

## A Implementation details

### A.1 Proposed algorithm

The key idea of this algorithm is to compute the proposed value function Equation 3 for each given class  $c$  (line 4). We implemented the approximation to the canonical CS-SHAPLEY suggested in Equation 6 and in practice we found that  $K = 500$  is sufficient for our experiments. Lines 10 – 18 is adopted from the truncated Monte Carlo algorithm [3] for computing the Shapley values for the training instances in class  $c$ , with  $\varepsilon = 10^{-4}$ . Note that, the approximation to the canonical CS-SHAPLEY also requires drawing different subsets  $S_{-c}^{(k)} \in T_{-c}$  for each  $k$  (line 7), which, based on our observation, is critical for estimating reliable CS-SHAPLEY values.

---

#### Algorithm 1: CS-SHAPLEY

---

```

1: Input: Train data  $T = \{1, \dots, n\}$ , label set  $C$ , development set  $D$ 
2: Output: Shapley value of training points:  $\phi_1, \dots, \phi_n$ 
3: Initialize  $\phi_i \leftarrow 0$  for  $i \in \{1, \dots, n\}$ 
4: for  $c \in C$  do
5:   for  $k \in \{1, \dots, K\}$  do
6:      $\pi_c^{(k)}$ : A random permutation of the training instances with label  $c$ ,  $T_c$ 
7:      $S_{-c}^{(k)} \subseteq T_{-c}$ : A random subset of training instances that in other classes
8:     Set  $v_c(\emptyset | S_{-c}^{(k)}) \leftarrow 0$ 
9:     Compute  $v_c(T_c | S_{-c}^{(k)})$  using Equation 3
10:    for  $j \in \{1, \dots, |T_c|\}$  do
11:      //  $\pi_c^{(k)}[1 : j - 1]$  represents the first  $j - 1$  examples in  $\pi_c^{(k)}$ 
12:      if  $|v_c(T_c | S_{-c}^{(k)}) - v_c(\pi_c^{(k)}[1 : j - 1] | S_{-c}^{(k)})| < \varepsilon$  then
13:         $v_c(\pi_c^{(k)}[1 : j] | S_{-c}^{(k)}) \leftarrow v_c(\pi_c^{(k)}[1 : j - 1] | S_{-c}^{(k)})$ 
14:      else
15:        Compute  $v_c(\pi_c^{(k)}[1 : j] | S_{-c}^{(k)})$  using Equation 3
16:      end if
17:       $\phi_{\pi_c^{(k)}[j]} \leftarrow \frac{k-1}{k} \phi_{\pi_c^{(k-1)}[j]} + \frac{1}{k} \{v_c(\pi_c^{(k)}[1 : j] | S_{-c}^{(k)}) - v_c(\pi_c^{(k)}[1 : j - 1] | S_{-c}^{(k)})\}$ 
18:    end for
19:  end for
20:  // Normalize the values of the instances in class  $c$  to satisfy the efficiency axiom
21:   $\sigma_c \leftarrow \sum_{j \in T_c} \phi_j$ 
22:   $\phi_j \leftarrow \frac{\phi_j}{\sigma_c} \cdot a_T(D_c)$  for  $j \in T_c$ 
23: end for

```

---

### A.2 Data usage

We use nine real datasets, including five tabular datasets and four image datasets: Diabetes, Click, CPU, Covertypes, Phoneme, FMNIST, CIFAR10, and MNIST (binarized and multi-class). All datasets are licensed under CC-BY. Dataset statistics and sources are provided in Table 3. Datasets do not contain offensive or personally identifiable information. For any consent information, if available, please see the cited data sources.

For each dataset, we maintain the original label distribution by performing stratified sampling. For all datasets other than Diabetes, we sample a 500/500/2000 train/dev/test split. Since the Diabetes dataset only contains 768 total instances, we use a 128/128/512 train/dev/test split. Split ratios are based on common ratios from prior works [14, 3].

For MNIST [16], CIFAR-10 [13], and Fashion MNIST [25] we follow the common procedure of prior work [14, 15]: we extract the learned image representations from ResNet-18 [8] using the pre-trained weights available from Pytorch [19]. Then, we fit a principal component analysis model on the extracted feature vectors and select the first 32 principal components. As prior work has

Table 3: Statistics for each dataset. Input dimensions refers to dimensionality of input ( $x$ ) without class label ( $y$ ).

Dataset	Data Type	Classes	Input Dims	Source
Diabetes	Tabular	2	8	<a href="https://www.openml.org/d/37">https://www.openml.org/d/37</a>
Click	Tabular	2	11	<a href="https://www.openml.org/d/1216">https://www.openml.org/d/1216</a>
CPU	Tabular	2	21	<a href="https://www.openml.org/d/197">https://www.openml.org/d/197</a>
Coverttype	Tabular	7	54	<a href="https://www.openml.org/d/1596">https://www.openml.org/d/1596</a>
Phoneme	Tabular	2	5	<a href="https://www.openml.org/d/1489">https://www.openml.org/d/1489</a>
FMNIST	Image	2	32	<a href="https://www.openml.org/d/40996">https://www.openml.org/d/40996</a>
CIFAR10	Image	2	32	<a href="https://www.openml.org/d/40927">https://www.openml.org/d/40927</a>
MNIST (binary)	Image	2	32	<a href="https://www.openml.org/d/554">https://www.openml.org/d/554</a>
MNIST (multi)	Image	10	32	<a href="https://www.openml.org/d/554">https://www.openml.org/d/554</a>

used binarized image datasets [14, 3], we binarize the datasets as follows:  $y = \{\text{automobile, truck}\}$  (CIFAR10),  $y = \{\text{t-shirt and tops, shirts}\}$  (FMNIST),  $y = \{1, 7\}$  (MNIST).

### A.3 Experiment settings

We use the model implementations from the Python module `scikit-learn`<sup>5</sup>, licensed under BSD 3-Clause License. Experiments primarily use `scikit-learn` default hyperparameters with some preliminary hyperparameter tuning. For example, as we use small datasets, we perform preliminary hyperparameter tuning of the logistic regression solver between the default 'lbfgs' and 'liblinear', selecting 'liblinear' as it had higher validation set performance across datasets. For the gradient boosting classifier, we reduced the size and number of estimators to maintain computational feasibility for all methods (see section 5 for example of computational scaling), using 40 estimators with a maximum depth of 2 and 6 required samples to split internal nodes. Experiments were run on 65 Intel(R) Xeon(R) Silver 4110 CPUs 2.10GHz.

For all reported results we use a development set to compute the train data values, then evaluate model performance on a held-out test set. Further, results represent the mean of 5 trials. In plots, the shaded region indicates standard deviation.

<sup>5</sup><https://scikit-learn.org/stable/>

## B Theoretical justification

### B.1 Additional discussion

**Development set level metrics:** In addition to the issues with development set level accuracy, using another development set level metric (e.g.  $F1$ ) does not resolve the issue of being able to distinguish between in-class and out-of-class contribution. In other words, the same issues are present with other performance metrics as with accuracy in the original utility function. We offer an example in Table 4 and Table 5:

Table 4: Confusion matrix prior to adding a noisy instance from the positive class.

	Predicted Label	
	+	-
Ground Truth	+	3   2
	-	2   3

Table 5: Confusion matrix after adding a noisy instance from the positive class.

	Predicted Label	
	+	-
Ground Truth	+	2   3
	-	1   4

In Table 4, prior to adding a noisy data point, we have precision= 0.6, recall= 0.6, so  $F1 = 0.6$ . In Table 5, we have precision= 0.7, recall= 0.6, so  $F1 = 0.65$ . The contribution of this instance is positive, even though it decreased the prediction performance of its own class, which is the same conceptual issue encountered with development set level accuracy that is exemplified in Figure 1.

## C Additional results

In the following subsections, we report the results across experiments for the KNN and gradient boosting classifiers. We also include the additional plots on removal and transferability for logistic regression and SVM with the RBF kernel.<sup>6</sup>

### C.1 Additional removal results

We report the weighted accuracy drop (WAD) for gradient boosting classifier and KNN and plot the prediction accuracy as a function of training instances removed for SVM-RBF, gradient boosting classifier, and KNN.

Table 6: Weighted accuracy drop for gradient boosting classifier and KNN using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Gradient Boosting				KNN			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.097</b>	0.027	0.053	-0.027	<b>0.108</b>	0.084	0.090	0.024
Click	<b>0.029</b>	0.017	0.026	0.012	<b>0.050</b>	0.008	0.007	0.002
Coverttype	<b>0.142</b>	0.101	0.029	-0.014	<b>0.232</b>	0.195	0.150	0.002
CPU	0.021	0.008	0.008	<b>0.046</b>	<b>0.053</b>	0.028	0.027	0.044
Diabetes	<b>0.117</b>	0.095	0.055	0.091	<b>0.226</b>	0.080	0.050	0.069
FMNIST	<b>0.117</b>	0.079	0.089	0.047	<b>0.079</b>	0.034	0.024	0.058
MNIST-2	<b>0.018</b>	0.006	0.012	0.012	0.010	<b>0.018</b>	0.012	0.013
MNIST-10	—	—	—	—	0.184	<b>0.241</b>	0.100	0.105
Phoneme	<b>0.073</b>	0.028	0.039	0.020	<b>0.064</b>	0.042	0.037	0.048

<sup>6</sup>Gradient boosting classifier on MNIST-10 is computationally prohibitive and therefore not reported.



