

1080 A. Connection to Example Difficulty

1081 A related line of work is identifying *difficult* examples,
 1082 or how well a model can generalize to a given unseen example.
 1083 Jiang *et al.* [20] introduce a per-instance empirical
 1084 consistency score which estimates the probability of predicting
 1085 the ground truth label with models trained on data sub-
 1086 samples of different sizes. Unlike our approach, however,
 1087 this requires training a large number of models. Toneva *et*
 1088 *al.* [32] quantifies example difficulty through the lens of a
 1089 forgetting event, in which the example is misclassified after
 1090 being correctly classified. However, the metrics that we
 1091 introduce in § 3, are based on the disagreement of the label
 1092 at each checkpoint with the final predicted label. Other
 1093 approaches estimate the example difficulty by: prediction
 1094 depth of the first layer at which a k -NN classifier correctly
 1095 classifies an example [2], the impact of pruning on model
 1096 predictions of the example [17], and estimating the leave-
 1097 one-out influence of each training example on the accuracy
 1098 of an algorithm by using influence functions [7]. Closest
 1099 to our method, the work of Agarwal *et al.* [1] utilizes
 1100 gradients of intermediate models during training to rank
 1101 examples by difficulty. In particular, they average pixel-wise
 1102 variance of gradients for each given input image. Notably,
 1103 this approach is more costly and less practical than our
 1104 approach and also does not study the accuracy/coverage trade-
 1105 off which is of paramount importance to SC.

1106 B. Alternative Metric Choices

1107 We briefly discuss additional potential metric choices
 1108 that we investigated but which lead to SC performance
 1109 worse than s_{avg} . A more elaborate discussion of these re-
 1110 sults is provided in Appendix C.3.

1111 B.1. Incorporating Estimates of e_t into s_{min} and s_{avg}

1112 Since the results in § 4.3 show that e_t is only nearly 0 almost
 1113 everywhere, we investigate whether incorporating an
 1114 estimate of e_t into s_{min} and s_{avg} leads to additional SC im-
 1115 provements. Our results demonstrate that neither an empirical
 1116 estimate nor a smooth decay function similar to v_t robustly
 1117 improves over s_{min} and s_{avg} (which assumed $e_t = 0$).

1118 B.2. Jump Score s_{jmp}

1119 We also consider a score which captures the level of dis-
 1120 agreement between the predicted label of two successive in-
 1121 termediate models (*i.e.* how much jumping occurred over
 1122 the course of training). For $j_t = 0$ iff $f_t(\mathbf{x}) = f_{t-1}(\mathbf{x})$
 1123 and $j_t = 1$ otherwise we can compute the jump score as

$$s_{\text{jmp}} = 1 - \sum v_t j_t$$
 and threshold it as in § 3.2 and § 3.3.

1124 B.3. Variance Score s_{var} for Continuous Metrics

1125 Finally, we consider monitoring the evolution of continuous
 1126 metrics that have been shown to be correlated with ex-

1127 ample difficulty. More specifically, Jiang *et al.* [20] show
 1128 that example difficulty is correlated with confidence and
 1129 negative entropy. Moreover, they find that difficult exam-
 1130 ples are learned later in the training process. This obser-
 1131 vation motivates designing a score based on these continu-
 1132 ous metrics that penalises changes later in the training pro-
 1133 cess more heavily. We consider the maximum softmax class
 1134 probability known as confidence, the negative entropy and
 1135 the gap between the most confident classes for each exam-
 1136 ple instead of the model predictions. Assume that any of
 1137 these metrics is given by a sequence $z = \{z_1, \dots, z_T\}$
 1138 obtained from T intermediate models. Then we can cap-
 1139 ture the uniformity of z via a (weighted) variance score

$$s_{\text{var}} = \sum_t w_t (z_t - \mu)^2$$
 for mean $\mu = \frac{1}{T} \sum_t z_t$ and an in-
 1140 creasing weighting sequence $w = \{w_1, \dots, w_T\}$.

1141 In order to show the effectiveness of the variance score
 s_{var} for continuous metrics, we provide a simple bound on
 1142 the variance of confidence $\max_{y \in \mathcal{Y}} f_t(\mathbf{x})$ in the final check-
 1143 points of the training. Assuming that the model has con-
 1144 verged to a local minima with a low learning rate, we can
 1145 assume that the distribution of model weights can be ap-
 1146 proximated by a Gaussian distribution.

1147 We consider a linear regression problem where the inputs
 1148 are linearly separable.

Lemma 2. *Assume that we have some Gaussian prior
 1149 on the model parameters in the logistic regression setting
 1150 across m final checkpoints. More specifically, given T total
 1151 checkpoints of model parameters $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T\}$ we
 1152 have $p(W = \mathbf{w}_t) = \mathcal{N}(\mathbf{w}_0 \mid \boldsymbol{\mu}, s\mathbf{I})$ for $t \in \{T-m+1, \dots, T\}$ and we assume that final checkpoints of the model
 1153 are sampled from this distribution. We show that the vari-
 1154 ance of model confidence $\max_{y \in \{-1, 1\}} p(y \mid \mathbf{x}_i, \mathbf{w}_t)$ for
 1155 a datapoint (\mathbf{x}_i, y_i) can be upper bounded by a factor of
 1156 probability of correctly classifying this example by the opti-
 1157 mal weights.*

Proof. We first compute the variance of model predictions
 $p(y_i \mid \mathbf{x}_i, W)$ for a given datapoint (\mathbf{x}_i, y_i) . Following pre-
 1158 vious work [4, 30], the variance of predictions over these
 1159 checkpoints can be estimated as follows:

1160 Taking two terms in Taylor expansion for model predic-
 1161 tions we have $p(y_i \mid \mathbf{x}_i, W) \simeq p(y_i \mid \mathbf{x}_i, \mathbf{w}) + g_i(\mathbf{w})^\top (W - \mathbf{w})$ where W and \mathbf{w} are current and the expected estimate
 1162 of the parameters and $g_i(\mathbf{w}) = p(y_i \mid \mathbf{x}_i, \mathbf{w})(1 - p(y_i \mid \mathbf{x}_i, \mathbf{w}))\mathbf{x}_i$ is the gradient vector. Now we can write the vari-
 1163 ance with respect to the model prior as:

$$\mathbb{V}(p(y_i \mid \mathbf{x}_i, W)) \simeq \mathbb{V}(g_i(\mathbf{w})^\top (W - \mathbf{w})) = g_i(\mathbf{w})^\top F^{-1} g_i(\mathbf{w})$$

1164 where F is the variance of posterior distribution $p(W \mid
 1165 X, Y) \sim \mathcal{N}(W \mid \mathbf{w}, F^{-1})$. This suggests that the variance
 1166 of probability of correctly classifying \mathbf{x}_i is proportional to

$$p(y_i \mid \mathbf{x}_i, \mathbf{w})^2(1 - p(y_i \mid \mathbf{x}_i, \mathbf{w}))^2$$
. Now we can bound

