

# Deep Imbalanced Regression

Yuzhe Yang<sup>1</sup> Kaiwen Zha<sup>1</sup> Ying-Cong Chen<sup>1</sup> Hao Wang<sup>2</sup>  
Dina Katabi<sup>1</sup>

<sup>1</sup>MIT Computer Science & Artificial Intelligence Laboratory

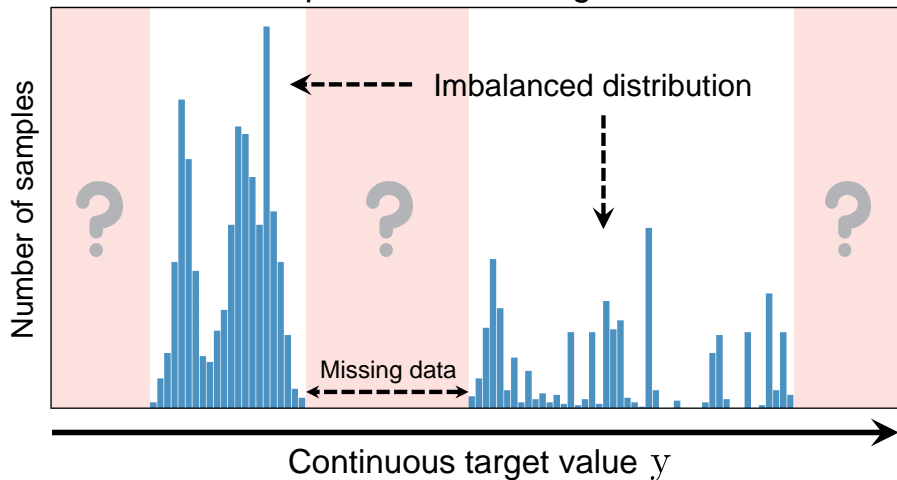
<sup>2</sup>Department of Computer Science, Rutgers University

ICML 2021

Presenter: Gianmarco Midena

26 November 2024

## Deep Imbalanced Regression



# Problem Settings

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ : training set
- $\mathbf{x}_i \in \mathbb{R}^d$ : input
- $y_i \in \mathcal{Y}$ : continuous label or target
- $b_i \in \mathcal{B}$ : discrete label or target
- $\mathcal{Y} \subset \mathbb{R}$ : continuous label space
- $\mathcal{B} = \{1, \dots, M\} \subset \mathbb{Z}^+$ : index space
  - ▶ divides  $\mathcal{Y}$  into  $M$  groups (bins) with equal intervals  $[t_j, t_{j+1})$
  - ▶  $\{[t_0, t_1), \dots, [t_{M-1}, t_M)\}$ : discrete label space
  - ▶  $t_k \in \mathcal{Y}$
  - ▶ minimum resolution
    - ★ e.g.,  $\delta y \triangleq t_{j+1} - t_j = 1$  in age estimation
- $\hat{y}_i = g(\mathbf{z}_i) \in \mathbb{R}$ : predicted continuous label
- $\mathbf{z}_i = f(\mathbf{x}_i; \theta) \in \mathbb{R}^{d'}$ : learned representation
- $\theta$ : trainable model parameters

# Evaluation

- Divide target space into disjoint regions (bins)

- ▶ *Many-shot*: > 100 training examples
- ▶ *Medium-shot*: 20-100 training examples
- ▶ *Few-shot*: < 20 training examples
- ▶ *Zero-shot*: 0 training examples
- Inspired by [Liu et al. 2019](#)

- Metrics

- ▶ Mean Absolute Error (MAE)
- ▶ Mean Squared Error (MSE)
- ▶ Pearson Correlation (PCC)
- ▶ Geometric Mean Error (GM)

$$GM = \sqrt[n]{\prod_{i=1}^n |y_i - \hat{y}_i|}$$

★ Pros: + fairness (uniformity) in prediction

# Datasets - Overview

Dataset	Target type	Target range	Bin size	Max bin density	Min bin density	# Training set	# Val. set	# Test set
IMDB-WIKI-DIR	Age	0 - 186	1	7,149	1	191,509	11,022	11,022
AgeDB-DIR	Age	0 - 101	1	353	1	12,208	2,140	2,140
STS-B-DIR	Text similarity score	0 - 5	0.1	428	1	5,249	1,000	1,000
NYUD2-DIR	Depth	0.7 - 10	0.1	$1.46 \times 10^8$	$1.13 \times 10^6$	50,688 ( $3.51 \times 10^9$ )	—	654 ( $8.70 \times 10^5$ )
SHHS-DIR	Health condition score	0 - 100	1	275	0	1,892	369	369

Table credit: [Yang et al. \(2021\)](#)

# (Training) Datasets - Label Distributions

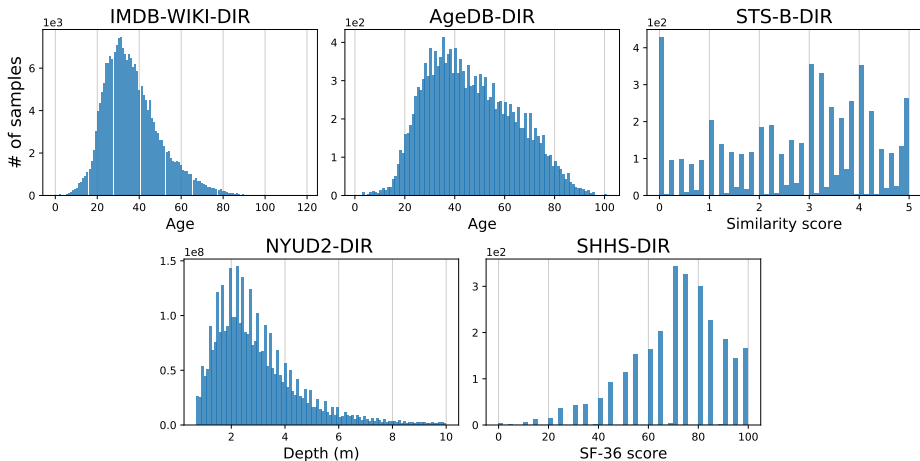
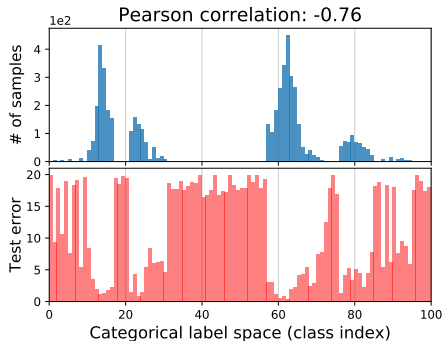


Image credit: [Yang et al. \(2021\)](#)

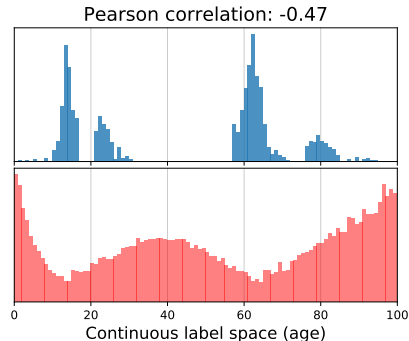
## Label Distribution Smoothing (LDS)

# Imbalanced Categorical vs. Continuous Label Space (1/3)



(a) Classification

- task: picture  $\longrightarrow$  class
- data source: CIFAR-100



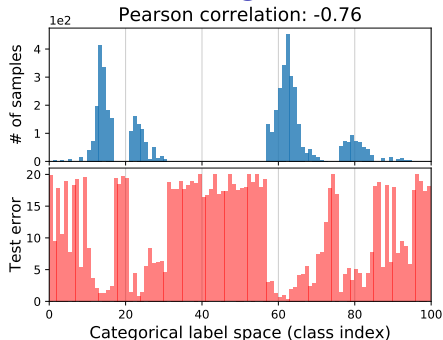
(b) Regression

- task:  
person's picture  $\longrightarrow$  person's age
- age subrange: 0-99
- data source: IMDB-WIKI
- Simulated label imbalance
- Label density distributions forced to be equal
- Balanced test sets

