



COMPSCI 389

Introduction to Machine Learning

Evaluation Part 3

Prof. Philip S. Thomas (pthomas@cs.umass.edu)

Overview

- In this lecture we will conduct a series of experiments to highlight:
 - The need for confidence intervals
 - How confidence intervals can be computed
 - The different types of intervals that are often reported
- The code covered in this lecture will be provided in:

`10 Evaluation Part 3.ipynb`

Load data, split train/test, compute predictions and errors

```
# Load the data set
df = pd.read_csv("data/GPA.csv", delimiter=',')

# We already loaded X and y, but do it again as a reminder
X = df.iloc[:, :-1]
y = df.iloc[:, -1]

# Split the data into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.8, shuffle=True)

# Train the NearestNeighbor model
model = KNearestNeighbors(k=1)
model.fit(X_train, y_train)

# Compute predictions for X_test
predictions = model.predict(X_test)

# Compute the sample squared errors
squared_errors = (predictions - y_test) ** 2
display("Squared errors: ", squared_errors)
```

The algorithm used won't be important for this discussion.

Let's use most of the data for evaluation for now.

We compute the squared error for all the testing points (80% of the data). We will "simulate" using less data by looking at portions of these errors.

Load data, split train/test, compute predictions and errors

```
# Load the data set
df = pd.read_csv("data/GPA.csv", delimiter=',')

# We already loaded X and y, but do it again as a reminder
X = df.iloc[:, :-1]
y = df.iloc[:, -1]

# Split the data into training and testing sets (20% train, 80% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.8, shuffle=True)

# Train the NearestNeighbor model
model = KNearestNeighbors(k=1)
model.fit(X_train, y_train)

# Compute predictions for X_test
predictions = model.predict(X_test)

# Compute the sample squared errors
squared_errors = (predictions - y_test) ** 2
display("Squared errors: ", squared_errors)
```

The algorithm used won't be important for this discussion.

Let's use most of the data for evaluation for now.

We compute the squared error for all the testing points (80% of the data). We will "simulate" using less data by looking at portions of these errors.

```
squared_errors
```

```
16809    2.170701
```

```
31412    0.030047
```

```
27099    0.277381
```

```
41341    0.440007
```

```
27941    0.523206
```

```
...
```

```
2874     0.006944
```

```
21169    0.017777
```

```
39080    0.572534
```

```
24225    3.635379
```

```
22171    0.005377
```

```
Name: gpa, Length: 34643, dtype: float64
```

Idea

- We would like to see what happens if we have test sets of different sizes.
- We would like to see what happens if we have different samples in the test set.
- Ideally, we would generate completely new test sets (and squared errors) for each experiment.
 - We only have 43,303 points total and cannot generate more.
- Instead, we use random subsets of the test set (80% of the data)
 - We vary the size of these subsets.
 - We pretend that random subsets are independent (they are not)
 - Note: This will be important at the end!

Compute the sample MSE from percentage% of the squared errors.