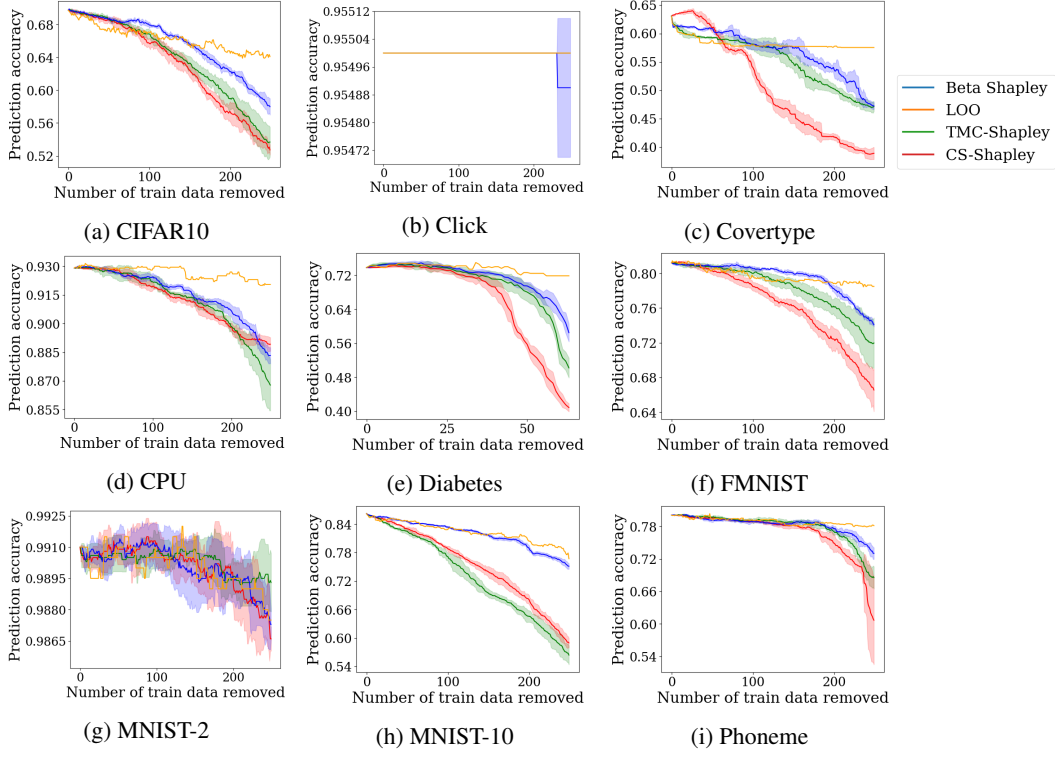


Figure 5: Performance across datasets when removing high-value instances for SVM with RBF Kernel.

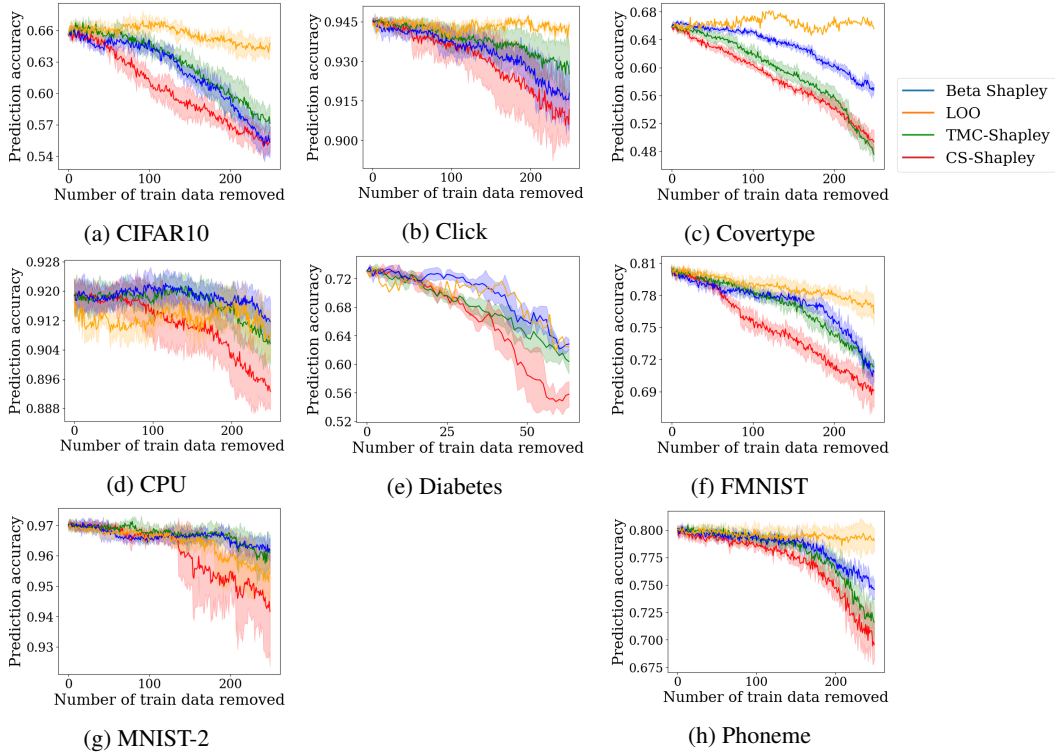


Figure 6: Performance across datasets when removing high-value instances for gradient boosting classifier.

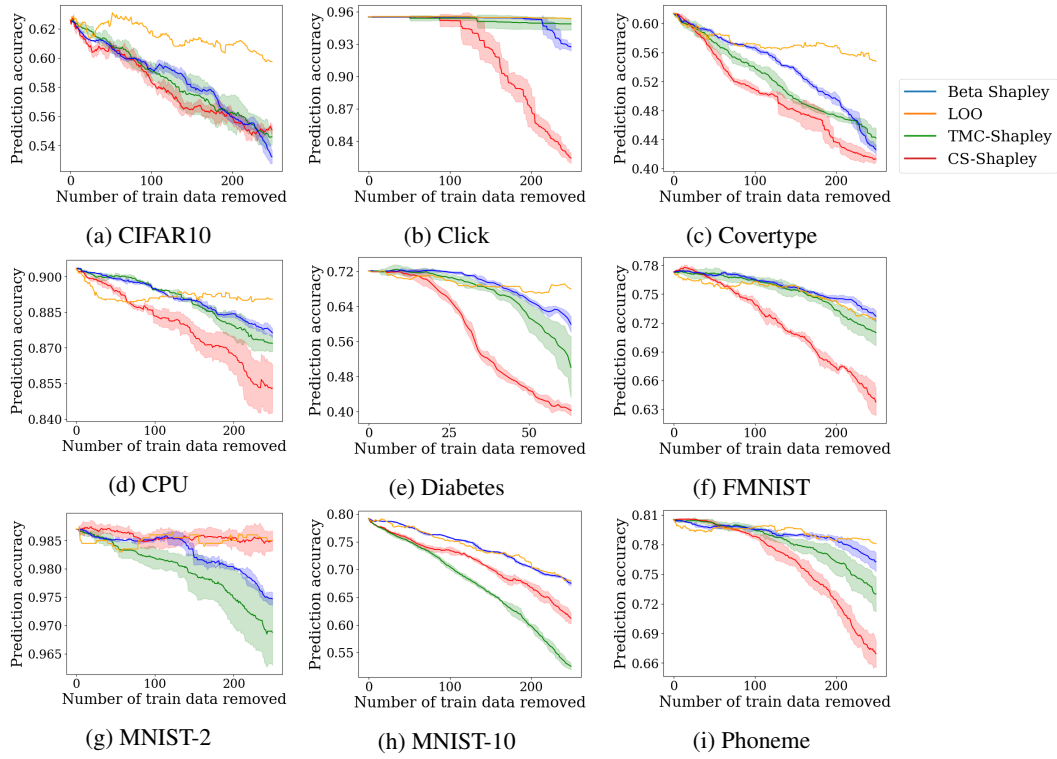


Figure 7: Performance across datasets when removing high-value instances for KNN.

## C.2 Additional noisy label detection results

We report the AUC for gradient boosting classifier and KNN. To help visualize the relationship between AUC and the PR-curve, we also include the PR-curves for the logistic regression and SVM-RBF AUC results reported in the main content. Overall, relative performance ranking of the methods is consistent for this task and high-removal task for gradient boosting classifier and is variable on KNN.