1188 the variance of maximum class probability or confidence as
 1189 below:
 1190

$$\begin{aligned} \mathbb{V} \left(\max_{y \in \{-1, 1\}} p(y | \mathbf{x}_i, W) \right) &\leq \mathbb{V}(p(y_i | \mathbf{x}_i, W)) \\ &+ \mathbb{V}(p(-y_i | \mathbf{x}_i, W)) \\ &\approx 2p(y_i | \mathbf{x}_i, \mathbf{w})^2 \\ &\cdot (1 - p(y_i | \mathbf{x}_i, \mathbf{w}))^2 \\ &\cdot \mathbf{x}_i^\top F^{-1} \mathbf{x}_i \end{aligned}$$

□

1200 We showed that if the probability of correctly classifying
 1201 an example given the final estimate of model parameters is
 1202 close to one, the variance of model predictions following a
 1203 Gaussian Prior gets close to zero, we expect a similar
 1204 behaviour for the variance of confidence under samples of this
 1205 distribution.
 1206

C. Extension of Empirical Evaluation

C.1. Comparison with One-Sided Prediction

1207 In addition to the main results, we also compare NNTD
 1208 to one-sided prediction (OSP). For this set of experiments,
 1209 we evaluate our proposed approach on image dataset bench-
 1210 marks that are common in the selective classification liter-
 1211 ature: CIFAR-10, SVHN, and Cats & Dogs. For each
 1212 dataset, we train a deep neural network following the
 1213 ResNet-32 architecture [16] and checkpoint each model af-
 1214 ter processing 50 mini-batches of size 128. All models are
 1215 trained over 200 epochs. SVHN, Cats & Dogs, and GTSRB
 1216 are trained using the Adam optimizer with learning rate
 1217 10^{-3} . The CIFAR-10 and CIFAR-100 models are trained
 1218 using momentum-based stochastic gradient descent with an
 1219 initial learning rate of 10^{-1} on a multi-step learning rate re-
 1220 duction schedule (reduction by 10^{-1} after epochs 100 and
 1221 150), momentum 0.9, and weight decay 10^{-4} . We report
 1222 these results in Tables 4 and 5.
 1223

C.2. Full Hyper-Parameters

1224 We document full hyper-parameter settings for our
 1225 method (NNTD) as well as all baseline approaches in Ta-
 1226 ble 6. Baseline hyper-parameters are consistent with Gan-
 1227 grade *et al.* [10] in the OSP setting and follow the setup
 1228 from Huang *et al.* [19] (Appendix A.4) for our main set of
 1229 experiments. Both DE and MC-DO use an ensemble of 10
 1230 models.
 1231

C.3. Additional Selective Classification Results

C.3.1 Variance Score Experiments on Continuous Metrics

1232 In addition to the evolution of predictions made across in-
 1233 termediate models, we also monitor a variety of easily com-
 1234 putable continuous metrics. In our work, we consider the
 1235 following metrics:
 1236

- Confidence (conf): $\max_{c \in \mathcal{Y}} f_t(\mathbf{x})$
- Confidence gap between top 2 most confident classes (gap): $\max_{c \in \mathcal{Y}} f_t(\mathbf{x}) - \max_{c \neq \hat{y}} f_t(\mathbf{x})$
- Entropy (ent): $-\sum_{c=1}^C f_t(\mathbf{x})_c \log(f_t(\mathbf{x})_c)$

1237 Note that for the purpose of this experiment, we adapt the
 1238 notation of f to map to a softmax-parameterized output in-
 1239 stead of the hard thresholded label, formally $f : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}$.
 1240

1241 We depict example evolution plots for these metrics in
 1242 Figures 7, 8, 9, 10, and 11 in the third row.
 1243

C.3.2 Individual Metric Evolution Plots

1244 We provide additional individual metric evolution plots for
 1245 SVHN (Figure 7), GTSRB (Figure 8), CIFAR-10 (Figure
 1246 9), CIFAR-100 (Figure 10), and Cats & Dogs (Figure
 1247 11).
 1248

C.3.3 Incorporating Estimates of e_t into s_{\min} and s_{avg}

1249 Our results in § 4.3 show that e_t is only nearly 0 almost
 1250 everywhere, we investigate whether incorporating an esti-
 1251 mate of e_t into s_{\min} and s_{avg} leads to additional SC im-
 1252 provements. Recall that in Lemma 1 we gave an upper-bound on
 1253 the probability that the test point is correctly classified as
 1254 $\frac{v_t}{|a_t - e_t|^2}$. In the case where e_t is not 0 everywhere, we can
 1255 adjust
 1256

$$s_{\min} = \min_{t \text{ s.t. } a_t=1} \frac{v_t}{|a_t - e_t|^2} \quad s_{\text{avg}} = \frac{\sum \frac{v_t}{|a_t - e_t|^2} a_t}{\sum a_t} \quad (4)$$

1257 accordingly. In Figure 12, we experimentally test whether
 1258 incorporating an empirical estimate (first row) or a smooth
 1259 decay function similar to v_t (second row) robustly improve
 1260 over s_{\min} and s_{avg} with $e_t = 0$. Our results show that neither
 1261 setting outperforms our default setting with $e_t = 0$.
 1262

C.3.4 Performance of Jump Score s_{jmp} and Weighted Variance s_{var}

1263 As discussed in Appendix B, we also investigated whether
 1264 the jump score s_{jmp} or the weighted variance of continuous
 1265 metrics s_{var} can be used as an effective score for SC. As we
 1266 show in Figure 13, none of these metrics robustly outper-
 1267 forms our main method NNTD($s_{\text{avg}}, 0.05$).
 1268

C.3.5 Concave Weighting for v_t

1269 As we have empirically analyzed in § 4.3, the variances v_t
 1270 follow a monotonically decreasing and convex trend. This
 1271 inspires our best performing method NNTD($s_{\text{avg}}, 0.05$) to
 1272

1296 Table 4. **Performance at low target errors for OSP based setup.** Our NNTD($s_{\text{avg}}, 0.05$) method either matches or provides higher test set
 1297 coverage at a fixed target error than other competing methods. Bold entries are within one standard deviation of each other over 5 random
 1298 runs.