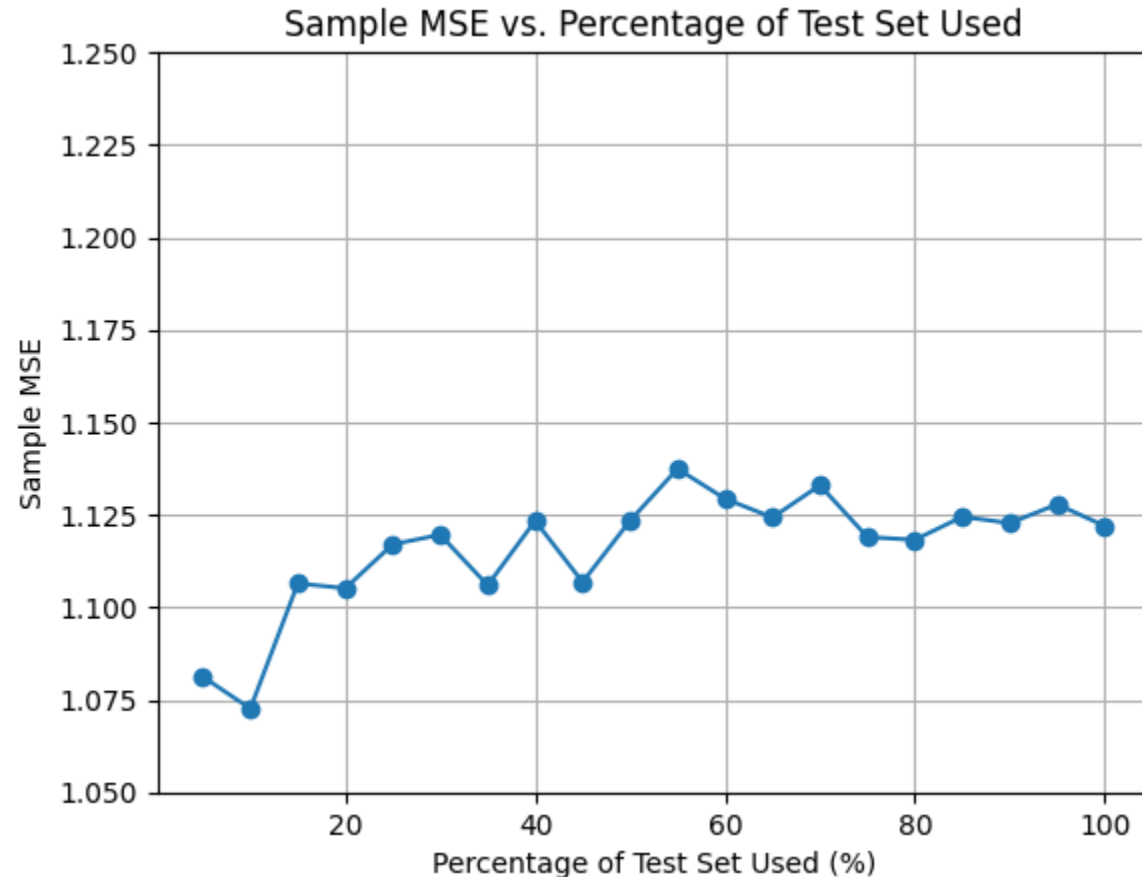
```
# Function to compute average MSE for a given percentage of squared errors
def compute_average_mse(percentage, squared_errors):
    # Get the number of samples that we should use from squared_errors
    subset_size = int(percentag * len(squared_errors))

    # Randomly select that many indices (without replacement)
    indices = np.random.choice(len(squared_errors), subset_size, replace=False)
    # Randomly select percentage% of the points, so each call produces results from a different simulated sample
    # Get the average of the squared errors at the selected indices
    average_mse = squared_errors.iloc[indices].mean()
    return average_mse
```

Plot: Sample MSE using different amounts of the test data.

Note: The amount of **training** data is not varied.

This represents the setting where a fixed model is given, and we are trying to evaluate that model, regardless of how it was created.