Image credit: [Yang et al. \(2021\)](#)

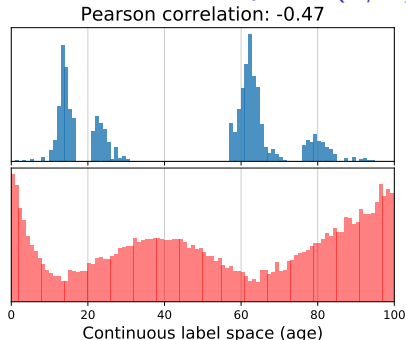


# Imbalanced Categorical vs. Continuous Label Space (2/3)



(a) Classification

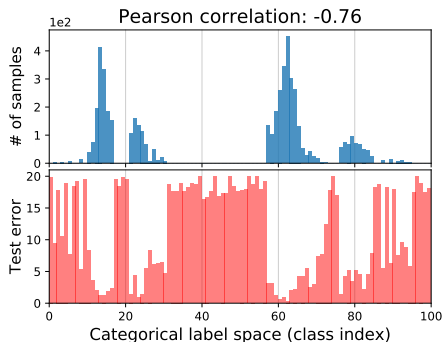
- the error distribution *correlates* with the label density distribution
- majority classes with more examples are better learned than minority classes



(b) Regression

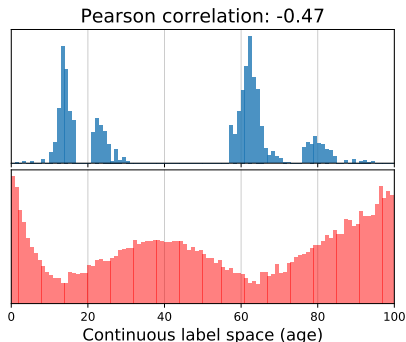
- the error distribution **DOES NOT** *correlate* well with the label density distribution
- smoother error distribution

# Imbalanced Categorical vs. Continuous Label Space (3/3)



(a) Classification

- Compensating for the imbalance in the empirical label density distribution WORKS WELL.

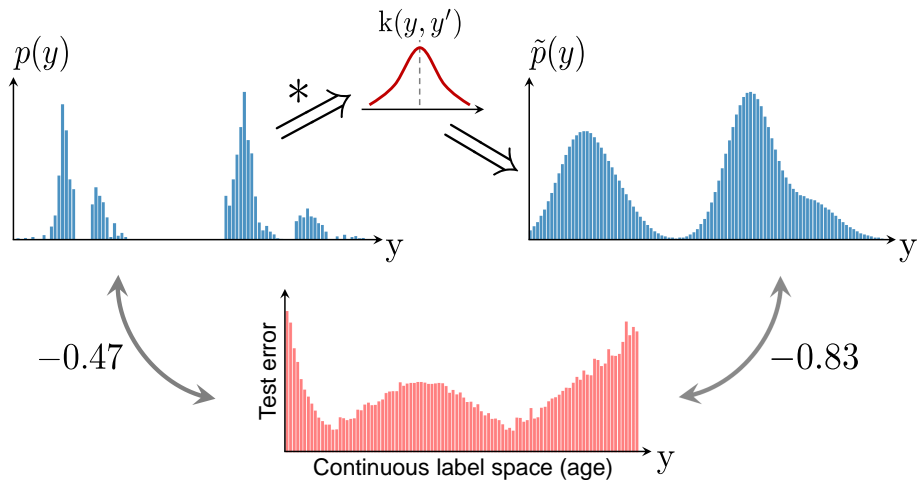


(b) Regression

- Compensating for the imbalance in the empirical label density distribution is INACCURATE.
- The empirical density does not accurately reflect the imbalance as seen by the model.
- Intuition: dependence between features at nearby labels.
- Proposed solution:  
**Label Distribution Smoothing (LDS)**

Image credit: [Yang et al. \(2021\)](#)

# Label Distribution Smoothing (LDS) - Overview



# Label Distribution Smoothing (LDS)

- Starting points
  - ▶ Dependence between features at nearby continuous labels
  - ▶ Expected density estimation
    - ★ Significant literature in statistics ([Parzen 1962](#))
    - ★ Kernel density estimation
- Functioning
  - ▶ Convolves a symmetric kernel with the empirical label density distribution.
  - ▶ Extracts a kernel-smoothed label density accounting for the feature overlap of neighbouring labels.
- Symmetric kernel
  - ▶ E.g., Gaussian or Laplacian kernel
  - ▶ Similarity between target values w.r.t. their distance in the target space.
- *Effective label density distribution*

$$\tilde{p}(y') \triangleq \int_{\mathcal{Y}} k(y, y') p(y) dy$$

where

- ▶  $p(y)$ : nr. occurrences of label  $y$  in training data
- How to use it in practice?
  - ▶ Possible direct adaptation of class imbalance techniques.
  - ▶ E.g., loss weighted by inverse effective label density

# Feature Distribution Smoothing (FDS) - Preliminaries

# Feature Distribution Smoothing (FDS) - Preliminaries

- Starting points

- 1 Continuity in the **target** space  $\longleftrightarrow$  Continuity in the **feature** space

# Feature Distribution Smoothing (FDS) - Preliminaries

- Starting points
  - ① Continuity in the **target** space  $\longleftrightarrow$  Continuity in the **feature** space
  - ② Data balance  $\implies$  close feature statistics of nearby targets

# Feature Distribution Smoothing (FDS) - Preliminaries

- Starting points
  - ① Continuity in the **target** space  $\longleftrightarrow$  Continuity in the **feature** space
  - ② Data balance  $\implies$  close feature statistics of nearby targets
- Feature statistics: mean and variance (or covariance) w.r.t. each bin

$$\{\boldsymbol{\mu}_b, \boldsymbol{\sigma}_b\}_{b=1}^B$$



# Feature Distribution Smoothing (FDS) - Preliminaries

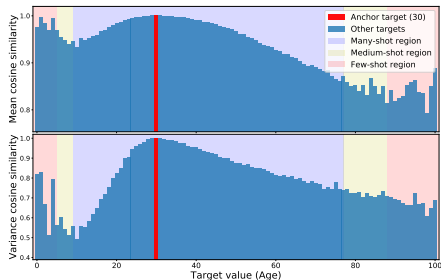
- Starting points
  - 1 Continuity in the **target** space  $\longleftrightarrow$  Continuity in the **feature** space
  - 2 Data balance  $\implies$  close feature statistics of nearby targets
- Feature statistics: mean and variance (or covariance) w.r.t. each bin

$$\{\boldsymbol{\mu}_b, \boldsymbol{\sigma}_b\}_{b=1}^B$$

- (next slides) Feature statistics similarity: cosine similarity of feature statistics between one anchor bin  $b_0$  and all other bins
  - ▶  $b_0 = 0, 30, 60, 90$  (age): chosen anchor bins
  - ▶ different target densities:  
many ( $>100$ ), medium (20-100), few ( $<20$ ) examples
  - ▶ task: person's picture  $\longrightarrow$  person's age
  - ▶ data source: IMDB-WIKI

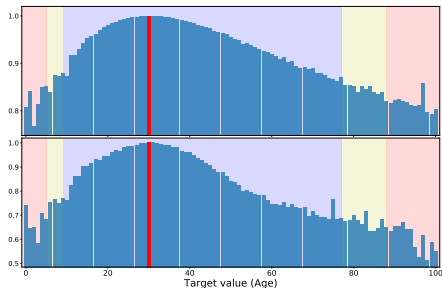
# Feature statistics similarity (1/4)

Anchor age 30



(a) Baseline

- High similarity in neighbourhood
- High similarities with further regions
- Lower similarities with some closer regions