Dataset	Target Error	NNTD		OSP		SR		SN		DG	
		Cov.	Error	Cov.	Error	Cov.	Error	Cov.	Error	Cov.	Error
CIFAR-10	2%	83.3	1.96	80.6	1.91	75.1	2.09	73.0	2.31	74.2	1.98
	1%	79.7	1.05	74.0	1.02	67.2	1.09	64.5	1.02	66.4	1.01
	0.5%	74.2	0.49	64.1	0.51	59.3	0.53	57.6	0.48	57.8	0.51
SVHN	2%	95.7	1.96	95.8	1.99	95.7	2.06	93.5	2.03	94.8	1.99
	1%	91.2	0.99	90.1	1.03	88.4	0.99	86.5	1.04	89.5	1.01
	0.5%	83.9	0.50	82.4	0.51	77.3	0.51	79.2	0.51	81.6	0.49
Cats & Dogs	2%	90.5	2.03	90.5	1.98	88.2	2.03	84.3	1.94	87.4	1.94
	1%	85.6	1.01	85.4	0.98	80.2	0.97	78.0	0.98	81.7	0.98
	0.5%	77.5	0.50	78.7	0.49	73.2	0.49	70.5	0.46	74.5	0.48

1313 Table 5. **Performance at high target coverage for OSP based setup.** Similar as demonstrated in Table 4, NNTD($s_{\text{avg}}, 0.05$) matches or
 1314 outperforms competing error rates at fixed coverage levels.

Dataset	Target Coverage	NNTD		OSP		SR		SN		DG	
		Cov.	Error	Cov.	Error	Cov.	Error	Cov.	Error	Cov.	Error
CIFAR-10	100%	100	9.57	100	9.74	99.99	9.58	100	11.07	100	10.81
	95%	95.1	6.13	95.1	6.98	95.2	8.74	94.7	8.34	95.1	8.21
	90%	90.1	4.16	90.0	4.67	90.5	6.52	89.6	6.45	90.1	6.14
SVHN	100%	100	4.11	100	4.27	99.97	3.86	100	4.27	100	4.03
	95%	94.8	1.80	95.1	1.83	95.1	1.86	95.1	2.53	95.0	2.05
	90%	90.1	0.77	90.1	1.01	90.0	1.04	90.1	1.31	90.0	1.06
Cats & Dogs	100%	100	5.18	100	5.93	100	5.72	100	7.36	100	6.16
	95%	94.9	2.99	95.1	2.97	95.0	3.46	95.2	5.1	95.1	4.28
	90%	90.0	1.87	90.0	1.74	90.0	2.28	90.2	3.3	90.0	2.50

1330 use $k = 0.05$ for $v_i = 1 - i^k$. As a sanity check, we examine whether a concave weighting yielded by $v_i = 1 - i^k$ for
 1331 $k \geq 1$. As we demonstrate in Figure 14, the best concave
 1332 weighting is given by $k = 1$. Therefore, we confirm that
 1333 no concave weighting outperforms the convex weightings
 1334 analyzed in Figure 4.

C.3.6 Limiting NNTD to a Subset of Last Checkpoints

1335 In order to determine which training stages are important
 1336 for selective classification, we perform an experiment
 1337 on CIFAR-10 and SVHN where we limit ourselves to a
 1338 subset of the last checkpoints. In particular, we examine
 1339 the coverage/error trade-off for only including the last
 1340 $\{10\%, 20\%, 50\%, 80\%, 90\%, 100\%\}$ of checkpoints. We
 1341 report our findings in Table 7. We see that taking into
 1342 account more than 80% of checkpoints does lead to diminishing
 1343 returns and that only taking into account 50% or less of
 1344 the total numbers of checkpoints does not lead to state-of-
 1345 the-art selective classification performance.

C.3.7 Applying OOD Scores to Selective Classification

Finally, we also compare popular out-of-distribution detection approaches to our NNTD method. In particular, we apply three state-of-the-art approaches, namely energy-based OOD detection (Energy) [26], Mahalanobis-distance-based OOD detection (Mahalanobis) [24], and ODIN (ODIN) [25] to the CIFAR-10 and SVHN test sets and document these results in Table 8. Overall, we find that SOTA out-of-distribution detection techniques are not well equipped to attain SOTA selective classification performance; NNTD outperforms these methods by a large margin.

Table 6. Hyper-parameters used for all SC algorithms in the OSP setup.

Dataset	SC Algorithm	Hyper-Parameters
CIFAR-10	Softmax Response (SR)	$t = 0.0445$
	Selective Net (SN)	$\lambda = 32, c = 0.51, t = 0.24$
	Deep Gamblers (DG)	$o = 1.179, t = 0.03$
	One-Sided Prediction (OSP)	$\mu = 0.49, t = 0.8884$
	Neural Network Training Dynamics (NNTD)	$T = 1593, k = 0.05$
SVHN	Softmax Response (SR)	$t = 0.0224$
	Selective Net (SN)	$\lambda = 32, c = 0.79, t = 0.86$
	Deep Gamblers (DG)	$o = 1.13, t = 0.23$
	One-Sided Prediction (OSP)	$\mu = 1.67, t = 0.9762$
	Neural Network Training Dynamics (NNTD)	$T = 1195, k = 0.05$
Cats & Dogs	Softmax Response (SR)	$t = 0.029$
	Selective Net (SN)	$\lambda = 32, c = 0.7, t = 0.73$
	Deep Gamblers (DG)	$o = 1.34, t = 0.06$
	One-Sided Prediction (OSP)	$\mu = 1.67, t = 0.9532$
	Neural Network Training Dynamics (NNTD)	$T = 797, k = 0.05$

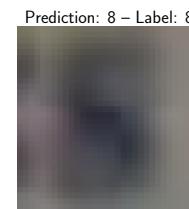
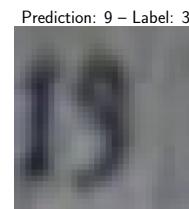
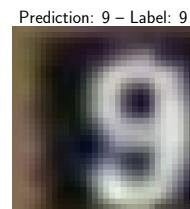
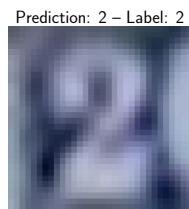
Table 7. Coverage/Error trade-off of NNTD as a function of the accounted data points (CIFAR-10 left, SVHN right). We observe that large parts of the training process contain valuable signals for selective classification.

Target Coverage	100%	90%	80%	50%	20%	10%	Target Coverage	100%	90%	80%	50%	20%	10%
100%	6.07	6.08	6.10	6.15	6.18	6.20	100%	2.68	2.70	2.73	2.83	2.88	2.92
95%	3.24	3.25	3.25	3.31	3.34	3.39	95%	0.88	0.90	0.93	1.01	1.09	1.11
90%	1.83	1.83	1.84	1.88	1.91	1.93	90%	0.55	0.57	0.59	0.62	0.68	0.73
80%	0.64	0.65	0.67	0.72	0.78	0.79	80%	0.38	0.39	0.40	0.45	0.51	0.53
70%	0.34	0.34	0.36	0.38	0.40	0.40	70%	0.33	0.33	0.35	0.41	0.44	0.45

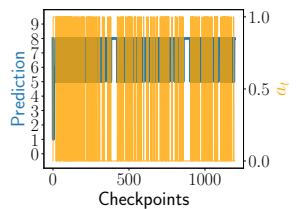
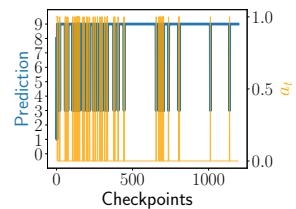
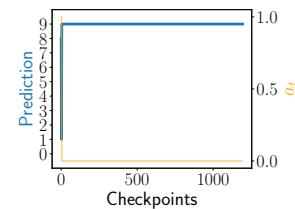
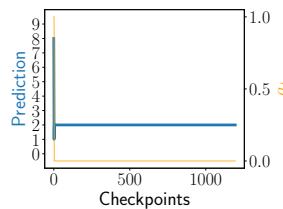
Table 8. Performance of NNTD in comparison with OOD scores. We find that NNTD significantly outperforms common OOD detection approaches for the purpose of selective classification.