Table 7: Area Under the Curve (AUC) for gradient boosting classifier and KNN using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Gradient Boosting				KNN			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.328</b>	0.318	0.311	0.297	0.324	0.313	<b>0.361</b>	0.299
Click	<b>0.774</b>	0.768	0.678	0.197	0.746	<b>0.808</b>	0.803	0.219
Coverttype	0.627	<b>0.711</b>	0.598	0.182	<b>0.685</b>	0.440	0.437	0.232
CPU	<b>0.664</b>	0.624	0.456	0.193	0.558	<b>0.632</b>	0.447	0.187
Diabetes	<b>0.327</b>	0.319	0.288	0.198	<b>0.330</b>	0.320	0.320	0.192
FMNIST	<b>0.416</b>	0.377	0.374	0.300	0.364	0.453	<b>0.503</b>	0.319
MNIST-2	0.584	<b>0.604</b>	0.500	0.277	0.705	0.728	<b>0.806</b>	0.263
MNIST-10	—	—	—	—	0.796	<b>0.839</b>	0.764	0.326
Phoneme	<b>0.522</b>	0.510	0.414	0.197	0.457	<b>0.479</b>	0.429	0.180

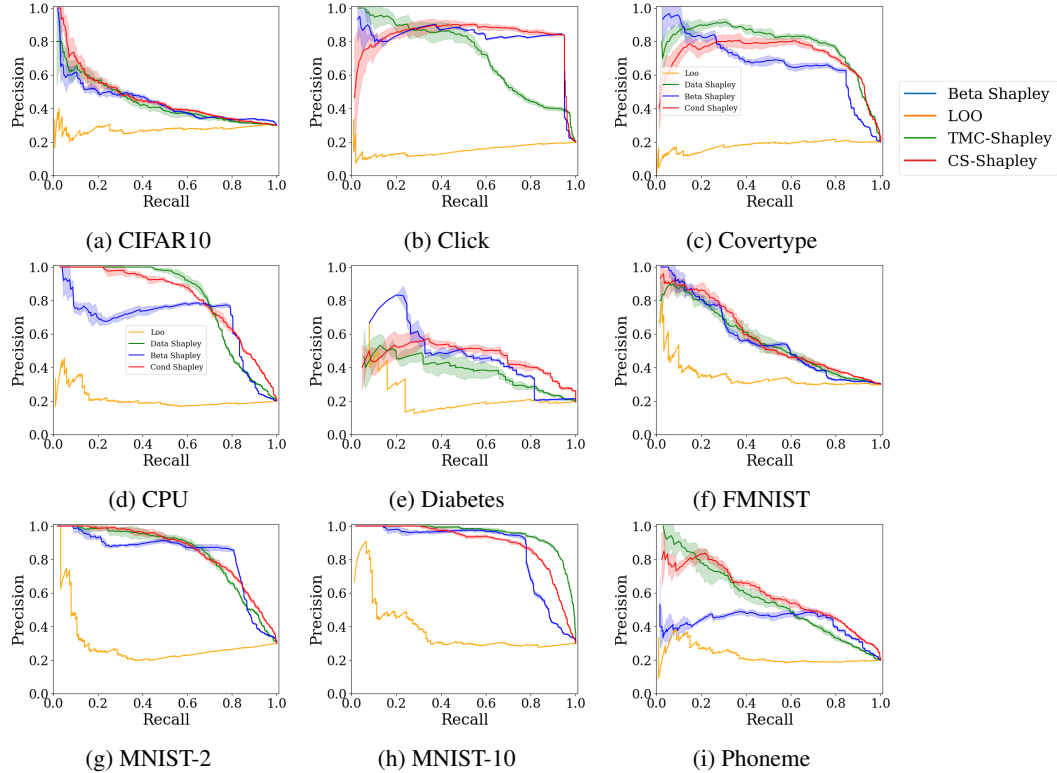


Figure 8: PR-Curve for noisy label detection using logistic regression.

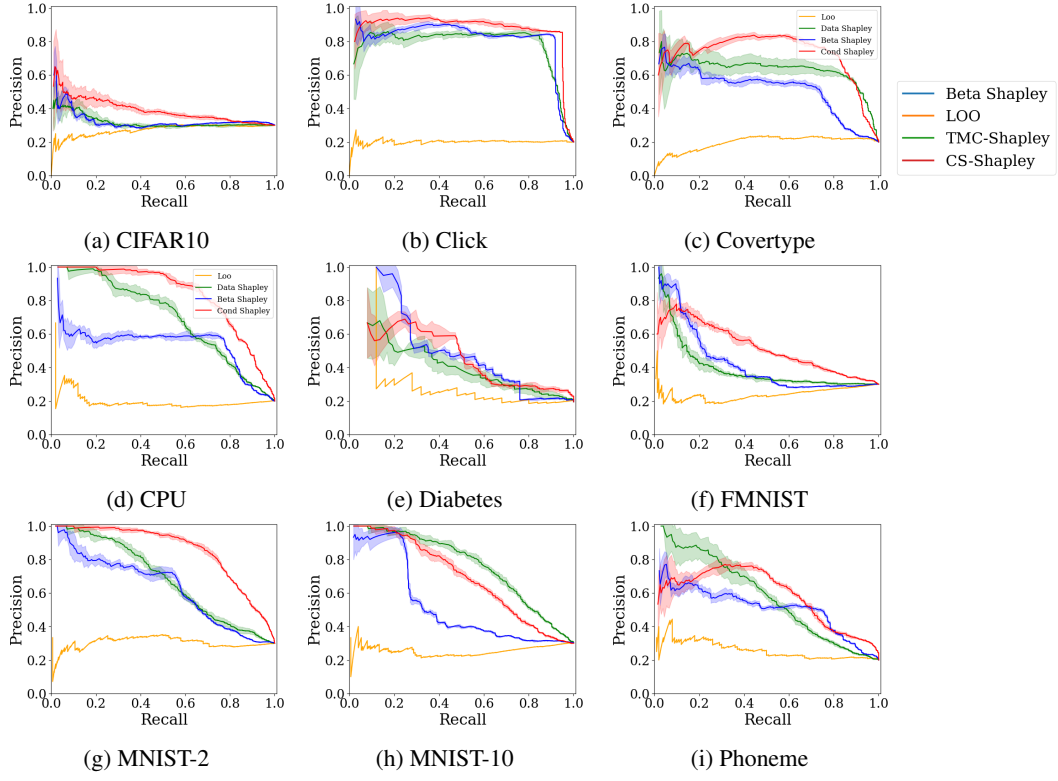


Figure 9: PR-Curve for noisy label detection using SVM with the RBF kernel.

### C.3 Additional transferability of data values results

While we are primarily interested in transfer performance to otherwise computationally infeasible classifiers (i.e. transfer to MLP), we report all transfer combinations to gauge general patterns of transferability across even relatively simple classifiers. Specifically, we report the weighted accuracy drop WAD for each source classifier (logistic regression, SVM-RBF, gradient boosting classifier, and KNN) when transferring to each target classifier (logistic regression, SVM-RBF, gradient boosting classifier, KNN, MLP). Additionally, we show these results visually by plotting the prediction accuracy as a function of training instances removed. Note that results where the source classifier and target classifier are the same are reported in the high-value data removal task. Overall, KNN shows the weakest potential for value transferability. Performance is generally weaker (including more instances where the LOO simple baseline outperforms all Shapley-based methods), and the accuracy drop is generally lower and less stable from the source classifier across most target classifiers in comparison to other source classifiers.

Table 8: Weighted accuracy drop for transferring logistic regression values to SVM-RBF and KNN using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	SVM-RBF				KNN			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.094</b>	0.092	0.071	0.035	0.055	0.070	<b>0.079</b>	-0.016
Click	0.004	0.004	0.004	0.004	<b>0.039</b>	0.005	0.004	0.002
Coverttype	<b>0.210</b>	0.187	0.115	0.073	<b>0.248</b>	0.187	0.109	0.126
CPU	<b>0.030</b>	0.014	0.013	0.022	<b>0.029</b>	0.012	0.008	-0.001
Diabetes	<b>0.113</b>	0.048	0.010	0.008	<b>0.182</b>	0.066	0.051	-0.008
FMNIST	<b>0.068</b>	0.060	0.025	0.033	0.015	<b>0.030</b>	0.015	0.005
MNIST-2	<b>0.103</b>	0.007	0.009	0.010	0.008	0.009	0.006	<b>0.010</b>
MNIST-10	0.186	<b>0.201</b>	0.080	0.076	0.140	<b>0.197</b>	0.083	0.052
Phoneme	<b>0.145</b>	0.023	0.021	0.041	<b>0.122</b>	0.022	0.019	0.024

Table 9: Weighted accuracy drop for transferring logistic regression values to gradient boosting classifier and multi-layer perceptron (MLP) using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Gradient Boosting				MLP			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.111</b>	0.082	0.045	0.001	<b>0.101</b>	0.096	0.078	0.036
Click	<b>0.022</b>	0.007	0.010	0.013	<b>0.095</b>	0.013	0.032	0.019
Coverttype	0.148	<b>0.162</b>	0.059	0.073	<b>0.233</b>	0.180	0.065	0.136
CPU	<b>0.149</b>	0.001	0.004	0.011	<b>0.038</b>	0.009	0.009	0.002
Diabetes	<b>0.109</b>	0.090	0.073	0.061	<b>0.210</b>	0.102	0.064	0.008
FMNIST	0.049	<b>0.056</b>	0.040	0.035	<b>0.083</b>	0.072	0.021	0.017
MNIST-2	0.005	0.004	0.013	<b>0.016</b>	0.004	0.002	0.005	<b>0.010</b>
MNIST-10	0.201	<b>0.211</b>	0.155	0.053	<b>0.191</b>	0.145	0.105	0.056
Phoneme	<b>0.135</b>	0.014	0.026	0.046	0.042	0.047	0.048	<b>0.090</b>

