples.

A weighted objective between the OOD and target samples can mitigate this non-monotonicity.

For more details and experiments, check out our paper on arXiv

[v É è ž â P , É — W v • Í W - - + 3 P , + 4 1 2](#)