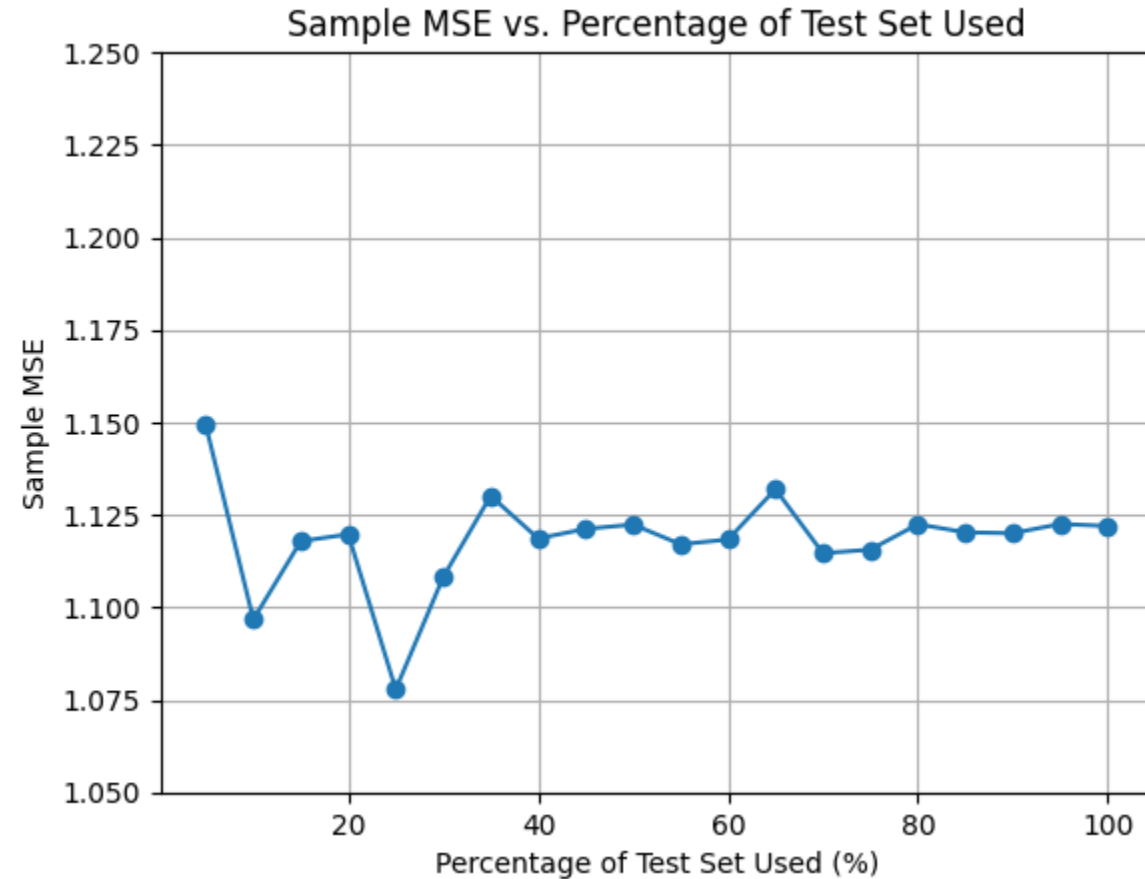


Note: In practice, we would only have *one point* from this plot: the sample MSE using our one testing set. This plot is meant to give further insight about how much we can (not) trust that value.

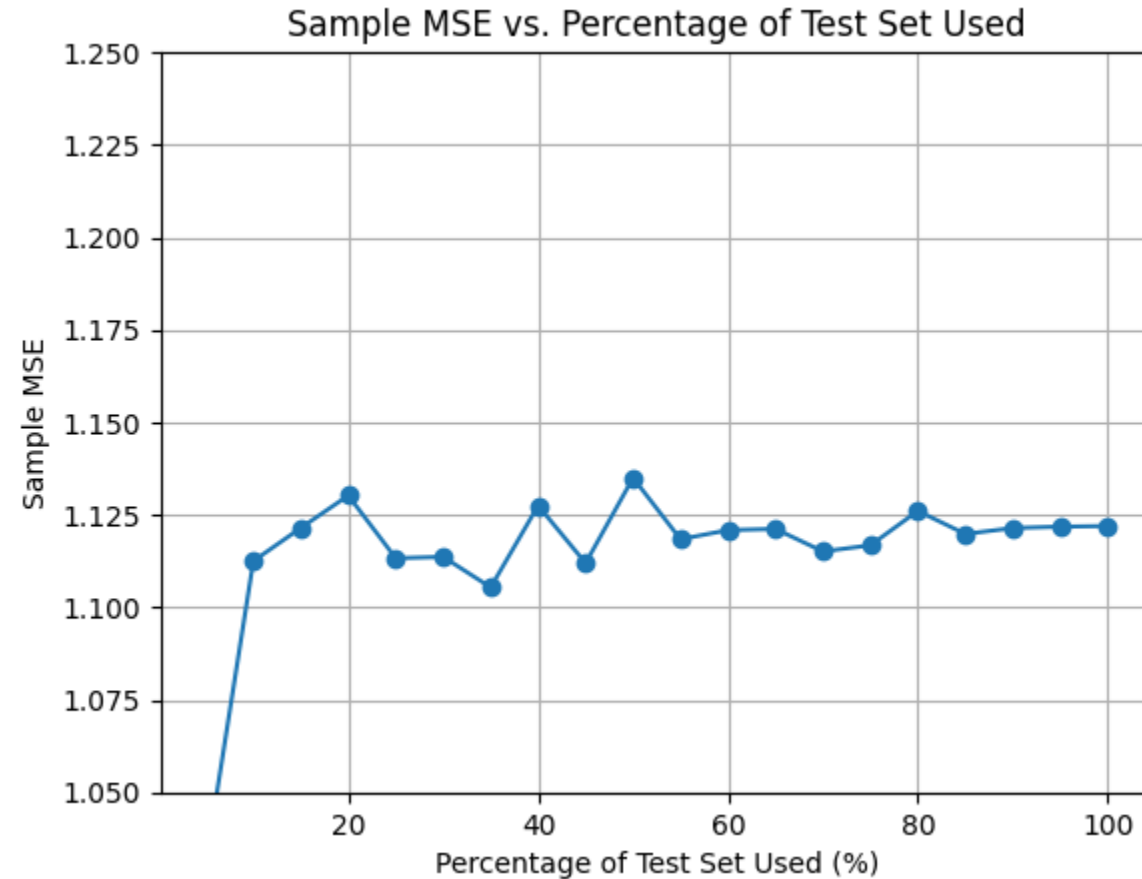
Let's rerun this several times and see what happens.