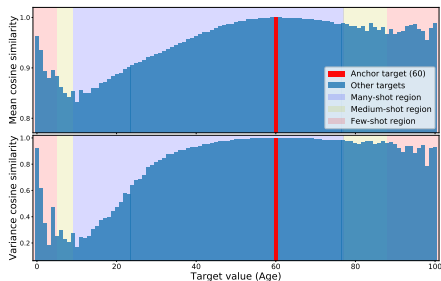
(b) FDS

- Improved feature statistics calibration:
  - ▶ High similarity only in neighbourhood
  - ▶ “The further the region the lower the similarity”
  - ▶ More gradual similarity change

Image credit: Yang et al. (2021)

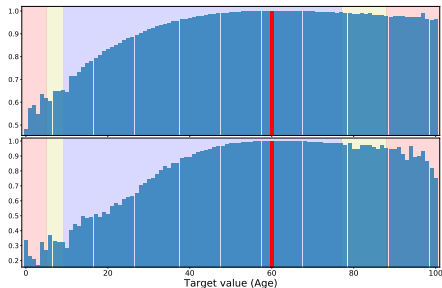
# Feature statistics similarity (2/4)

Anchor age 60



(a) Baseline

- High similarity in neighbourhood
- High similarities with further regions
- Lower similarities with some closer regions



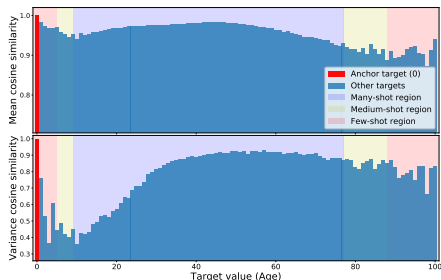
(b) FDS

- Improved feature statistics calibration:
  - ▶ High similarity only in neighbourhood
  - ▶ “The further the region the lower the similarity”
  - ▶ More gradual similarity change

Image credit: Yang et al. (2021)

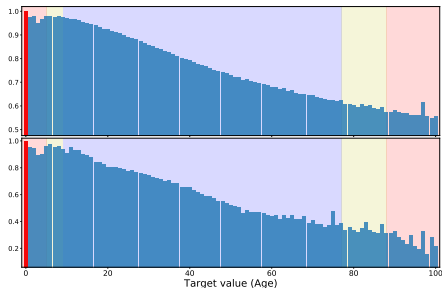
# Feature statistics similarity (3/4)

Anchor age 0



(a) Baseline

- High similarity in neighbourhood for mean
- High similarities with further regions
- Lower similarities with some closer regions, e.g., variance neighbourhood



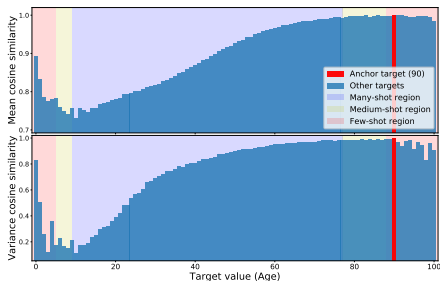
(b) FDS

- Improved feature statistics calibration:
  - ▶ High similarity only in neighbourhood
  - ▶ “The further the region the lower the similarity”
  - ▶ More gradual similarity change

Image credit: Yang et al. (2021)

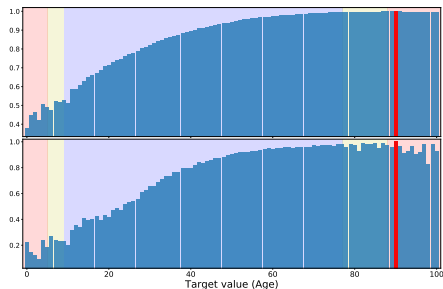
# Feature statistics similarity (4/4)

Anchor age 90



(a) Baseline

- High similarity in neighbourhood, esp. for mean
- High similarities with further regions
- Lower similarities with some closer regions

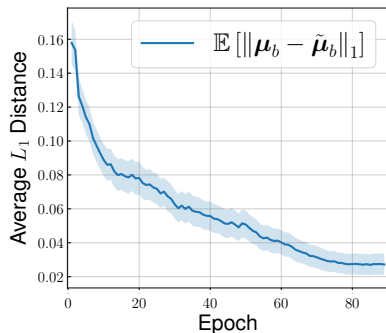


(b) FDS

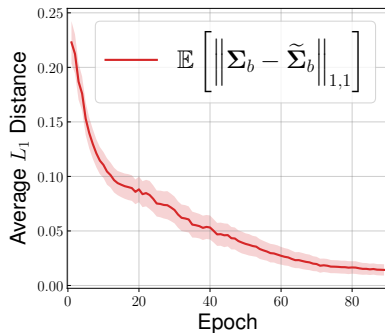
- Improved feature statistics calibration:
  - ▶ High similarity only in neighbourhood
  - ▶ “The further the region the lower the similarity”
  - ▶ More gradual similarity change

Image credit: Yang et al. (2021)

# Change of feature statistics w.r.t. epoch



(a) Mean



(b) Covariance

- $\mu, \Sigma$ : Running mean and covariance
- $\tilde{\mu}, \tilde{\Sigma}$ : Smoothed mean and covariance

# Feature Distribution Smoothing (FDS) - Overview

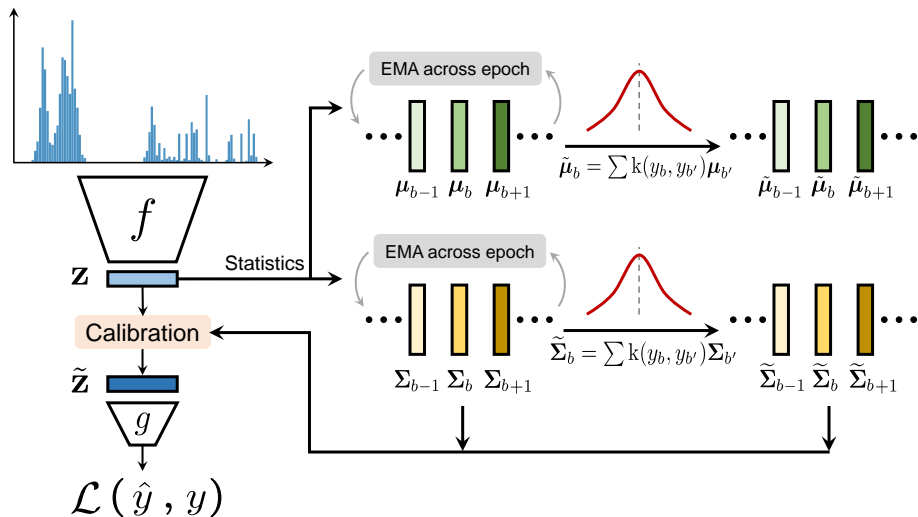


Image credit: Yang et al. (2021)

# Feature Distribution Smoothing (FDS)

- Transfers the feature statistics between nearby bins.
- Aim: calibrate the potentially biased estimates of feature distribution, esp. for underrepresented target values in training data.
- General functioning
  - ▶ Estimates mean  $\mu_b$  and covariance  $\Sigma_b$  feature statistics by each target bin.
  - ▶ Smooths the feature statistics over the target bins  $\mathcal{B}$  by a symmetric kernel  $k(y_b, y'_b)$ . Obtains the smoothed mean  $\tilde{\mu}_b$  and covariance  $\tilde{\Sigma}_b$  feature statistics.
  - ▶ Whitens features (Sun et al. 2016):

$$\mathbf{z}^w = \Sigma_b^{-\frac{1}{2}}(\mathbf{z} - \mu_b)$$

- ▶ Re-colors whitened features (Sun et al. 2016):

$$\mathbf{z}^r = \tilde{\Sigma}_b \mathbf{z}^w + \tilde{\mu}_b$$

- Integration into deep learning
  - ▶ Feature calibration layer after the final feature map.
  - ▶ Momentum update running statistics  $\{\mu_b, \Sigma_b\}$  across each epoch.
    - ★ Exponential Moving Average (EMA)
  - ▶ Smoothed statistics  $\{\tilde{\Sigma}_b, \tilde{\mu}_b\}$  updated across different but fixed within each training epoch.