Table 10: Weighted accuracy drop for transferring SVM-RBF values to logistic regression and KNN using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Logistic Regression				KNN			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.084</b>	0.078	0.053	0.071	0.101	<b>0.104</b>	0.083	0.056
Click	0.006	0.008	<b>0.010</b>	0.007	<b>0.005</b>	0.004	0.004	0.002
Coverttype	<b>0.228</b>	0.199	0.084	0.174	<b>0.245</b>	0.171	0.111	0.106
CPU	0.017	0.017	<b>0.038</b>	0.007	0.045	<b>0.047</b>	0.041	-0.011
Diabetes	0.067	0.079	<b>0.096</b>	0.069	<b>0.166</b>	0.073	0.053	0.019
FMNIST	<b>0.060</b>	0.035	0.035	0.054	<b>0.050</b>	0.032	0.021	0.004
MNIST-2	0.007	0.006	0.007	<b>0.010</b>	0.004	<b>0.007</b>	0.005	0.005
MNIST-10	0.110	<b>0.131</b>	0.061	0.065	0.166	<b>0.232</b>	0.089	0.066
Phoneme	<b>0.060</b>	0.004	0.011	-0.032	<b>0.044</b>	0.029	0.035	0.009

Table 11: Weighted accuracy drop for transferring SVM-RBF values to gradient boosting classifier and multi-layer perceptron (MLP) using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Gradient Boosting				MLP			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	0.101	<b>0.126</b>	0.050	0.094	<b>0.130</b>	0.093	0.079	0.071
Click	0.007	0.006	<b>0.009</b>	0.007	<b>0.030</b>	0.013	0.021	-0.002
Coverttype	0.138	<b>0.139</b>	0.056	0.096	<b>0.210</b>	0.142	0.083	0.119
CPU	0.011	0.008	<b>0.021</b>	0.007	0.008	<b>0.018</b>	0.014	-0.001
Diabetes	0.062	<b>0.086</b>	0.051	0.036	<b>0.180</b>	0.104	0.047	0.109
FMNIST	<b>0.075</b>	0.036	0.047	0.025	<b>0.098</b>	0.027	0.031	0.062
MNIST-2	0.008	0.011	0.012	<b>0.018</b>	<b>0.012</b>	0.008	0.002	0.005
MNIST-10	0.226	<b>0.263</b>	0.125	0.104	<b>0.165</b>	0.159	0.081	0.091
Phoneme	<b>0.045</b>	0.041	0.013	0.016	<b>0.055</b>	0.044	0.046	0.017

Table 12: Weighted accuracy drop for transferring KNN values to logistic regression and SVM-RBF using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Logistic Regression				SVM-RBF			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	0.014	0.033	<b>0.035</b>	0.014	0.040	0.051	<b>0.057</b>	0.018
Click	0.002	0.006	<b>0.132</b>	0.007	0.004	0.004	0.004	0.004
Coverttype	<b>0.173</b>	0.163	0.100	0.091	0.160	<b>0.174</b>	0.067	0.109
CPU	0.006	0.011	<b>0.017</b>	0.006	<b>0.044</b>	0.007	0.007	0.007
Diabetes	0.070	0.080	<b>0.107</b>	0.089	<b>0.118</b>	0.051	0.009	0.036
FMNIST	0.021	0.026	0.027	<b>0.046</b>	0.024	0.027	0.029	<b>0.045</b>
MNIST-2	0.010	0.010	0.009	<b>0.013</b>	0.011	<b>0.013</b>	0.010	0.010
MNIST-10	0.023	<b>0.060</b>	0.052	0.059	0.087	<b>0.149</b>	0.077	0.109
Phoneme	<b>0.027</b>	0.005	0.003	0.017	<b>0.049</b>	0.040	0.029	0.010

Table 13: Weighted accuracy drop for transferring KNN values to gradient boosting classifier and multi-layer perceptron (MLP) using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Gradient Boosting				MLP			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	0.044	0.039	<b>0.062</b>	0.032	<b>0.040</b>	0.038	0.016	0.002
Click	0.005	0.006	0.009	0.009	0.025	0.012	<b>0.043</b>	-0.008
Covertypes	<b>0.124</b>	0.109	0.046	0.032	0.179	<b>0.191</b>	0.124	0.080
CPU	<b>0.015</b>	0.004	0.006	0.001	<b>0.027</b>	0.0001	0.012	-0.0001
Diabetes	<b>0.084</b>	0.061	0.065	0.069	<b>0.188</b>	0.093	0.077	0.134
FMNIST	0.042	0.013	<b>0.043</b>	0.025	0.032	0.044	0.025	<b>0.056</b>
MNIST-2	0.009	0.003	0.010	<b>0.015</b>	0.008	0.006	0.008	0.008
MNIST-10	0.131	<b>0.164</b>	0.096	0.147	0.077	0.080	0.098	<b>0.105</b>
Phoneme	0.033	<b>0.036</b>	0.021	0.009	0.053	0.038	<b>0.064</b>	0.055

Table 14: Weighted accuracy drop for transferring gradient boosting classifier values to logistic regression and SVM-RBF using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	Logistic Regression				SVM-RBF			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.091</b>	0.040	0.062	0.020	<b>0.073</b>	0.048	0.070	0.028
Click	0.011	0.006	<b>0.014</b>	0.003	0.004	0.004	0.004	0.004
Covertypes	<b>0.143</b>	0.141	0.070	0.026	0.095	<b>0.103</b>	0.084	0.047
CPU	0.006	<b>0.012</b>	0.004	0.007	<b>0.021</b>	0.013	0.010	0.005
Diabetes	<b>0.097</b>	0.070	0.089	0.009	<b>0.084</b>	0.051	0.009	-0.044
FMNIST	<b>0.060</b>	0.033	0.038	0.021	<b>0.050</b>	0.033	0.023	0.021
MNIST-2	0.008	0.006	0.009	0.009	0.004	0.004	0.003	<b>0.006</b>
MNIST-10	—	—	—	—	—	—	—	—
Phoneme	<b>0.013</b>	-0.002	0.002	0.001	<b>0.042</b>	0.029	0.033	0.020

Table 15: Weighted accuracy drop for transferring gradient boosting classifier values to KNN and multi-layer perceptron (MLP) using CS-SHAPLEY (CS), TMC-Shapley (TMC), Beta Shapley (Beta), and Leave-One-Out (LOO).

Dataset	KNN				MLP			
	CS	TMC	Beta	LOO	CS	TMC	Beta	LOO
CIFAR10	<b>0.067</b>	0.056	0.061	0.019	0.033	<b>0.114</b>	0.019	0.060
Click	<b>0.009</b>	0.005	0.004	0.003	<b>0.035</b>	0.021	0.022	-0.009
Covertypes	0.089	<b>0.114</b>	0.071	0.053	0.138	<b>0.222</b>	0.023	-0.008
CPU	<b>0.019</b>	0.005	0.009	0.017	<b>0.017</b>	0.008	0.008	0.004
Diabetes	<b>0.117</b>	0.087	0.029	0.022	<b>0.178</b>	0.110	0.063	0.071
FMNIST	<b>0.037</b>	0.023	0.017	0.010	<b>0.075</b>	0.028	0.047	0.021
MNIST-2	<b>0.007</b>	0.005	0.005	0.006	<b>0.007</b>	0.006	0.002	0.006
MNIST-10	—	—	—	—	—	—	—	—
Phoneme	<b>0.044</b>	0.028	0.034	0.022	<b>0.080</b>	0.058	0.068	0.025

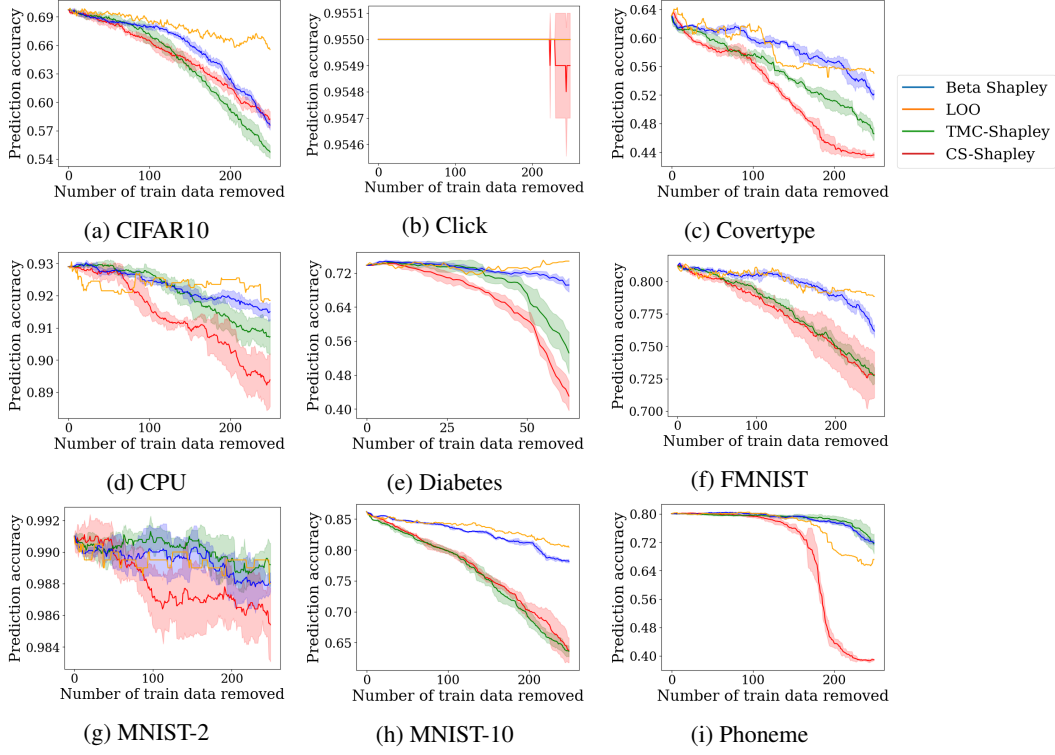


Figure 10: Performance across datasets when transferring values from logistic regression to SVM-RBF for the high-value removal task.