Simulates the sample MSE we would get with different amounts of testing data

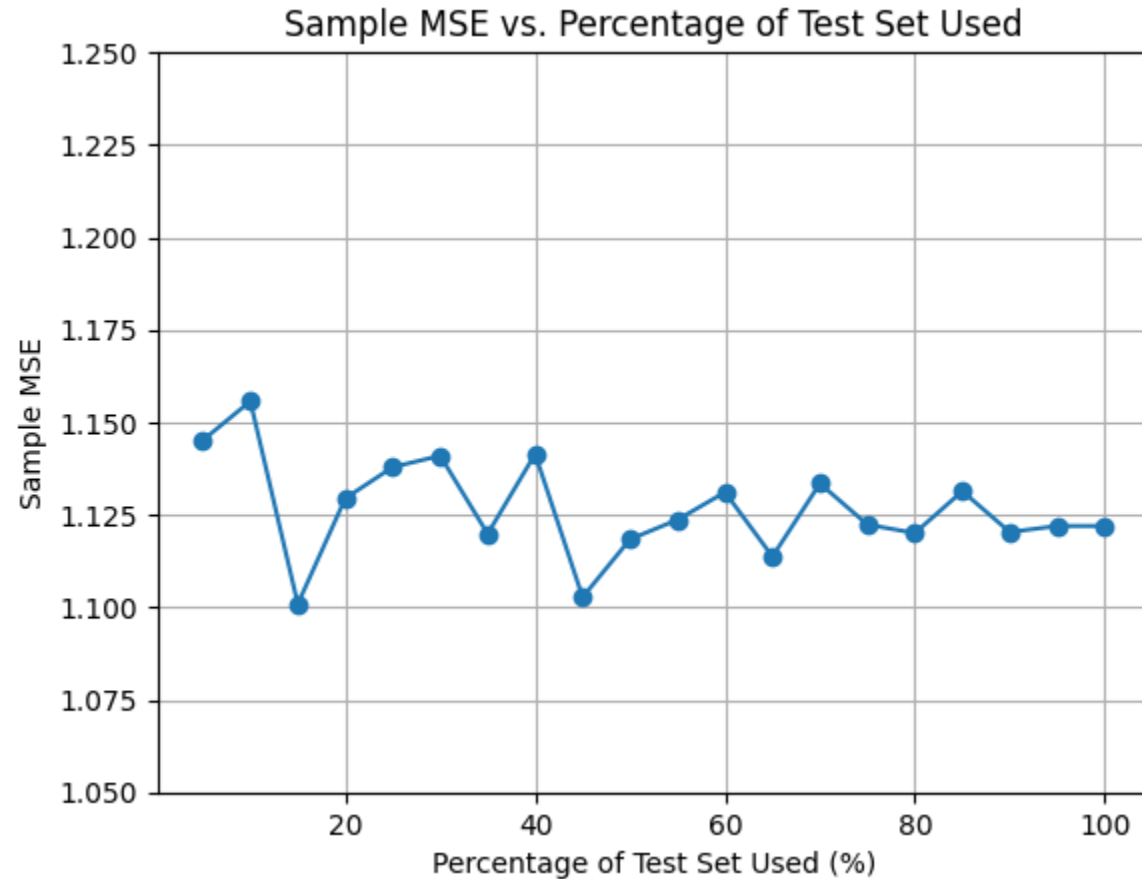
Plot: Sample MSE using different amounts of the test data.