# Baselines (1/2)

- Vanilla: neglects data imbalance
- Synthetic samples
  - ▶ SMOTER ([Torgo et al. 2013](#))
    - ① Defines frequent and rare regions using label density.
    - ② Creates synthetic samples for pre-defined rare regions by linearly interpolating both inputs and labels.
  - ▶ SMOGN ([Branco et al. 2017](#)): augments SMOTER with Gaussian noise
- Focal-R

$$\frac{1}{n} \sum_{i=1}^n \sigma(|\beta e_i|)^{\gamma} e_i$$

- ▶ Error-aware loss
- ▶ Maps the absolute error into  $[0, 1]$ .
- ▶  $e_i$ :  $L_1$  error for the  $i$ -th sample
- ▶  $\beta, \gamma$ : hyper-parameters
- ▶ Inspired by Focal Loss ([Lin 2017](#)) for classification

# Baselines (2/2)

- Regressor re-training (RRT)
  - ▶ Two-stage training
    - ① Train encoder
    - ② Re-train regressor with inverse re-weighting and frozen encoder.
  - ▶ Inspired by [Kang et al. 2019](#)
- Cost-sensitive re-weighting: re-weighting schemes based on label distribution
  - ▶ Inverse-frequency weighting (INV)
  - ▶ Square-root weighting variant (SQINV)

## Results

# Inferring Age from Images

IMDB-WIKI

Metrics	MAE ↓				GM ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	8.06	7.23	15.12	26.33	4.57	4.17	10.59	20.46
SMOTER (Torgo et al. 2013)	8.14	7.42	14.15	25.28	4.64	<b>4.30</b>	9.05	19.46
SMOGLN (Branco et al. 2017)	8.03	<b>7.30</b>	14.02	25.93	4.63	<b>4.30</b>	8.74	20.12
SMOGLN + LDS	8.02	7.39	13.71	23.22	4.63	4.39	8.71	15.80
SMOGLN + FDS	8.03	7.35	14.06	23.44	4.65	4.33	8.87	16.00
SMOGLN + LDS + FDS	<b>7.97</b>	7.38	<b>13.22</b>	<b>22.95</b>	<b>4.59</b>	4.39	<b>7.84</b>	<b>14.94</b>
FOCAL-R	7.97	7.12	15.14	26.96	4.49	4.10	10.37	21.20
FOCAL-R + LDS	7.90	<b>7.10</b>	14.72	25.84	<b>4.47</b>	<b>4.09</b>	10.11	19.14
FOCAL-R + FDS	7.96	7.14	14.71	26.06	4.51	4.12	10.16	19.56
FOCAL-R + LDS + FDS	<b>7.88</b>	<b>7.10</b>	<b>14.08</b>	<b>25.75</b>	<b>4.47</b>	4.11	<b>9.32</b>	<b>18.67</b>
RRT	7.81	7.07	14.06	25.13	4.35	4.03	8.91	16.96
RRT + LDS	7.79	7.08	13.76	24.64	4.34	<b>4.02</b>	8.72	16.92
RRT + FDS	<b>7.65</b>	<b>7.02</b>	12.68	23.85	<b>4.31</b>	4.03	7.58	16.28
RRT + LDS + FDS	<b>7.65</b>	7.06	<b>12.41</b>	<b>23.51</b>	<b>4.31</b>	4.07	<b>7.17</b>	<b>15.44</b>
SQINV	7.87	7.24	12.44	22.76	4.47	4.22	7.25	15.10
SQINV + LDS	7.83	7.31	<b>12.43</b>	22.51	4.42	4.19	7.00	13.94
SQINV + FDS	7.83	7.23	12.60	22.37	4.42	4.20	<b>6.93</b>	13.48
SQINV + LDS + FDS	<b>7.78</b>	<b>7.20</b>	12.61	<b>22.19</b>	<b>4.37</b>	<b>4.12</b>	7.39	<b>12.61</b>
Ours (best) vs. VANILLA	<b>+0.41</b>	<b>+0.21</b>	<b>+2.71</b>	<b>+4.14</b>	<b>+0.26</b>	<b>+0.15</b>	<b>+3.66</b>	<b>+7.85</b>

- Either LDS, FDS, or both leads to performance gains.
- LDS + FDS often achieves the best results:
  - ▶ maintains or improves performance overall and on many-shot regions,
  - ▶ boosts performance for medium-shot and few-shot regions.

Table credit: Yang et al. (2021)

# Inferring Age from Images

AgeDB

Metrics	MAE ↓				GM ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	7.77	6.62	9.55	13.67	5.05	4.23	7.01	10.75
SMOTER (Torgo et al. 2013)	8.16	7.39	8.65	12.28	5.21	4.65	5.69	8.49
SMOBN (Branco et al. 2017)	8.26	7.64	9.01	12.09	5.36	4.90	6.19	8.44
SMOBN + LDS	7.96	7.44	8.64	11.77	5.03	4.68	5.69	7.98
SMOBN + FDS	8.06	7.52	8.75	11.89	5.02	4.66	5.63	8.02
SMOBN + LDS + FDS	<b>7.90</b>	<b>7.32</b>	<b>8.51</b>	<b>11.19</b>	<b>4.98</b>	<b>4.64</b>	<b>5.41</b>	<b>7.35</b>
FOCAL-R	7.64	6.68	9.22	13.00	4.90	4.26	6.39	9.52
FOCAL-R + LDS	7.56	<b>6.67</b>	8.82	12.40	4.82	4.27	5.87	8.83
FOCAL-R + FDS	7.65	6.89	8.70	<b>11.92</b>	4.83	4.32	5.89	<b>8.04</b>
FOCAL-R + LDS + FDS	<b>7.47</b>	6.69	<b>8.30</b>	12.55	<b>4.71</b>	<b>4.25</b>	<b>5.36</b>	8.59
RRT	7.74	6.98	8.79	11.99	5.00	4.50	5.88	8.63
RRT + LDS	7.72	7.00	8.75	11.62	4.98	4.54	5.71	8.27
RRT + FDS	7.70	<b>6.95</b>	8.76	11.86	4.82	<b>4.32</b>	5.83	8.08
RRT + LDS + FDS	<b>7.66</b>	6.99	<b>8.60</b>	<b>11.32</b>	<b>4.80</b>	4.42	<b>5.53</b>	<b>6.99</b>
SQINV	7.81	7.16	8.80	11.20	4.99	4.57	5.73	7.77
SQINV + LDS	7.67	<b>6.98</b>	8.86	10.89	4.85	4.39	5.80	7.45
SQINV + FDS	7.69	7.10	8.86	<b>9.98</b>	4.83	4.41	5.97	<b>6.29</b>
SQINV + LDS + FDS	<b>7.55</b>	7.01	<b>8.24</b>	10.79	<b>4.72</b>	<b>4.36</b>	<b>5.45</b>	6.79
Ours (best) vs. VANILLA	<b>+0.30</b>	<b>-0.05</b>	<b>+1.31</b>	<b>+3.69</b>	<b>+0.34</b>	<b>-0.02</b>	<b>+1.65</b>	<b>+4.46</b>

- Either LDS, FDS, or both leads to performance gains.
- LDS + FDS often achieves the best results:
  - ▶ maintains or improves performance overall and on many-shot regions,
  - ▶ boosts performance for medium-shot and few-shot regions.