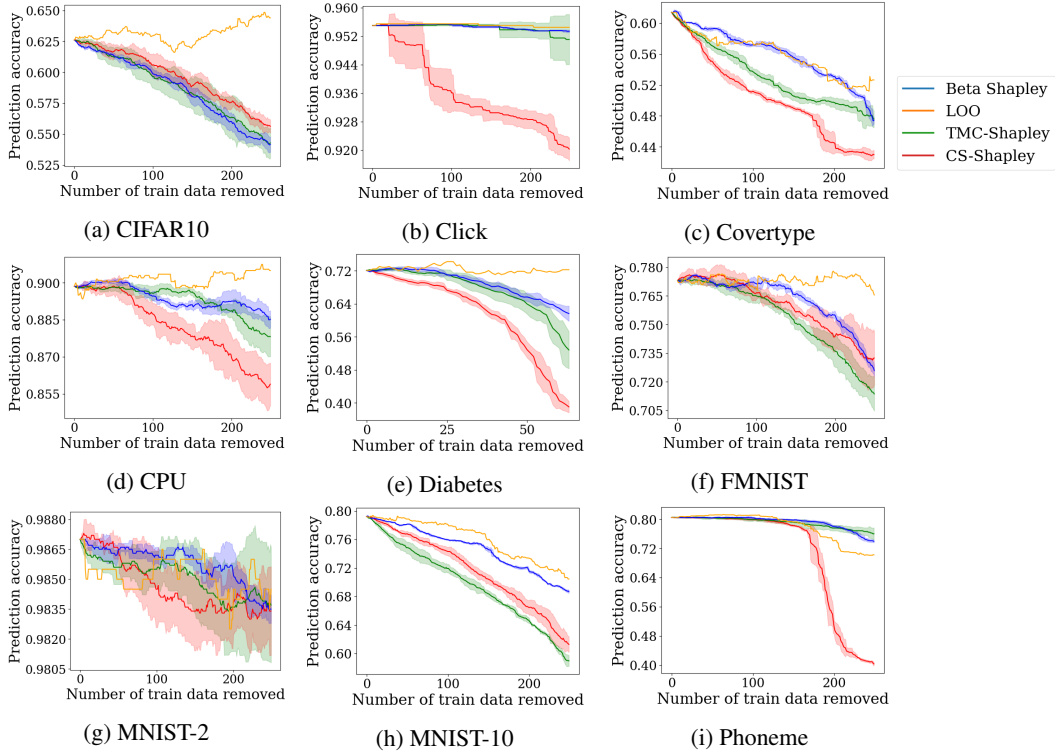


Figure 11: Performance across datasets when transferring values from logistic regression to KNN for the high-value removal task.



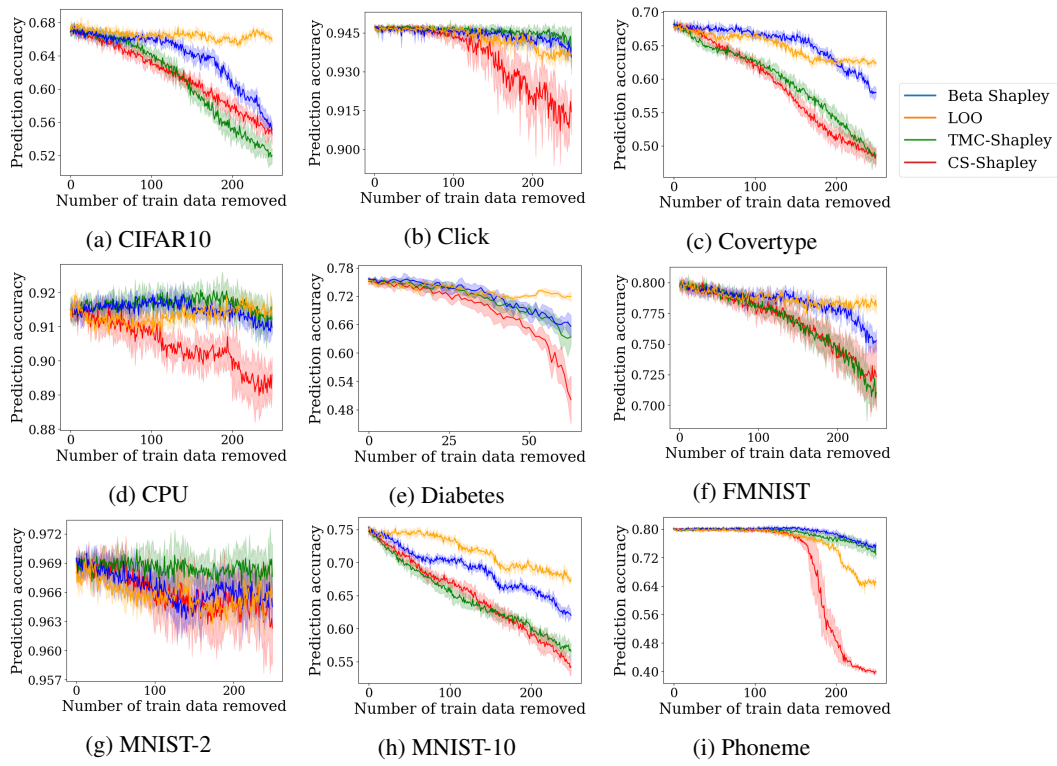


Figure 12: Performance across datasets when transferring values from logistic regression to gradient boosting classifier for the high-value removal task.

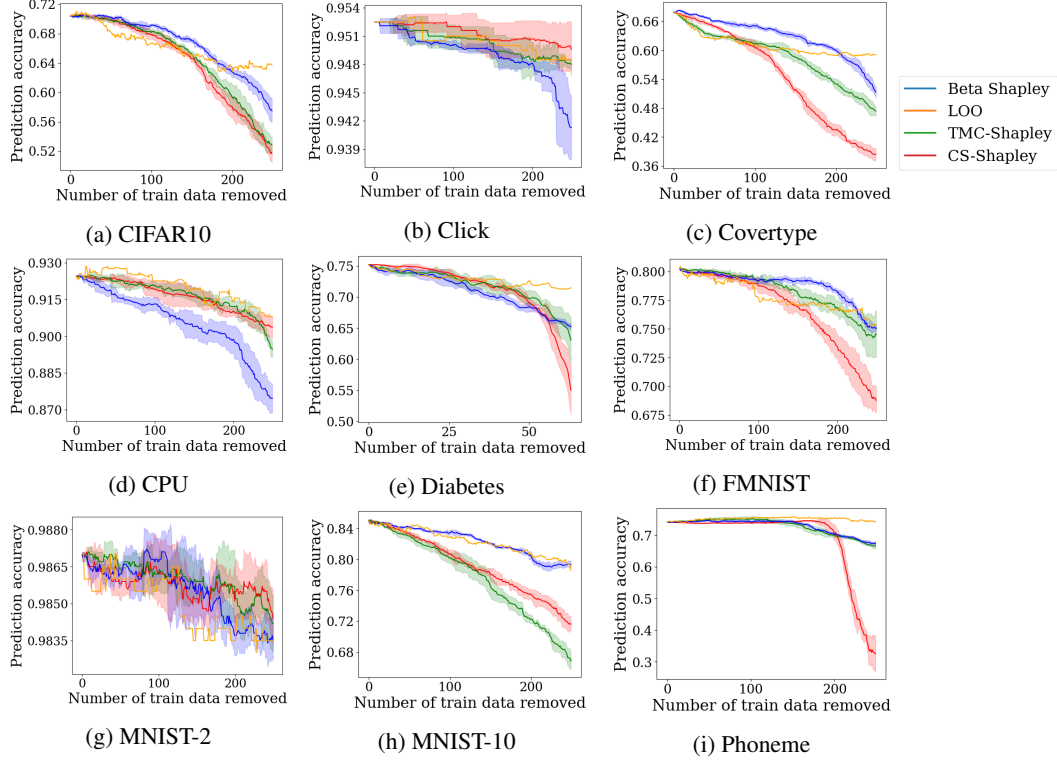


Figure 13: Performance across datasets when transferring values from SVM-RBF to logistic regression for the high-value removal task.