Dataset	Target Coverage	NNTD		Energy		Mahalanobis		ODIN	
		Cov	Err \downarrow	Cov	Err \downarrow	Cov	Err \downarrow	Cov	Err \downarrow
CIFAR-10	100%	100	6.07	100	6.07	100	6.07	100	6.07
	95%	95.0	3.24	95.0	3.99	95.0	4.44	95.1	6.04
	90%	90.1	1.83	90.0	2.67	90.0	2.93	88.9	6.01
	80%	79.9	0.64	80.0	1.10	80.1	1.20	88.9	6.01
	70%	69.8	0.34	70.0	0.83	69.9	0.82	68.3	4.41
SVHN	100%	100	2.68	100	2.68	100	2.68	100	2.68
	95%	95.0	0.88	95.1	1.35	95.1	1.51	95.0	2.67
	90%	90.1	0.55	89.9	0.92	90.0	1.04	89.9	2.64
	80%	79.9	0.38	79.9	0.70	80.3	0.70	81.9	2.51
	70%	69.8	0.33	70.0	0.58	69.8	0.58	74.9	2.18

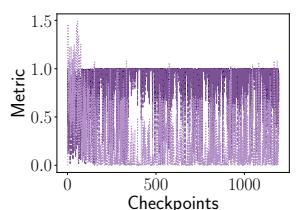
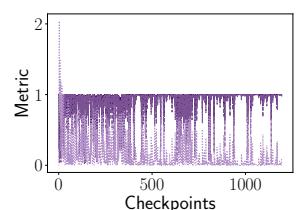
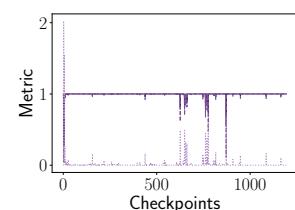
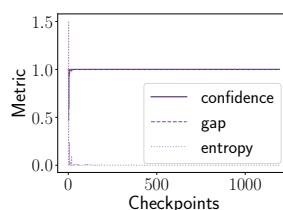
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521



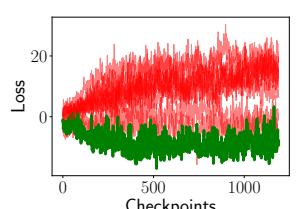
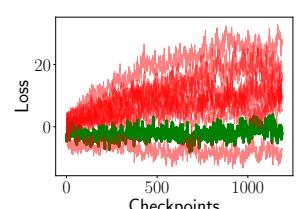
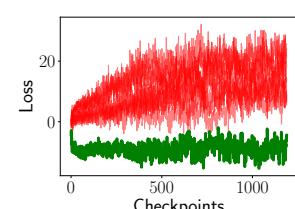
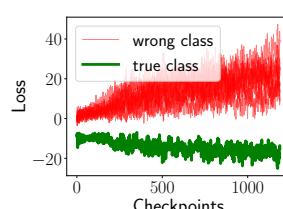
1522
1523
1524
1525
1526
1527
1528



1529



1536



1544

1545
1546
1547
1548
1549
1550
1551
1552

Figure 7. Individual SVHN examples. We extend Figure 2 by also including the evolution of the confidence, the gap, and the entropy over the course of training in the third row. In the final row, we further plot the loss evolution for all candidate classes. We note that the first two examples are easy to classify as indicated by mostly stationary metrics and a clear loss separation, while the last two examples are harder to classify since they exhibit highly erratic metrics and ambiguous loss curves.