Table credit: Yang et al. (2021)

# Inferring Text Similarity Score

STS-B

Metrics	MSE ↓				Pearson correlation (%) ↑			
Shot	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	0.974	0.851	1.520	0.984	74.2	72.0	62.7	75.2
SMOTER (Torgo et al. 2013)	1.046	0.924	1.542	1.154	72.6	69.3	65.3	70.6
SMOEN (Branco et al. 2017)	0.990	0.896	1.327	1.175	73.2	70.4	65.5	69.2
SMOEN + LDS	0.962	0.880	1.242	1.155	74.0	71.5	65.2	69.8
SMOEN + FDS	0.987	0.945	<b>1.101</b>	1.153	73.0	69.6	<b>68.5</b>	69.9
SMOEN + LDS + FDS	<b>0.950</b>	<b>0.851</b>	1.327	<b>1.095</b>	<b>74.6</b>	<b>72.1</b>	65.9	<b>71.7</b>
FOCAL-R	0.951	0.843	1.425	0.957	74.6	72.3	61.8	76.4
FOCAL-R + LDS	0.930	<b>0.807</b>	1.449	0.993	<b>75.7</b>	<b>73.9</b>	62.4	75.4
FOCAL-R + FDS	<b>0.920</b>	0.855	<b>1.169</b>	1.008	75.1	72.6	<b>66.4</b>	74.7
FOCAL-R + LDS + FDS	0.940	0.849	1.358	<b>0.916</b>	74.9	72.2	66.3	<b>77.3</b>
RRT	0.964	0.842	1.503	0.978	74.5	72.4	62.3	75.4
RRT + LDS	0.916	0.817	1.344	0.945	75.7	73.5	64.1	76.6
RRT + FDS	0.929	0.857	<b>1.209</b>	1.025	74.9	72.1	<b>67.2</b>	74.0
RRT + LDS + FDS	<b>0.903</b>	<b>0.806</b>	1.323	<b>0.936</b>	<b>76.0</b>	<b>73.8</b>	65.2	<b>76.7</b>
INV	1.005	0.894	1.482	1.046	72.8	70.3	62.5	73.2
INV + LDS	0.914	0.819	1.319	0.955	75.6	73.4	63.8	76.2
INV + FDS	0.927	0.851	<b>1.225</b>	1.012	75.0	72.4	<b>66.6</b>	74.2
INV + LDS + FDS	<b>0.907</b>	<b>0.802</b>	1.363	<b>0.942</b>	<b>76.0</b>	<b>74.0</b>	65.2	<b>76.6</b>
Ours (best) vs. VANILLA	<b>+0.071</b>	<b>+0.049</b>	<b>+0.419</b>	<b>+0.068</b>	<b>+1.8</b>	<b>+2.0</b>	<b>+5.8</b>	<b>+2.1</b>

- Both LDS and FDS improve results for various methods, esp. medium- and few-shot regions.

Table credit: Yang et al. (2021)

# Inferring Depth

NYUD2

Metrics	RMSE ↓				$\delta_1$ ↑			
Shot	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	1.477	0.591	0.952	2.123	0.677	0.777	0.693	0.570
VANILLA + LDS	1.387	0.671	0.913	1.954	0.672	0.701	0.706	0.630
VANILLA + FDS	1.442	<b>0.615</b>	0.940	2.059	0.681	<b>0.760</b>	0.695	0.596
VANILLA + LDS + FDS	<b>1.338</b>	0.670	<b>0.851</b>	<b>1.880</b>	<b>0.705</b>	0.730	<b>0.764</b>	<b>0.655</b>
Ours (best) vs. VANILLA	<b>+0.139</b>	<b>-0.024</b>	<b>+0.101</b>	<b>+0.243</b>	<b>+0.028</b>	<b>-0.017</b>	<b>+0.071</b>	<b>+0.085</b>

## FDS and LDS

- alleviates overfitting on many-shot regions,
- generalizes better to all regions,
- slightly degrades many-shot region,
- boosts other regions.

Table credit: Yang et al. (2021)

# Inferring Health Score

## SHHS-DIR

Metrics	MAE ↓				GM ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	15.36	12.47	13.98	16.94	10.63	8.04	9.59	12.20
FOCAL-R	14.67	11.70	13.69	17.06	9.98	7.93	8.85	11.95
FOCAL-R + LDS	14.49	12.01	12.43	16.57	9.98	7.89	8.59	11.40
FOCAL-R + FDS	14.18	<b>11.06</b>	13.56	15.99	9.45	<b>6.95</b>	8.81	11.13
FOCAL-R + LDS + FDS	<b>14.02</b>	11.08	<b>12.24</b>	<b>15.49</b>	<b>9.32</b>	7.18	<b>8.10</b>	<b>10.39</b>
RRT	14.78	12.43	14.01	16.48	10.12	8.05	9.71	11.96
RRT + LDS	14.56	12.08	13.44	16.45	9.89	7.85	9.18	11.82
RRT + FDS	14.36	11.97	13.33	16.08	9.74	7.54	9.20	11.31
RRT + LDS + FDS	<b>14.33</b>	<b>11.96</b>	<b>12.47</b>	<b>15.92</b>	<b>9.63</b>	<b>7.35</b>	<b>8.74</b>	<b>11.17</b>
INV	14.39	11.84	13.12	16.02	9.34	7.73	8.49	11.20
INV + LDS	14.14	11.66	12.77	16.05	9.26	7.64	8.18	11.32
INV + FDS	13.91	<b>11.12</b>	12.29	15.53	8.94	<b>6.91</b>	7.79	10.65
INV + LDS + FDS	<b>13.76</b>	<b>11.12</b>	<b>12.18</b>	<b>15.07</b>	<b>8.70</b>	6.94	<b>7.60</b>	<b>10.18</b>
Ours (best) vs. VANILLA	<b>+1.60</b>	<b>+1.41</b>	<b>+1.80</b>	<b>+1.87</b>	<b>+1.93</b>	<b>+1.13</b>	<b>+1.99</b>	<b>+2.02</b>

- Both FDS and LDS are effective.
- FDS + LDS often get highest gains over all tested regions.
- Note: SMOTER and SMOGN not directly applicable.

Table credit: Yang et al. (2021)



# Ablation: kernel type

IMDB-WIKI

Metrics	MSE ↓				MAE ↓				GM ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	138.06	108.70	366.09	964.92	8.06	7.23	15.12	26.33	4.57	4.17	10.59	20.46
<b>LDS:</b>												
GAUSSIAN KERNEL	131.65	109.04	298.98	834.08	7.83	7.31	12.43	22.51	4.42	4.19	7.00	13.94
TRIANGULAR KERNEL	133.77	110.24	309.70	850.74	7.89	7.30	12.72	22.80	4.50	4.24	7.75	14.91
LAPLACIAN KERNEL	132.87	109.27	312.10	829.83	7.87	7.29	12.68	22.38	4.50	4.26	7.29	13.71
<b>FDS:</b>												
GAUSSIAN KERNEL	133.81	107.51	332.90	916.18	7.85	7.18	13.35	24.12	4.47	4.18	8.18	15.18
TRIANGULAR KERNEL	134.09	110.49	301.18	927.99	7.97	7.41	12.20	23.99	4.64	4.41	7.06	14.28
LAPLACIAN KERNEL	133.00	104.26	352.95	968.62	8.05	7.25	14.78	26.16	4.71	4.33	10.19	19.09

- All kernel types lead to gains
- Often best results with Gaussian kernel

Table credit: Yang et al. (2021)

# Ablation: kernel type

STS-B

Metrics	MSE ↓				MAE ↓				Pearson correlation (%) ↑				Spearman correlation (%) ↑			
Shot	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA	0.974	0.851	1.520	0.984	0.794	0.740	1.043	0.771	74.2	72.0	62.7	75.2	74.4	68.8	50.5	75.0
<b>LDS:</b>																
GAUSSIAN KERNEL	0.914	0.819	1.319	0.955	0.773	0.729	0.970	0.772	75.6	73.4	63.8	76.2	76.1	70.4	55.6	74.3
TRIANGULAR KERNEL	0.938	0.870	1.193	1.039	0.786	0.754	0.929	0.784	74.8	72.4	64.1	74.0	75.2	69.3	54.1	73.9
LAPLACIAN KERNEL	0.938	0.829	1.413	0.962	0.782	0.731	1.014	0.773	75.7	73.0	65.8	76.5	76.0	70.0	52.3	75.2
<b>FDS:</b>																
GAUSSIAN KERNEL	0.916	0.875	1.027	1.086	0.767	0.746	0.840	0.811	75.5	73.0	67.0	72.8	75.8	69.9	54.4	72.0
TRIANGULAR KERNEL	0.935	0.863	1.239	0.966	0.762	0.725	0.912	0.788	74.6	72.4	64.8	75.9	74.4	69.1	48.4	75.4
LAPLACIAN KERNEL	0.925	0.843	1.247	1.020	0.771	0.733	0.929	0.800	75.0	72.6	64.7	74.2	75.4	70.1	53.5	73.5