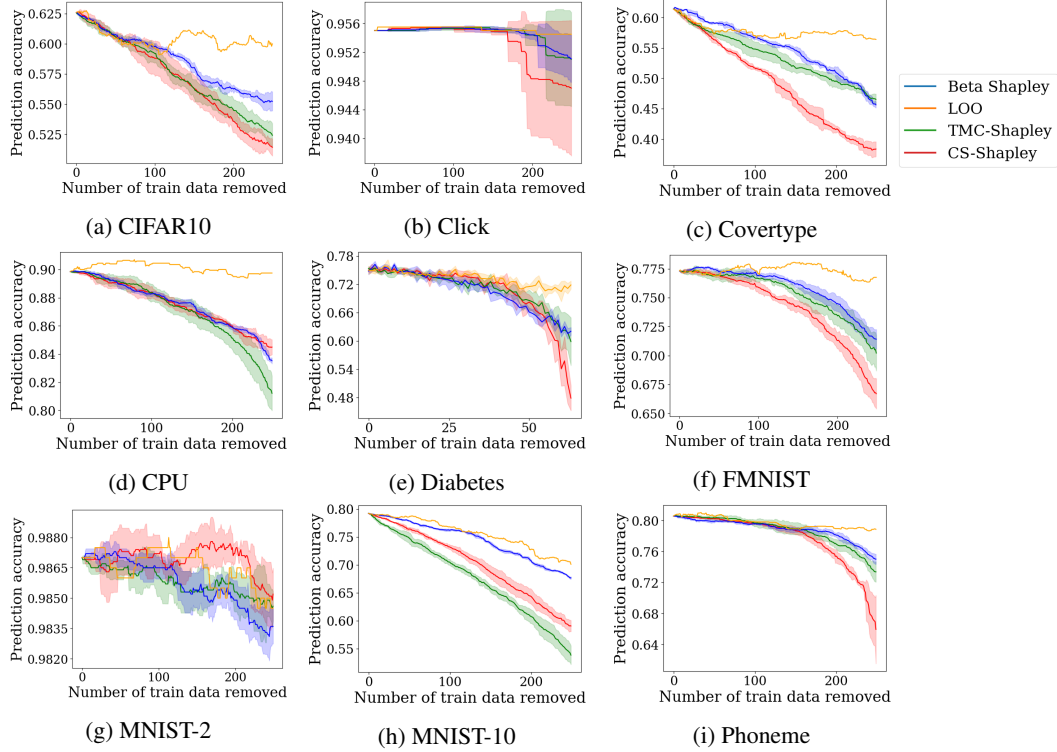


Figure 14: Performance across datasets when transferring values from SVM-RBF to KNN for the high-value removal task.

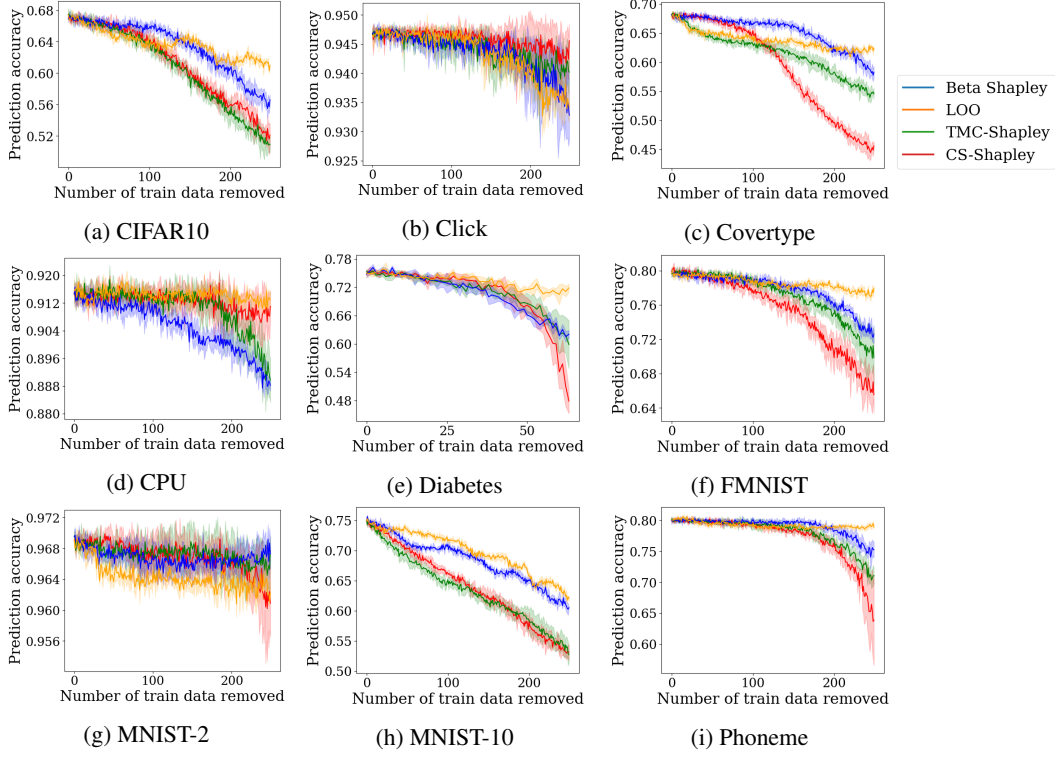


Figure 15: Performance across datasets when transferring values from SVM-RBF to gradient boosting classifier for the high-value removal task.

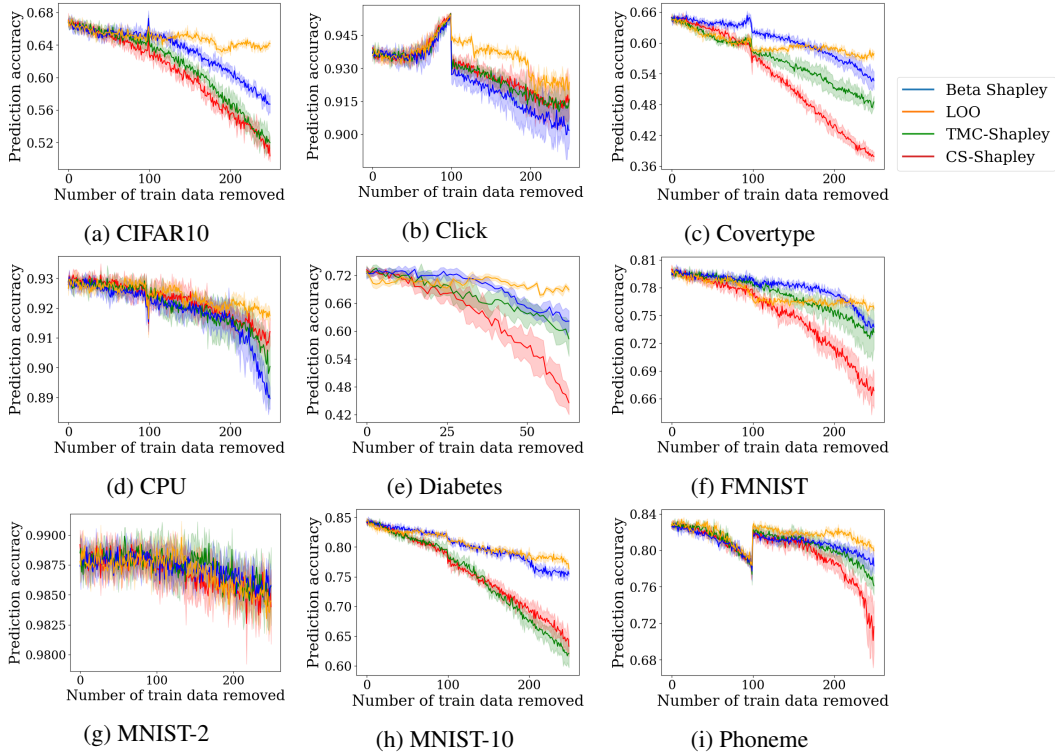


Figure 16: Performance across datasets when transferring values from SVM-RBF to MLP for the high-value removal task.

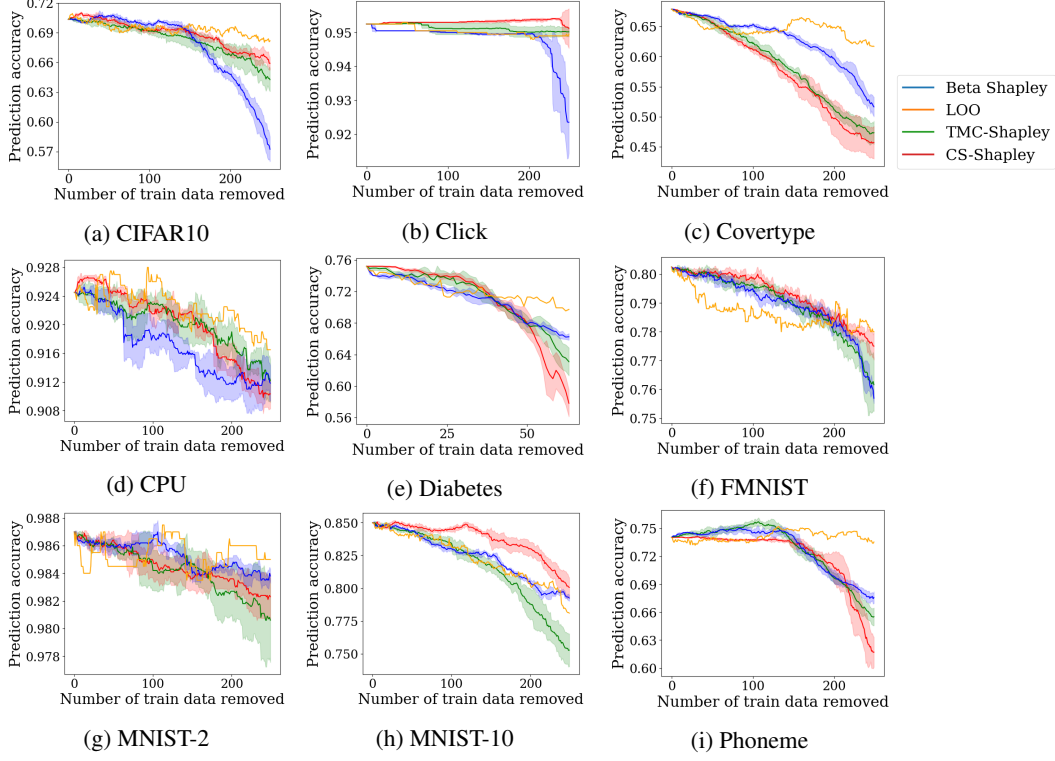


Figure 17: Performance across datasets when transferring values from KNN to logistic regression for the high-value removal task.