1566

1567

1568

1569

1570

1571

1572

1573

1574

1575

1576

1577

1578

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

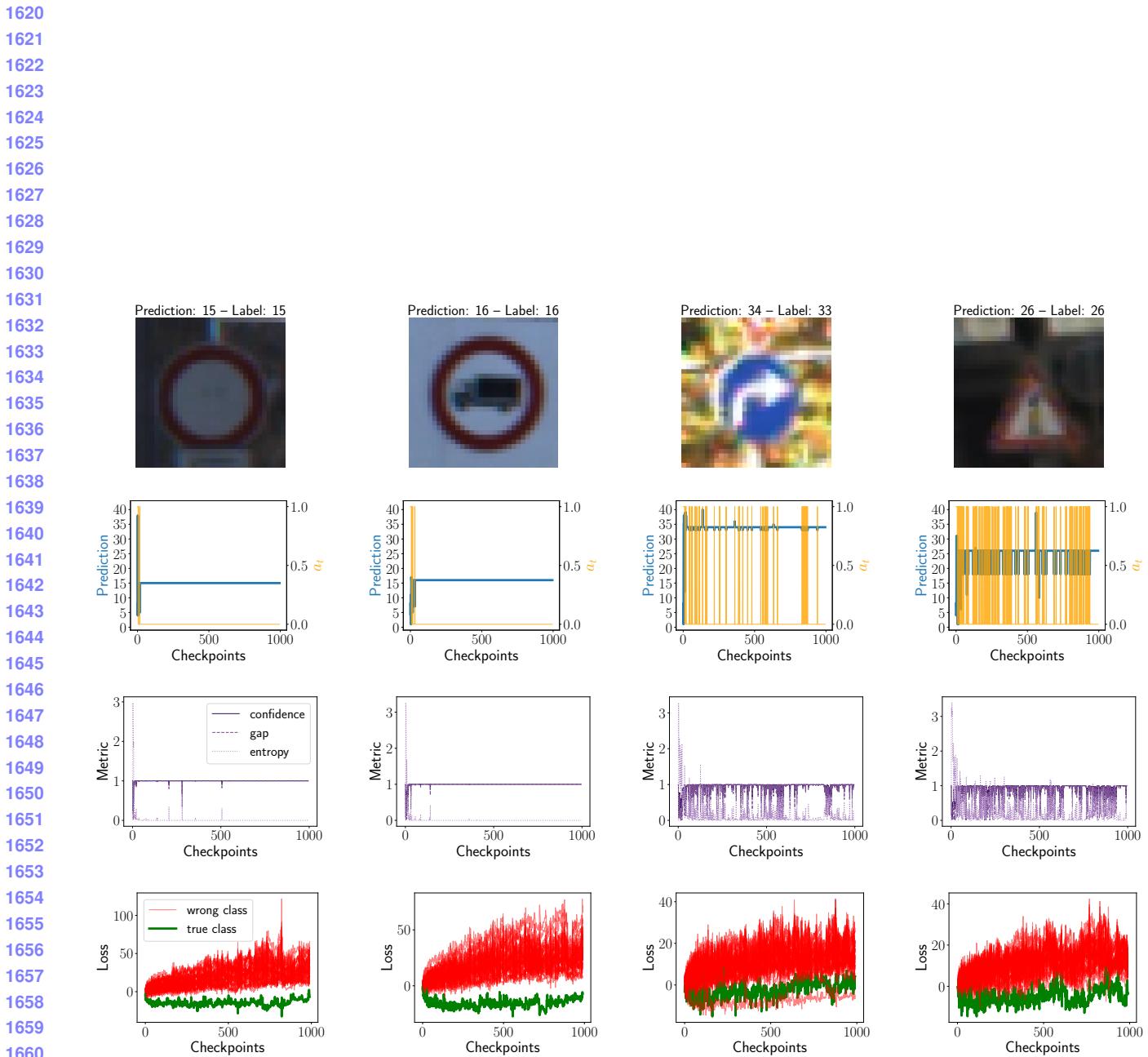
1615

1616

1617

1618

1619

Figure 8. **Individual GTSRB examples.** Similar as Figure 7.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738

1739 Prediction: 9 – Label: 9



1740 Prediction: 6 – Label: 6



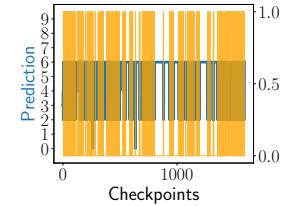
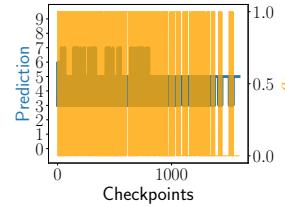
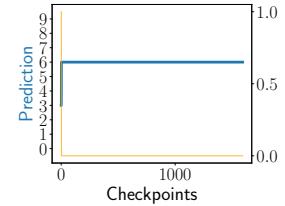
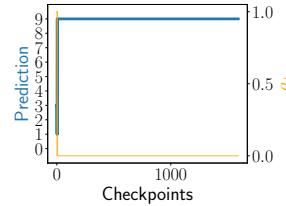
1741 Prediction: 5 – Label: 3



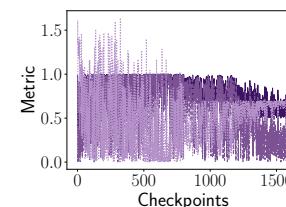
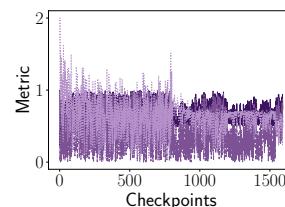
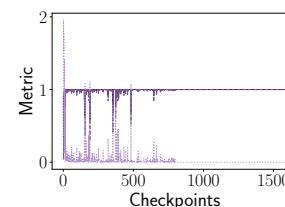
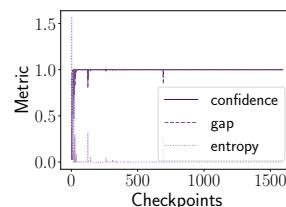
1742 Prediction: 6 – Label: 6



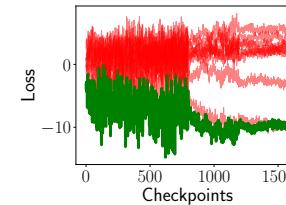
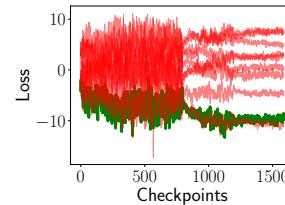
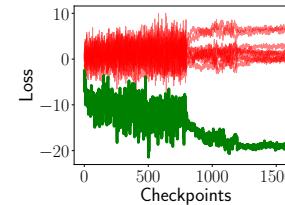
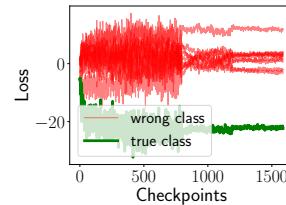
1743
1744
1745
1746



1747
1748
1749
1750
1751
1752
1753



1754
1755
1756
1757
1758
1759
1760
1761



1762
1763
1764
1765
1766
1767
1768
1769

Figure 9. Individual CIFAR-10 examples. Similar as Figure 7.