- All kernel types lead to gains
- Often best results with Gaussian kernel

Table credit: [Yang et al. \(2021\)](#)

# Ablation: Gaussian kernel hyper-parameters

IMDB-WIKI

Metrics		MSE ↓				MAE ↓				GM ↓			
Shot		All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA		138.06	108.70	366.09	964.92	8.06	7.23	15.12	26.33	4.57	4.17	10.59	20.46
$l$   $\sigma$													
LDS:													
5	1	132.08	108.53	309.03	843.53	7.80	7.22	12.61	22.33	4.42	4.19	7.16	12.54
9	1	135.04	112.32	307.90	803.15	7.97	7.39	12.74	22.19	4.55	4.30	7.53	14.11
15	1	134.06	110.49	308.83	864.30	7.84	7.28	12.35	22.81	4.44	4.22	6.95	14.22
5	2	131.65	109.04	298.98	834.08	7.83	7.31	12.43	22.51	4.42	4.19	7.00	13.94
9	2	136.78	112.41	322.65	850.47	8.02	7.41	13.00	23.23	4.55	4.29	7.55	15.65
15	2	135.66	111.68	319.20	833.02	7.98	7.40	12.74	22.27	4.60	4.37	7.30	12.92
5	3	137.56	113.50	322.47	831.38	8.07	7.47	13.06	22.85	4.63	4.36	7.87	15.11
9	3	138.91	114.89	319.40	863.16	8.18	7.57	13.19	23.33	4.71	4.44	8.09	15.17
15	3	138.86	114.25	326.97	856.27	8.18	7.54	13.53	23.17	4.77	4.47	8.52	15.25
FDS:													
5	1	133.63	104.80	354.24	972.54	7.87	7.06	14.71	25.96	4.42	4.04	9.95	18.47
9	1	134.34	105.97	356.54	919.16	7.95	7.18	14.58	24.80	4.54	4.20	9.56	15.13
15	1	136.32	107.47	355.84	948.71	7.97	7.23	14.81	25.59	4.60	4.23	9.99	17.60
5	2	133.81	107.51	332.90	916.18	7.85	7.18	13.35	24.12	4.47	4.18	8.18	15.18
9	2	133.99	105.01	357.31	963.79	7.94	7.11	14.95	25.97	4.48	4.09	10.49	18.19
15	2	136.61	107.93	361.08	973.56	7.98	7.23	14.68	25.21	4.61	4.24	10.14	17.91
5	3	136.81	107.76	359.08	953.16	7.98	7.18	14.85	24.94	4.53	4.15	10.27	17.33
9	3	133.48	104.14	359.80	972.29	7.94	7.09	15.04	25.87	4.48	4.09	10.40	16.85
15	3	132.55	103.08	360.39	970.43	8.03	7.22	14.86	25.40	4.67	4.33	10.04	13.86

- Gaussian kernel size  $l \in \{5, 9, 15\}$  and standard deviation  $\sigma \in \{1, 2, 3\}$
- LDS
  - ▶ Smaller  $\sigma$  usually leads to slightly better results over all regions.
  - ▶ Larger gains w.r.t. the performance in medium-shot and few-shot regions.
  - ▶ Minor degradation in many-shot regions.
- FDS
  - ▶ Smaller  $l$  often obtains slightly higher improvements over all regions.
  - ▶ Equally boosts all the regions, with slightly smaller improvements in medium-shot and few-shot regions.
- 3.3-6.2% overall MSE gain
- Best results with  $l = 5$  and  $\sigma = 2$
- Robust to different hyper-parameters

Table credit: Yang et al. (2021)

# Ablation: Gaussian kernel hyper-parameters

STS-B

Metrics		MSE ↓				MAE ↓				Pearson correlation (%) ↑				Spearman correlation (%) ↑			
Shot		All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
VANILLA		0.974	0.851	1.520	0.984	0.794	0.740	1.043	0.771	74.2	72.0	62.7	75.2	74.4	68.8	50.5	75.0
$l$   $\sigma$																	
<b>LDS:</b>																	
5	1	0.942	0.825	1.431	1.023	0.781	0.726	1.016	0.809	75.1	73.2	61.8	74.5	75.3	70.2	52.2	72.5
9	1	0.931	0.840	1.323	0.962	0.785	0.744	0.972	0.773	75.0	72.7	63.3	75.8	75.6	70.1	53.6	74.8
15	1	0.941	0.833	1.413	0.953	0.781	0.728	1.014	0.776	75.0	72.8	62.6	76.3	75.5	70.2	52.0	74.6
5	2	0.914	0.819	1.319	0.955	0.773	0.729	0.970	0.772	75.6	73.4	63.8	76.2	76.1	70.4	55.6	74.3
9	2	0.926	0.823	1.379	0.944	0.782	0.733	1.003	0.764	75.5	73.4	63.6	76.8	76.0	70.5	53.5	76.2
15	2	0.949	0.831	1.452	1.005	0.788	0.735	1.023	0.782	74.9	72.9	63.0	74.7	75.4	70.1	52.5	73.6
5	3	0.928	0.845	1.250	1.041	0.775	0.733	0.951	0.798	75.1	73.3	63.2	73.8	75.3	70.4	51.4	72.6
9	3	0.939	0.816	1.462	1.000	0.786	0.732	1.030	0.783	75.3	73.5	62.6	74.7	75.9	70.9	53.0	73.7
15	3	0.927	0.824	1.348	1.010	0.774	0.726	0.982	0.780	75.2	73.4	62.2	74.6	75.7	70.7	53.0	72.3
<b>FDS:</b>																	
5	1	0.943	0.869	1.217	1.066	0.776	0.742	0.914	0.799	74.4	71.7	65.6	72.5	74.2	68.4	51.1	71.2
9	1	0.927	0.851	1.193	1.096	0.770	0.736	0.896	0.822	74.9	72.8	65.8	71.6	74.8	69.7	52.3	68.3
15	1	0.926	0.854	1.202	1.029	0.776	0.743	0.914	0.800	74.9	72.6	66.1	74.0	75.1	69.8	49.5	73.6
5	2	0.916	0.875	1.027	1.086	0.767	0.746	0.840	0.811	75.5	73.0	67.0	72.8	75.8	69.9	54.4	72.0
9	2	0.933	0.888	1.068	1.081	0.776	0.752	0.855	0.839	74.8	72.0	67.9	72.2	74.9	68.9	53.3	72.0
15	2	0.944	0.890	1.125	1.078	0.783	0.761	0.864	0.822	74.4	71.8	65.8	72.2	74.5	68.9	53.1	70.9
5	3	0.924	0.860	1.190	0.964	0.771	0.740	0.897	0.790	75.0	72.7	64.4	76.1	75.1	69.4	53.8	76.5
9	3	0.932	0.878	1.149	0.982	0.770	0.746	0.876	0.780	74.8	72.5	63.8	75.3	74.8	69.3	50.2	75.6
15	3	0.956	0.915	1.110	1.016	0.784	0.767	0.855	0.803	74.4	72.1	63.7	75.5	74.3	68.7	50.0	74.6

- Gaussian kernel size  $l \in \{5, 9, 15\}$  and standard deviation  $\sigma \in \{1, 2, 3\}$
- 3.3-6.2% overall MSE gain
- Best results with  $l = 5$  and  $\sigma = 2$
- Robust to different hyper-parameters

Table credit: Yang et al. (2021)