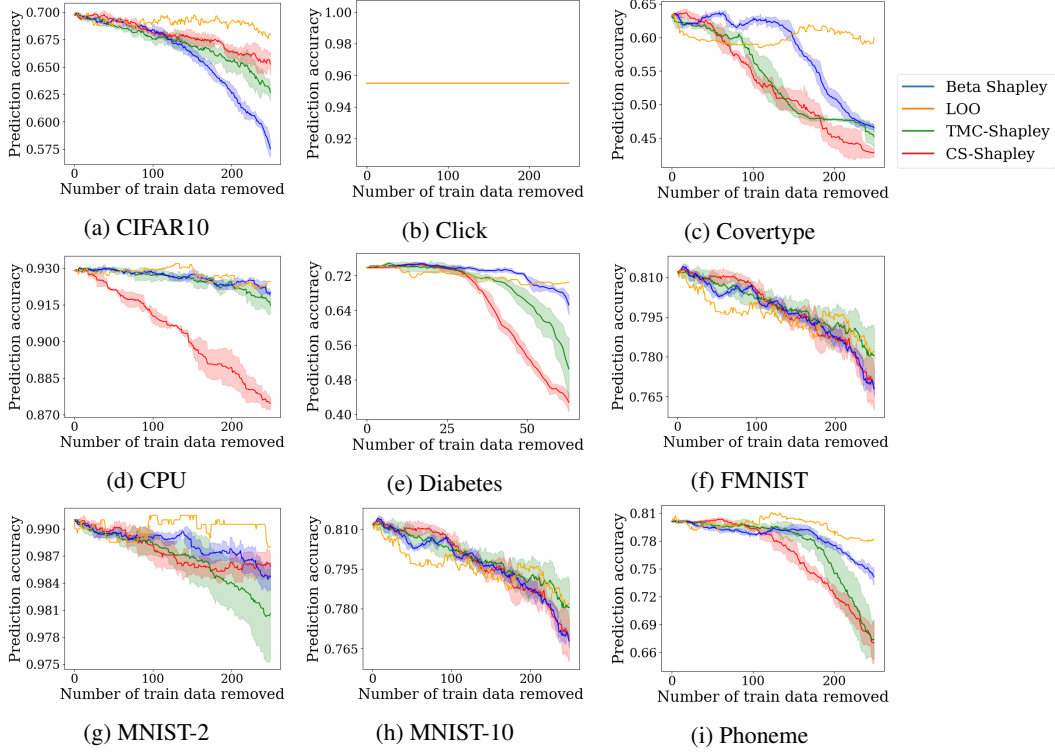


Figure 18: Performance across datasets when transferring values from KNN to SVM-RBF for the high-value removal task.

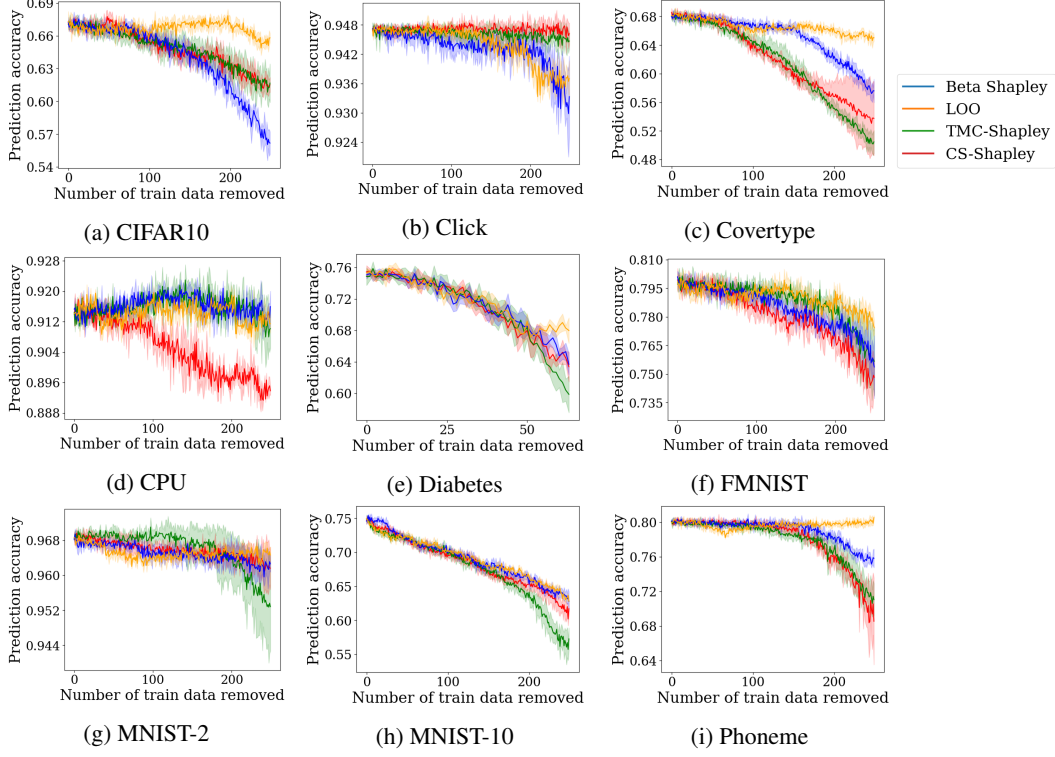


Figure 19: Performance across datasets when transferring values from KNN to gradient boosting classifier for the high-value removal task.

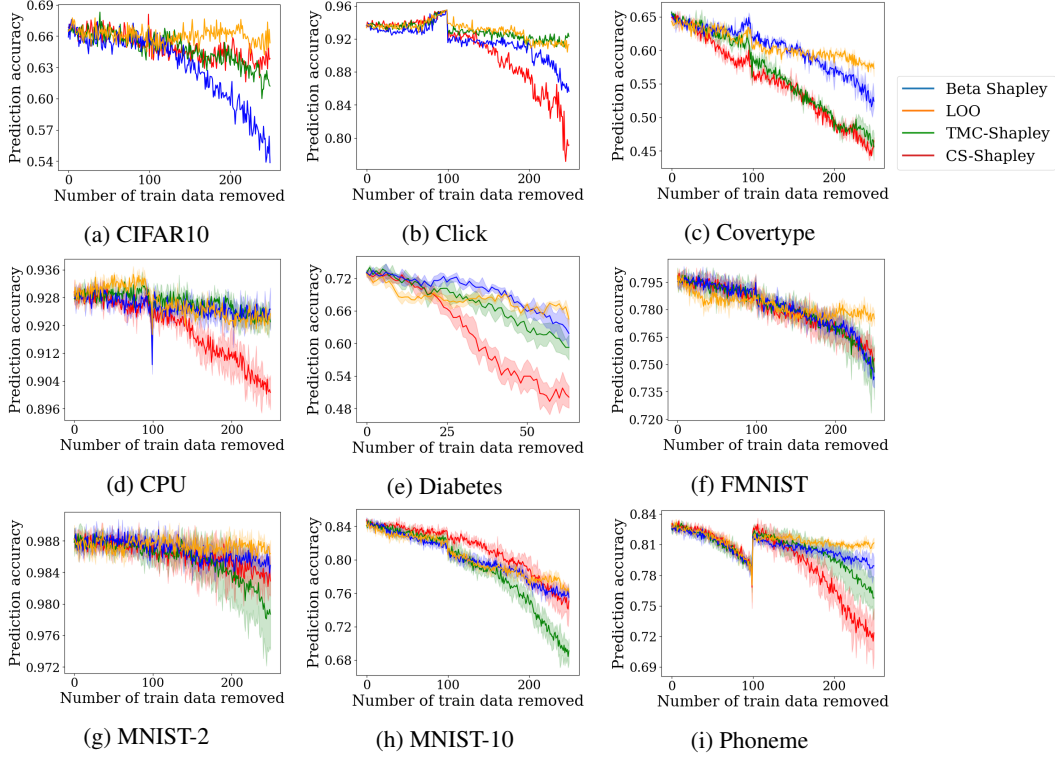


Figure 20: Performance across datasets when transferring values from KNN to MLP for the high-value removal task.