1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

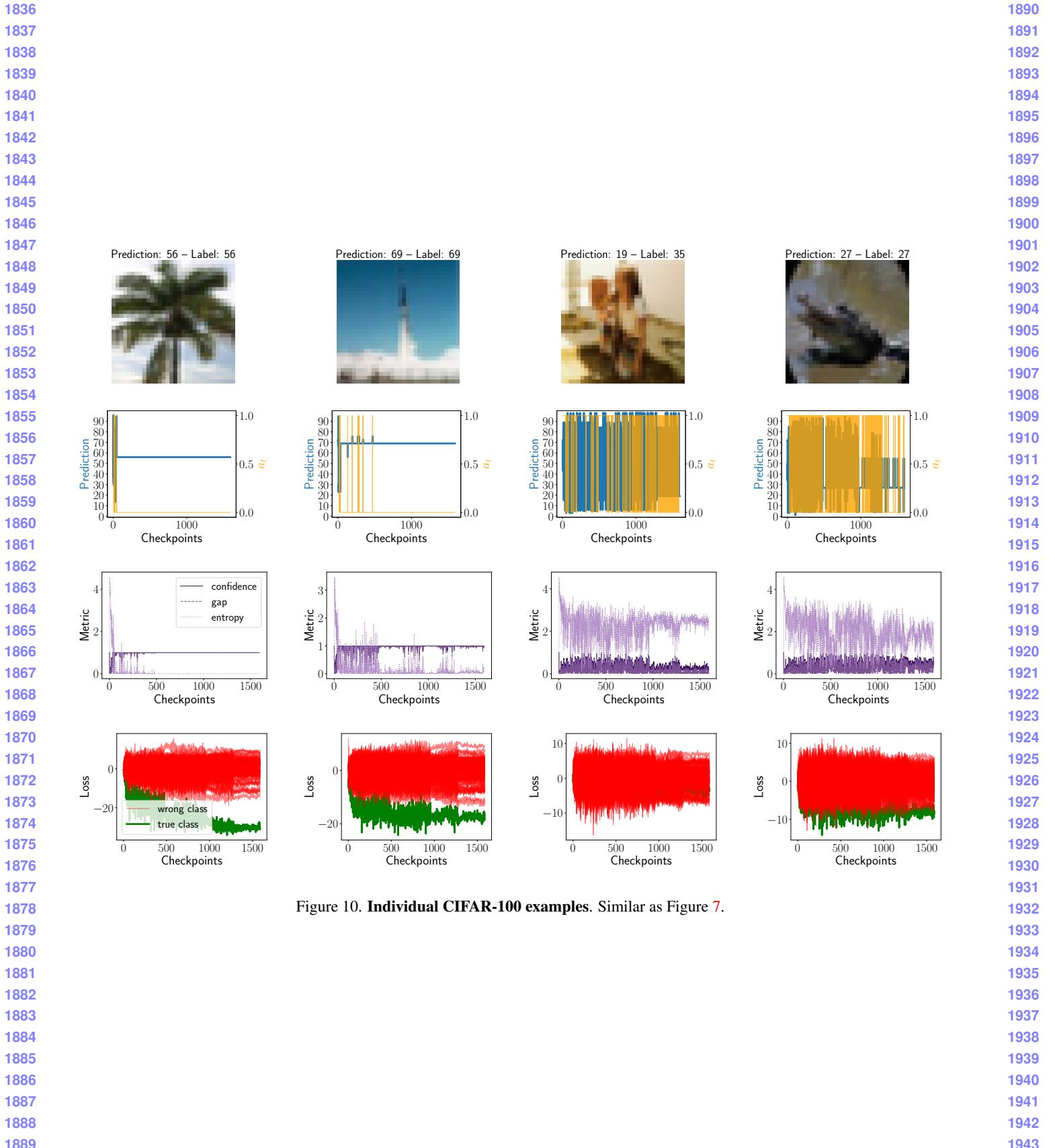


Figure 10. Individual CIFAR-100 examples. Similar as Figure 7.

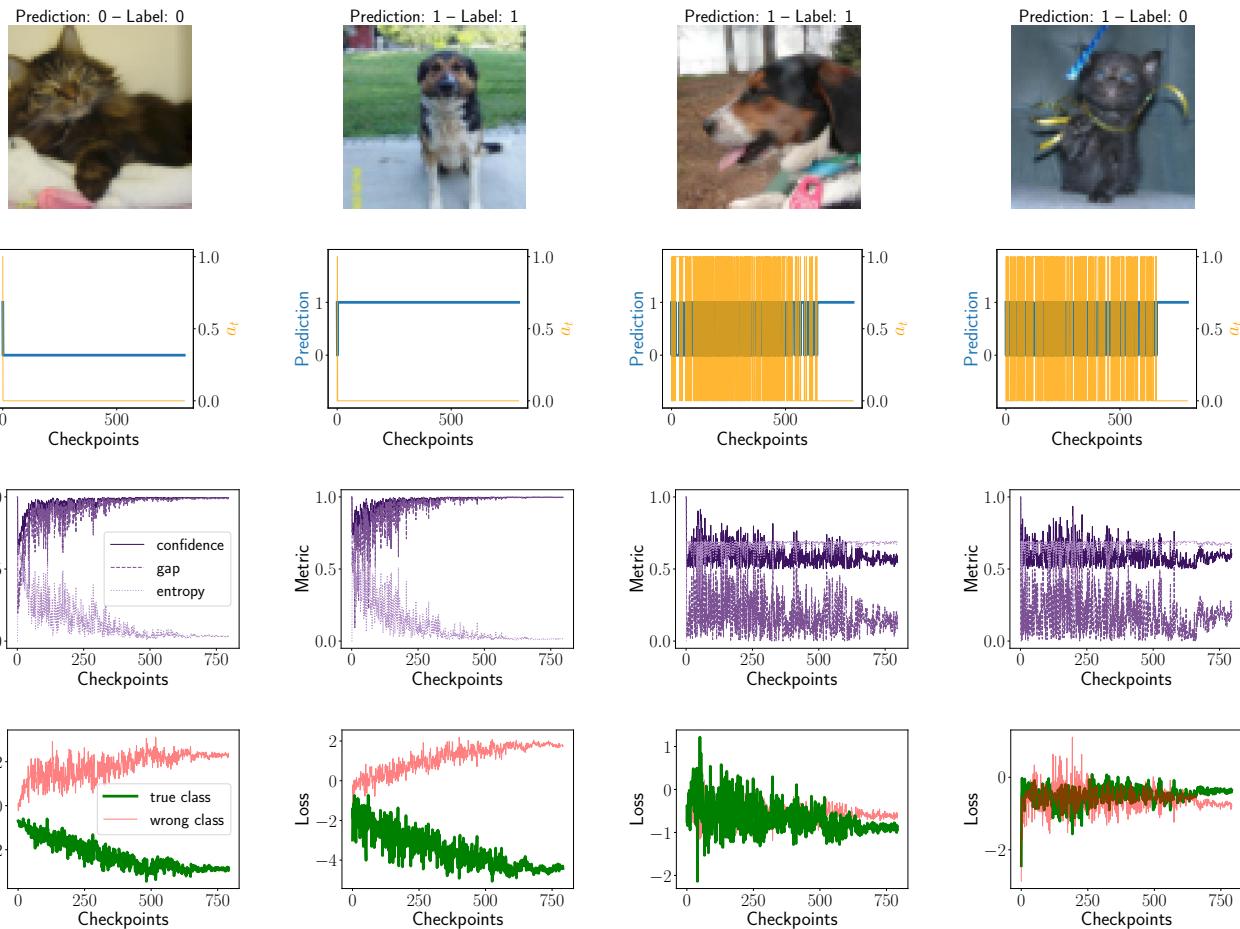


Figure 11. Individual Cats & Dogs examples. Similar as Figure 7.