# Ablation: loss function

## STS-B

Metrics	MSE ↓				MAE ↓				Pearson correlation (%) ↑				Spearman correlation (%) ↑			
Shot	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
<b>LDS:</b>																
MAE (L1)	0.893	0.808	1.241	0.964	0.765	0.727	0.938	0.758	76.3	73.9	66.0	75.9	76.7	71.1	54.5	75.6
MSE (L2)	0.914	0.819	1.319	0.955	0.773	0.729	0.970	0.772	75.6	73.4	63.8	76.2	76.1	70.4	55.6	74.3
HUBER LOSS (sL1)	0.902	0.811	1.276	0.978	0.761	0.718	0.954	0.751	76.1	74.2	64.7	75.5	76.5	71.6	52.9	74.3
<b>FDS:</b>																
MAE (L1)	0.918	0.860	1.105	1.082	0.762	0.733	0.859	0.833	75.5	73.7	65.3	72.3	75.6	70.9	52.1	71.5
MSE (L2)	0.916	0.875	1.027	1.086	0.767	0.746	0.840	0.811	75.5	73.0	67.0	72.8	75.8	69.9	54.4	72.0
HUBER LOSS (sL1)	0.920	0.867	1.097	1.052	0.765	0.741	0.858	0.800	75.3	72.9	66.6	73.6	75.3	69.7	52.3	73.6

- Similar results for all losses
- Robust to different losses

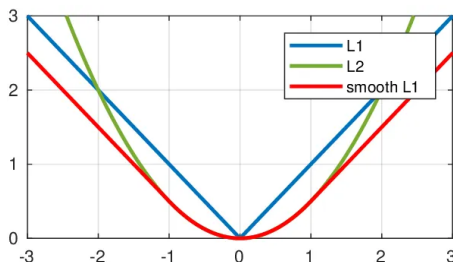


Table credit: [Yang et al. \(2021\)](#), Image credit: <https://medium.com/artificialis/loss-functions-361b2ad439a>

# Could LDS + FDS help when the label distribution is skewed with one or more Gaussian peaks?

- Experimental setup

- ▶ Curated skewed label distributions with 1-4 Gaussian peaks on IMDB-WIKI-DIR
- ▶ Compared with the vanilla model

- Findings

- ▶ Robustness to distribution change
- ▶ Brings improvement

# Skewed label distribution with one Gaussian peak

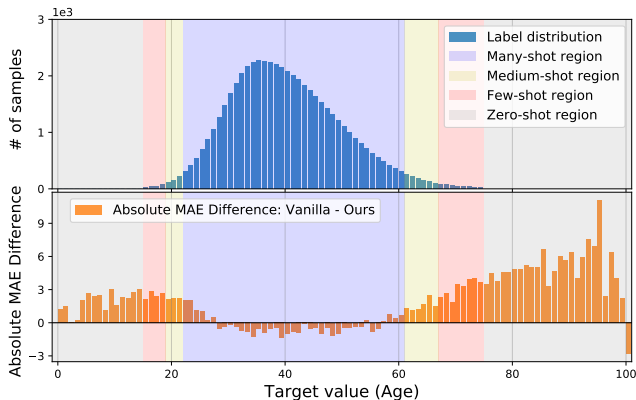


Figure: MAE gains of LDS + FDS over the vanilla model.

- Performance gains, esp. for extrapolation & interpolation

# Skewed label distribution with two Gaussian peaks

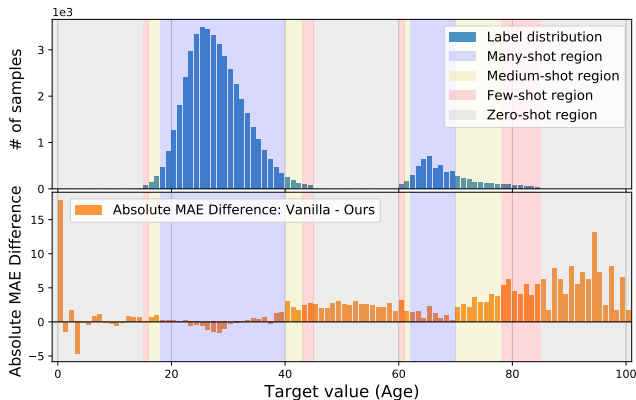


Figure: MAE gains of LDS + FDS over the vanilla model.

- Performance gains, esp. for extrapolation & interpolation



# Skewed label distribution with three Gaussian peaks

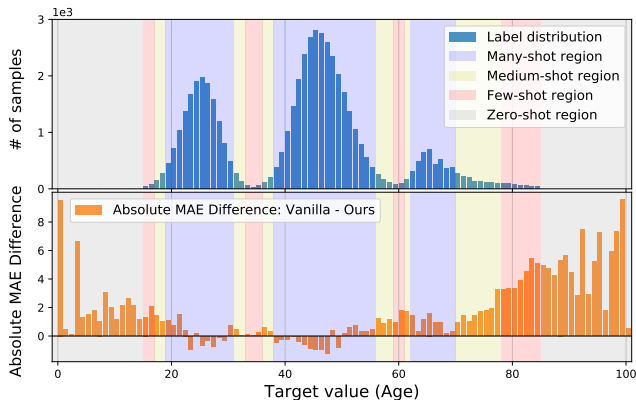


Figure: MAE gains of LDS + FDS over the vanilla model.

- Performance gains, esp. for extrapolation & interpolation

# Skewed label distribution with four Gaussian peaks

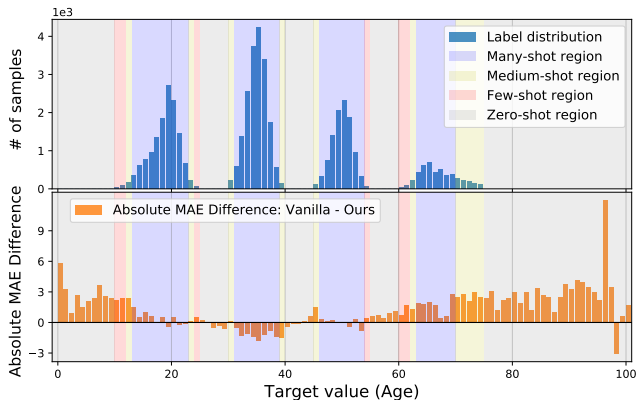


Figure: MAE gains of LDS + FDS over the vanilla model.

- Performance gains, esp. for extrapolation & interpolation

# Skewed label distribution with two Gaussian peaks on IMDB-WIKI-DIR

Metrics	MAE ↓				GM ↓			
Shot	All	w/ data	Interp.	Extrap.	All	w/ data	Interp.	Extrap.
VANILLA	11.72	9.32	16.13	18.19	7.44	5.33	14.41	16.74
VANILLA + <b>LDS</b>	10.54	8.31	14.14	17.38	6.50	4.67	12.13	15.36
VANILLA + <b>FDS</b>	11.40	8.97	15.83	18.01	7.18	5.12	14.02	16.48
VANILLA + <b>LDS</b> + <b>FDS</b>	<b>10.27</b>	<b>8.11</b>	<b>13.71</b>	<b>17.02</b>	<b>6.33</b>	<b>4.55</b>	<b>11.71</b>	<b>15.13</b>
<b>Ours (best)</b> vs. VANILLA	<b>+1.45</b>	<b>+1.21</b>	<b>+2.42</b>	<b>+1.17</b>	<b>+1.11</b>	<b>+0.78</b>	<b>+2.70</b>	<b>+1.61</b>