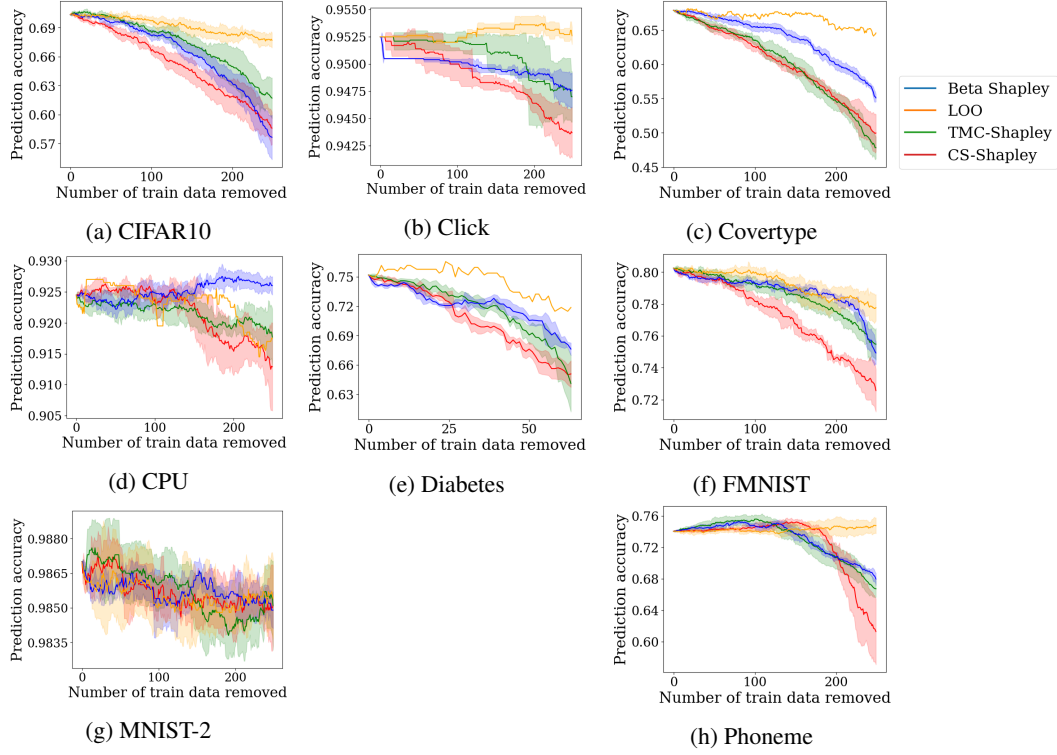


Figure 21: Performance across datasets when transferring values from gradient boosting classifier to logistic regression for the high-value removal task.

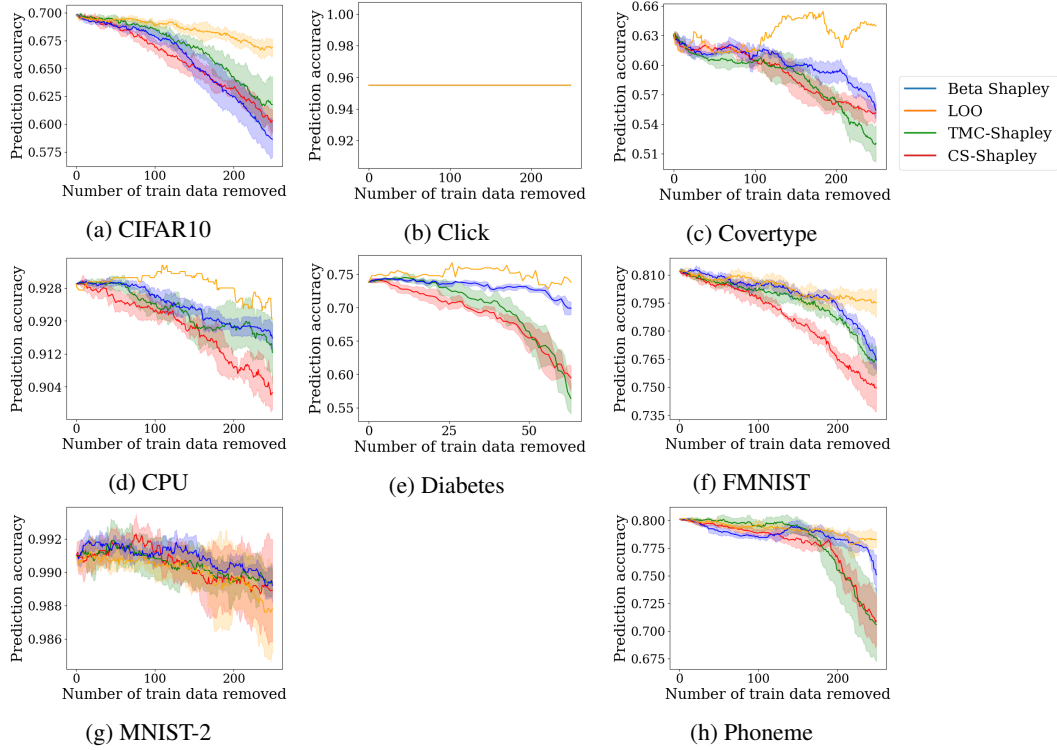


Figure 22: Performance across datasets when transferring values from gradient boosting classifier to SVM-RBF for the high-value removal task.



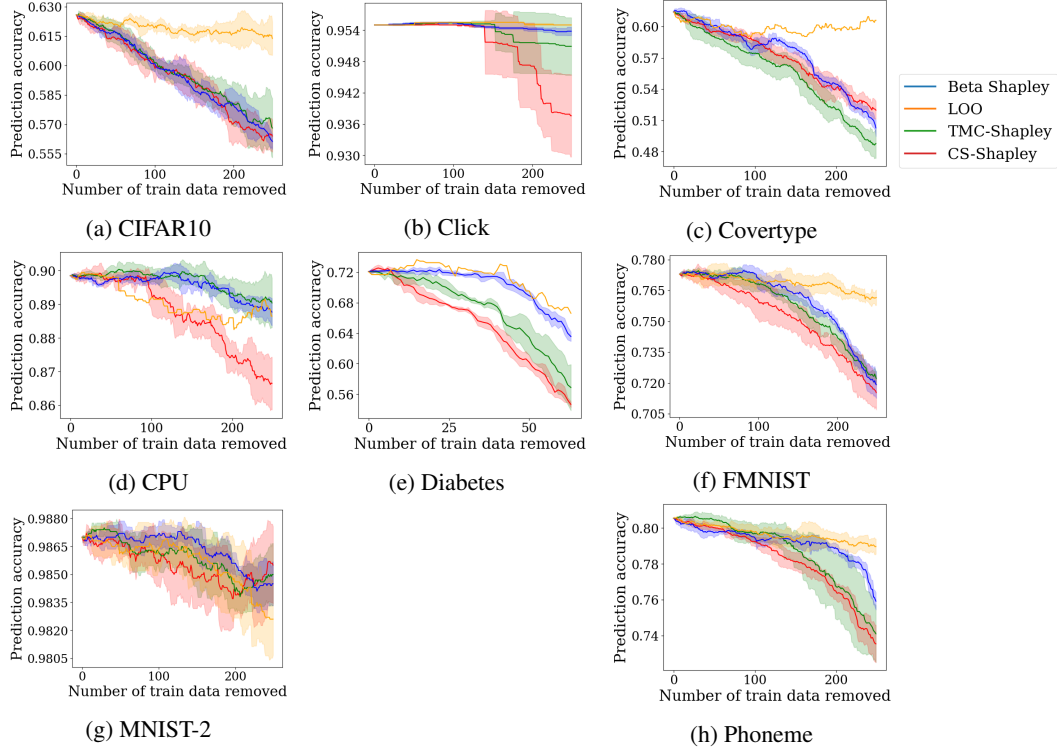


Figure 23: Performance across datasets when transferring values from gradient boosting classifier to KNN for the high-value removal task.

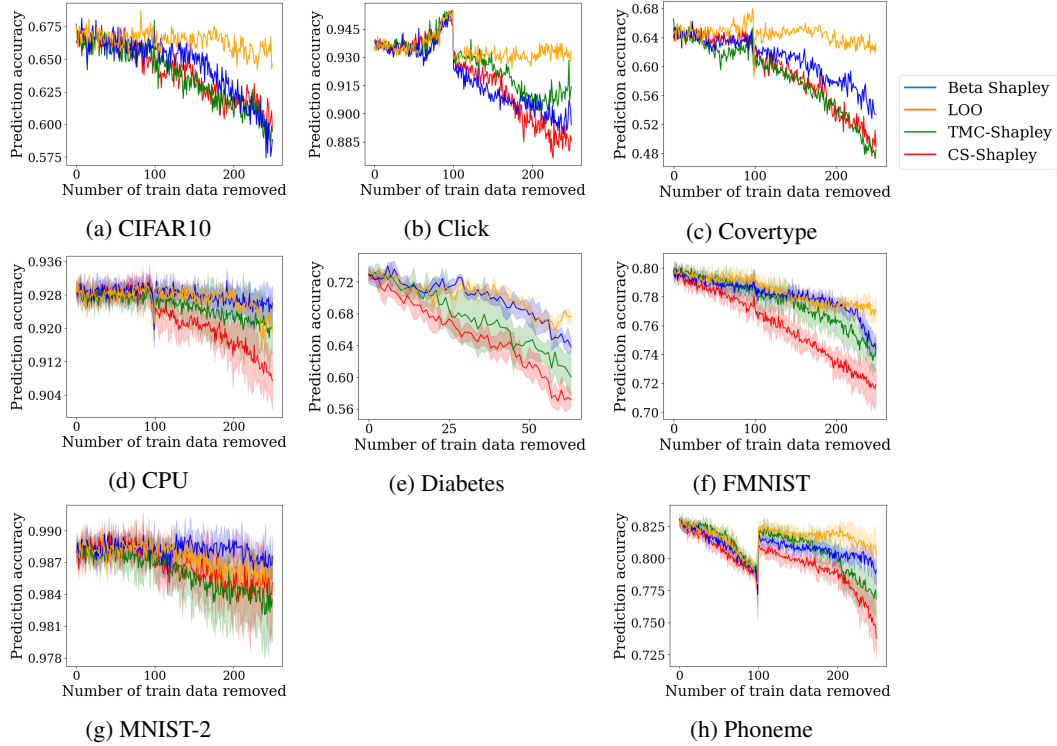


Figure 24: Performance across datasets when transferring values from gradient boosting classifier to MLP for the high-value removal task.