2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

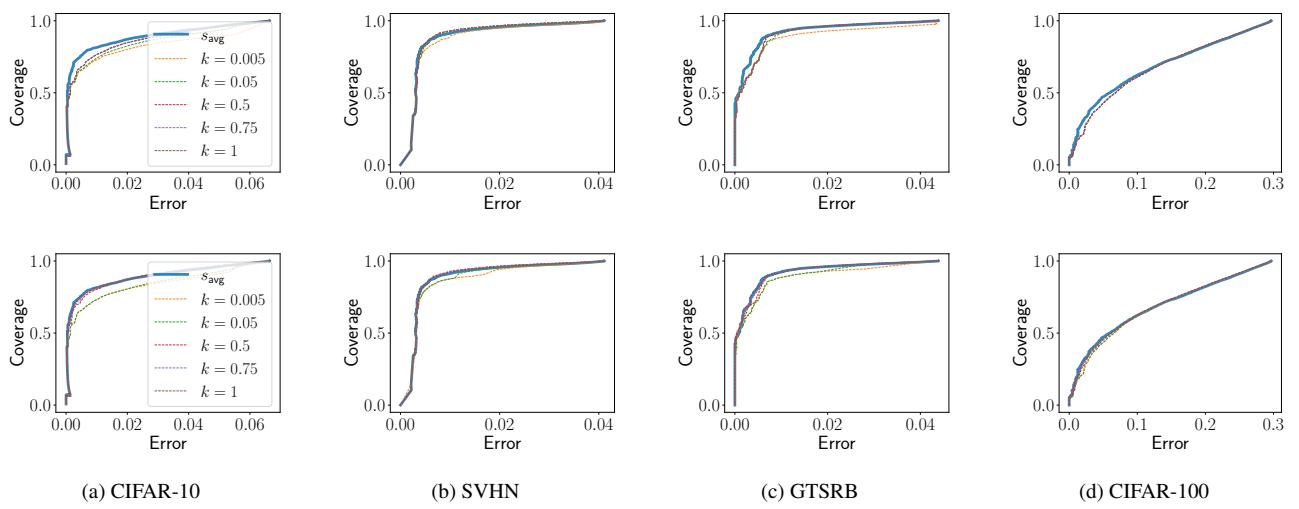


Figure 12. **Coverage/error trade-off when incorporating e_t into s_{\min} and s_{avg} .** In the first row, the solid blue line corresponds to $\text{NNTD}(s_{\text{avg}}, 0.05)$ while the dashed lines show the performance when we incorporate an empirical estimate of e_t for various weightings v_t . In the second row, maintaining the solid blue line as $\text{NNTD}(s_{\text{avg}}, 0.05)$, we fix v_t at $k = 0.05$ and test various continuous approximations of e_t as $e_t = 1 - t^k$ for $k \in (0, 1]$. Overall, we find that introducing an adaptable e_t does not further improve the performance of s_{\min} and s_{avg} with $e_t = 0$.

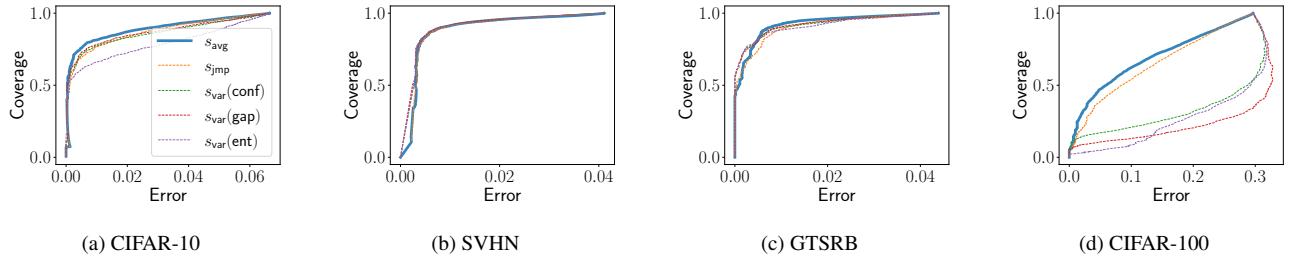


Figure 13. **Coverage/error trade-off of NNTD for alternate scores.** We find that neither the jump score, nor any of the weighted variance metrics at their optimal k for v_i outperform $\text{NNTD}(s_{\text{avg}}, 0.05)$.

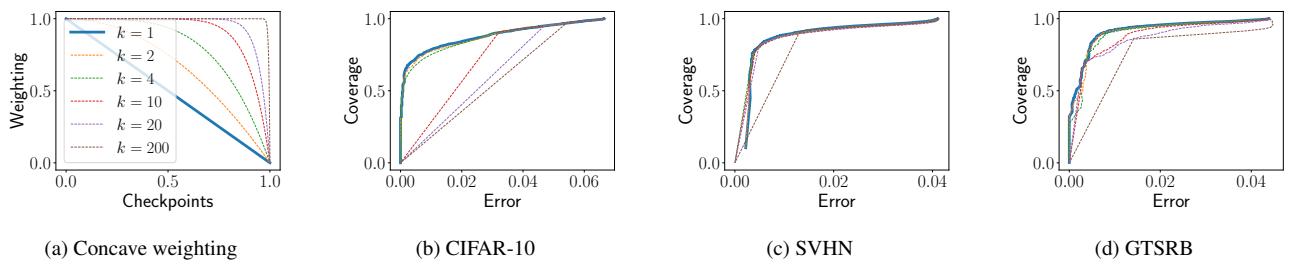


Figure 14. **Concave weighting used in $v_t = 1 - t^k$.** Overall, since the best concave weighting is given by $k = 1$, we conclude that no concave weighting outperforms any convex weighting.

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

Table 9. Performance at low target errors. Same results as in Table 1 but includes standard deviations.

Dataset	Target	NNTD		SAT		DG		SN		SR		MC-DO	
		Error	Cov ↑	Err	Cov ↑	Err	Cov ↑	Err	Cov ↑	Err	Cov ↑	Err	Cov ↑
CIFAR-10	2%	91.2 (± 0.05)	1.99 (± 0.01)	90.3 (± 0.03)	1.97 (± 0.02)	89.1 (± 0.04)	2.02 (± 0.01)	88.3 (± 0.03)	2.03 (± 0.02)	85.8 (± 0.08)	1.98 (± 0.02)	86.1 (± 0.07)	2.01 (± 0.01)
	1%	86.4 (± 0.03)	1.00 (± 0.01)	86.1 (± 0.02)	1.02 (± 0.02)	85.5 (± 0.07)	1.03 (± 0.02)	84.4 (± 0.06)	0.98 (± 0.02)	79.1 (± 0.08)	1.01 (± 0.03)	79.9 (± 0.08)	1.01 (± 0.03)
	0.5%	75.9 (± 0.04)	0.49 (± 0.02)	76.0 (± 0.05)	0.51 (± 0.02)	75.2 (± 0.08)	0.5 (± 0.03)	74.7 (± 0.09)	0.49 (± 0.01)	71.2 (± 0.05)	0.51 (± 0.02)	72.0 (± 0.03)	0.50 (± 0.03)
SVHN	2%	98.5 (± 0.05)	1.98 (± 0.03)	98.2 (± 0.04)	1.99 (± 0.02)	97.8 (± 0.03)	2.06 (± 0.05)	97.7 (± 0.06)	2.03 (± 0.02)	97.6 (± 0.05)	1.99 (± 0.03)	97.9 (± 0.03)	2.00 (± 0.04)
	1%	96.3 (± 0.03)	0.99 (± 0.02)	95.7 (± 0.02)	1.03 (± 0.03)	94.8 (± 0.04)	0.99 (± 0.01)	94.5 (± 0.03)	1.04 (± 0.02)	93.5 (± 0.05)	1.01 (± 0.03)	94.1 (± 0.06)	0.97 (± 0.02)
	0.5%	88.1 (± 0.02)	0.50 (± 0.02)	87.9 (± 0.05)	0.51 (± 0.01)	86.4 (± 0.04)	0.51 (± 0.01)	86.0 (± 0.04)	0.51 (± 0.01)	70.0 (± 0.09)	0.50 (± 0.04)	70.1 (± 0.04)	0.49 (± 0.03)
Cats & Dogs	2%	97.7 (± 0.09)	2.01 (± 0.03)	98.2 (± 0.02)	1.98 (± 0.03)	98.0 (± 0.05)	2.03 (± 0.02)	97.4 (± 0.05)	1.98 (± 0.03)	95.1 (± 0.12)	1.99 (± 0.04)	95.7 (± 0.10)	1.99 (± 0.03)
	1%	93.1 (± 0.03)	1.01 (± 0.02)	93.6 (± 0.02)	0.98 (± 0.03)	92.6 (± 0.08)	0.97 (± 0.04)	92.2 (± 0.05)	0.98 (± 0.02)	86.9 (± 0.13)	0.98 (± 0.04)	88.6 (± 0.09)	1.01 (± 0.02)
	0.5%	85.7 (± 0.06)	0.51 (± 0.02)	86.0 (± 0.04)	0.49 (± 0.01)	85.3 (± 0.02)	0.49 (± 0.02)	84.8 (± 0.03)	0.46 (± 0.05)	68.4 (± 0.10)	0.48 (± 0.02)	70.1 (± 0.08)	0.51 (± 0.03)