Table: Interpolation & extrapolation results

- Best results by smoothing both label & feature distributions

# Different skewed label distributions on IMDB-WIKI-DIR

Metrics	MAE ↓							GM ↓						
Shot	All	Many	Med.	Few	Zero	Interp.	Extrap.	All	Many	Med.	Few	Zero	Interp.	Extrap.
<b>1 peak:</b>														
VANILLA	11.20	6.05	11.43	14.76	22.67	—	22.67	7.02	<b>3.84</b>	8.67	12.26	21.07	—	21.07
VANILLA + LDS	10.09	6.26	9.91	12.12	19.37	—	19.37	6.14	3.92	6.50	8.30	16.35	—	16.35
VANILLA + FDS	11.04	<b>5.97</b>	11.19	14.54	22.35	—	22.35	6.96	<b>3.84</b>	8.54	12.08	20.71	—	20.71
VANILLA + LDS + FDS	<b>10.00</b>	6.28	<b>9.66</b>	<b>11.83</b>	<b>19.21</b>	—	<b>19.21</b>	<b>6.09</b>	3.96	<b>6.26</b>	<b>8.14</b>	<b>15.89</b>	—	<b>15.89</b>
<b>2 peaks:</b>														
VANILLA	11.72	6.83	11.78	15.35	16.86	16.13	18.19	7.44	3.61	8.06	12.94	15.21	14.41	16.74
VANILLA + LDS	10.54	6.72	9.65	12.60	15.30	14.14	17.38	6.50	3.65	<b>5.65</b>	9.30	13.20	12.13	15.36
VANILLA + FDS	11.40	6.69	11.02	14.85	16.61	15.83	18.01	7.18	<b>3.50</b>	7.49	12.73	14.86	14.02	16.48
VANILLA + LDS + FDS	<b>10.27</b>	<b>6.61</b>	<b>9.46</b>	<b>11.96</b>	<b>14.89</b>	<b>13.71</b>	<b>17.02</b>	<b>6.33</b>	3.54	5.68	<b>8.80</b>	<b>12.83</b>	<b>11.71</b>	<b>15.13</b>
<b>3 peaks:</b>														
VANILLA	9.83	7.01	9.81	11.93	20.11	—	20.11	6.04	3.93	6.94	9.84	17.77	—	17.77
VANILLA + LDS	9.08	<b>6.77</b>	8.82	10.48	18.43	—	18.43	<b>5.35</b>	<b>3.78</b>	5.63	7.49	15.46	—	15.46
VANILLA + FDS	9.65	6.88	9.58	11.75	19.80	—	19.80	5.86	3.83	6.68	9.48	17.43	—	17.43
VANILLA + LDS + FDS	<b>8.96</b>	6.88	<b>8.62</b>	<b>10.08</b>	<b>17.76</b>	—	<b>17.76</b>	5.38	3.90	<b>5.61</b>	<b>7.36</b>	<b>14.65</b>	—	<b>14.65</b>
<b>4 peaks:</b>														
VANILLA	9.49	7.23	9.73	10.85	12.16	8.23	18.78	5.68	3.45	6.95	8.20	9.43	6.89	16.02
VANILLA + LDS	8.80	<b>6.98</b>	8.26	10.07	11.26	8.31	<b>16.22</b>	5.10	<b>3.33</b>	<b>5.07</b>	7.08	8.47	6.66	<b>12.74</b>
VANILLA + FDS	9.28	7.11	9.16	10.88	11.95	8.30	18.11	5.49	3.36	6.35	8.15	9.21	6.82	15.30
VANILLA + LDS + FDS	<b>8.76</b>	7.07	<b>8.23</b>	<b>9.54</b>	<b>11.13</b>	<b>8.05</b>	16.32	<b>5.05</b>	3.36	<b>5.07</b>	<b>6.56</b>	<b>8.30</b>	<b>6.34</b>	13.10

Table credit: Yang et al. (2021)

# Balanced vs. Imbalanced Test Label Distribution

## IMDB-WIKI

Metrics	MSE ↓				MAE ↓				GM ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
<b>Balanced:</b>												
VANILLA	138.06	108.70	366.09	964.92	8.06	7.23	15.12	26.33	4.57	4.17	10.59	20.46
VANILLA + LDS + FDS	<b>129.35</b>	<b>106.52</b>	<b>311.49</b>	<b>811.82</b>	<b>7.78</b>	<b>7.20</b>	<b>12.61</b>	<b>22.19</b>	<b>4.37</b>	<b>4.12</b>	<b>7.39</b>	<b>12.61</b>
<b>Same as training set:</b>												
VANILLA	<b>68.44</b>	<b>62.10</b>	320.52	1350.01	<b>5.84</b>	<b>5.72</b>	15.11	30.54	<b>3.44</b>	<b>3.40</b>	11.76	24.06
VANILLA + LDS + FDS	69.86	63.43	<b>161.97</b>	<b>1067.89</b>	5.90	5.77	<b>9.94</b>	<b>25.17</b>	3.48	3.44	<b>7.03</b>	<b>15.95</b>

- Skewed label distribution for training set
- Case: balanced label distribution for test set.
  - ▶ LDS and FDS can improve the performance of all the regions.
- Case: skewed label distribution for test set, same label distribution for training set.
  - ▶ Minor degradation in many-shot region.
  - ▶ Boosts in medium-shot and few-shot regions.
  - ▶ Note: overall performance dominated by many-shot region, potentially biased and undesired evaluation.

Table credit: Yang et al. (2021)

# Comparison to imbalanced classification methods

Dataset	IMDB-WIKI-DIR (subsampling)				STS-B-DIR				NYUD2-DIR			
Metric	MAE ↓				MSE ↓				RMSE ↓			
Shot	All	Many	Med.	Few	All	Many	Med.	Few	All	Many	Med.	Few
<b>Imbalanced Classification:</b>												
CLS-VANILLA	15.94	15.64	18.95	30.21	1.926	1.906	2.022	1.907	1.576	0.596	1.011	2.275
CB (Cui et al. 2019)	22.41	22.32	22.05	32.90	2.159	2.194	2.028	2.107	1.664	0.592	1.044	2.415
CRT (Kang et al. 2019)	15.65	15.33	17.52	29.54	1.891	1.906	1.930	1.650	1.488	0.659	1.032	2.107
<b>Imbalanced Regression:</b>												
REG-VANILLA	14.64	13.98	17.47	30.29	0.974	0.851	1.520	0.984	1.477	<b>0.591</b>	0.952	2.123
LDS	14.03	13.72	15.93	26.71	0.914	0.819	1.319	0.955	1.387	0.671	0.913	1.954
FDS	13.97	13.55	16.42	24.64	0.916	0.875	<b>1.027</b>	1.086	1.442	0.615	0.940	2.059
LDS + FDS	<b>13.32</b>	<b>13.14</b>	<b>15.06</b>	<b>23.87</b>	<b>0.907</b>	<b>0.802</b>	1.363	<b>0.942</b>	<b>1.338</b>	0.670	<b>0.851</b>	<b>1.880</b>

- Imbalanced regression methods outperform classification ones.
- Can reduce error up to 50-60% in few-shot regions
- Imbalanced classification methods can perform worse than vanilla regression.
- Main finding: imbalance regression requires something different than just imbalance classification methods, which
  - ▶ can ignore similarity between nearby targets,
  - ▶ can ignore similarity between features linked to nearby targets,
  - ▶ cannot interpolate & extrapolate in the continuous label space, so cannot deal with zero-shot label regions.

Table credit: Yang et al. (2021)

# References I

- Branco, Paula, Luís Torgo, and Rita P Ribeiro (2017). "SMOGL: a pre-processing approach for imbalanced regression". In: *First international workshop on learning with imbalanced domains: Theory and applications*. PMLR, pp. 36–50.
- Cui, Yin et al. (2019). "Class-balanced loss based on effective number of samples". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9268–9277.
- Kang, Bingyi et al. (2019). "Decoupling representation and classifier for long-tailed recognition". In: *arXiv preprint arXiv:1910.09217*.
- Lin, T (2017). "Focal Loss for Dense Object Detection". In: *arXiv preprint arXiv:1708.02002*.
- Liu, Ziwei et al. (2019). "Large-scale long-tailed recognition in an open world". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2537–2546.

# References II

- Parzen, Emanuel (1962). “On estimation of a probability density function and mode”. In: *The annals of mathematical statistics* 33.3, pp. 1065–1076.
- Sun, Baochen, Jiashi Feng, and Kate Saenko (2016). “Return of frustratingly easy domain adaptation”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 30. 1.
- Torgo, Luís et al. (2013). “Smote for regression”. In: *Portuguese conference on artificial intelligence*. Springer, pp. 378–389.
- Yang, Yuzhe et al. (2021). “Delving into deep imbalanced regression”. In: *International conference on machine learning*. PMLR, pp. 11842–11851.