Table 10. Performance at high target coverage. Same results as in Table 2 but includes standard deviations.

Dataset	Target	NNTD		SAT		DG		SN		SR		MC-DO	
		Coverage	Cov	Err ↓	Cov								
CIFAR-10	100%	100 (± 0.00)	6.07 (± 0.05)	100 (± 0.00)	6.06 (± 0.03)	100 (± 0.00)	6.11 (± 0.05)	100 (± 0.00)	6.13 (± 0.03)	100 (± 0.00)	6.07 (± 0.05)	100 (± 0.00)	6.07 (± 0.05)
	95%	95.0 (± 0.01)	3.24 (± 0.03)	95.1 (± 0.01)	3.32 (± 0.05)	95.1 (± 0.02)	3.47 (± 0.06)	95.0 (± 0.01)	4.08 (± 0.08)	94.9 (± 0.03)	4.48 (± 0.07)	95.1 (± 0.04)	4.48 (± 0.09)
	90%	90.1 (± 0.02)	1.83 (± 0.04)	89.9 (± 0.02)	1.90 (± 0.02)	90.0 (± 0.01)	2.19 (± 0.05)	90.1 (± 0.02)	2.29 (± 0.03)	90.1 (± 0.02)	2.78 (± 0.06)	90.0 (± 0.02)	2.87 (± 0.07)
	80%	79.9 (± 0.02)	0.64 (± 0.03)	80.0 (± 0.00)	0.65 (± 0.04)	80.1 (± 0.03)	0.66 (± 0.04)	80.1 (± 0.01)	0.81 (± 0.07)	79.8 (± 0.02)	1.05 (± 0.08)	79.9 (± 0.01)	1.01 (± 0.03)
	70%	69.8 (± 0.03)	0.34 (± 0.04)	69.9 (± 0.03)	0.32 (± 0.05)	69.8 (± 0.04)	0.41 (± 0.04)	70.2 (± 0.03)	0.30 (± 0.02)	70.0 (± 0.01)	0.47 (± 0.07)	70.1 (± 0.02)	0.42 (± 0.05)
SVHN	100%	100 (± 0.00)	2.68 (± 0.02)	100 (± 0.00)	2.71 (± 0.03)	100 (± 0.00)	2.72 (± 0.05)	100 (± 0.00)	2.77 (± 0.06)	100 (± 0.00)	2.68 (± 0.02)	100 (± 0.00)	2.68 (± 0.02)
	95%	95.0 (± 0.01)	0.88 (± 0.02)	95.1 (± 0.02)	0.95 (± 0.03)	95.1 (± 0.01)	1.01 (± 0.06)	95.0 (± 0.02)	1.07 (± 0.05)	94.9 (± 0.03)	1.15 (± 0.11)	95.1 (± 0.02)	1.12 (± 0.07)
	90%	90.1 (± 0.02)	0.55 (± 0.4)	89.9 (± 0.01)	0.58 (± 0.01)	90.0 (± 0.00)	0.63 (± 0.06)	90.1 (± 0.01)	0.71 (± 0.04)	90.1 (± 0.02)	0.82 (± 0.06)	90.0 (± 0.02)	0.76 (± 0.03)
	80%	79.9 (± 0.01)	0.38 (± 0.02)	80.0 (± 0.01)	0.37 (± 0.02)	80.1 (± 0.01)	0.43 (± 0.01)	80.1 (± 0.00)	0.48 (± 0.02)	79.8 (± 0.03)	0.55 (± 0.09)	79.9 (± 0.02)	0.53 (± 0.08)
	70%	69.8 (± 0.03)	0.33 (± 0.03)	69.9 (± 0.02)	0.33 (± 0.01)	69.8 (± 0.04)	0.35 (± 0.02)	70.2 (± 0.02)	0.45 (± 0.06)	70.0 (± 0.01)	0.50 (± 0.05)	70.1 (± 0.02)	0.49 (± 0.06)
Cats & Dogs	100%	100 (± 0.00)	3.48 (± 0.01)	100 (± 0.00)	3.45 (± 0.01)	100 (± 0.00)	3.41 (± 0.03)	100 (± 0.00)	3.56 (± 0.04)	100 (± 0.00)	3.48 (± 0.01)	100 (± 0.00)	3.48 (± 0.01)
	95%	95.1 (± 0.02)	1.51 (± 0.03)	95.1 (± 0.03)	1.45 (± 0.02)	95.0 (± 0.01)	1.43 (± 0.02)	94.9 (± 0.02)	1.61 (± 0.05)	94.8 (± 0.03)	1.92 (± 0.12)	95.1 (± 0.02)	1.95 (± 0.08)
	90%	90.1 (± 0.03)	0.60 (± 0.03)	89.9 (± 0.02)	0.57 (± 0.03)	90.0 (± 0.01)	0.69 (± 0.04)	90.1 (± 0.02)	0.95 (± 0.11)	90.1 (± 0.03)	1.13 (± 0.07)	90.0 (± 0.01)	1.09 (± 0.06)
	80%	79.9 (± 0.01)	0.42 (± 0.04)	80.0 (± 0.01)	0.41 (± 0.03)	80.1 (± 0.02)	0.56 (± 0.02)	80.1 (± 0.03)	0.39 (± 0.05)	79.8 (± 0.02)	0.69 (± 0.07)	79.9 (± 0.02)	0.58 (± 0.05)
	70%	69.8 (± 0.02)	0.36 (± 0.05)	69.9 (± 0.03)	0.33 (± 0.03)	69.8 (± 0.03)	0.45 (± 0.05)	70.2 (± 0.03)	0.33 (± 0.03)	70.0 (± 0.01)	0.62 (± 0.06)	70.1 (± 0.02)	0.51 (± 0.04)