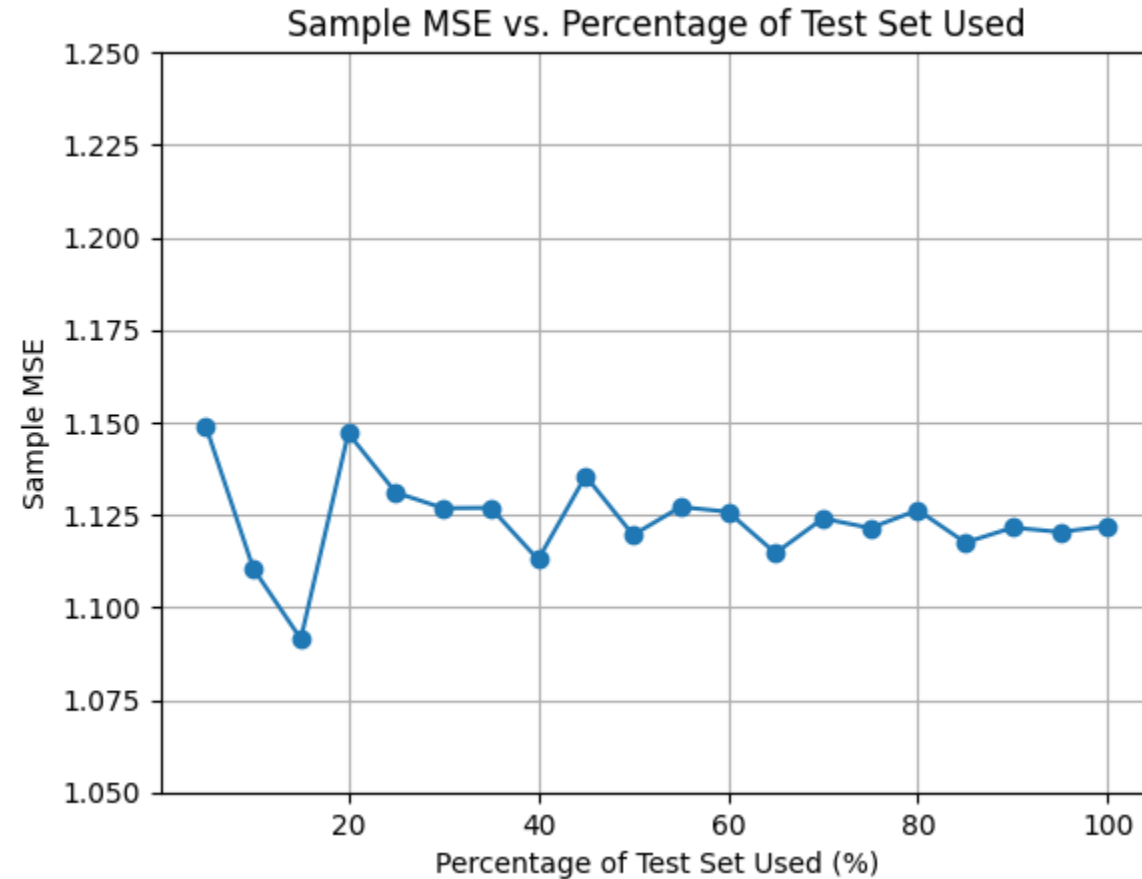
Plot: Sample MSE using different amounts of the test data.



Plot: Sample MSE using different amounts of the test data.

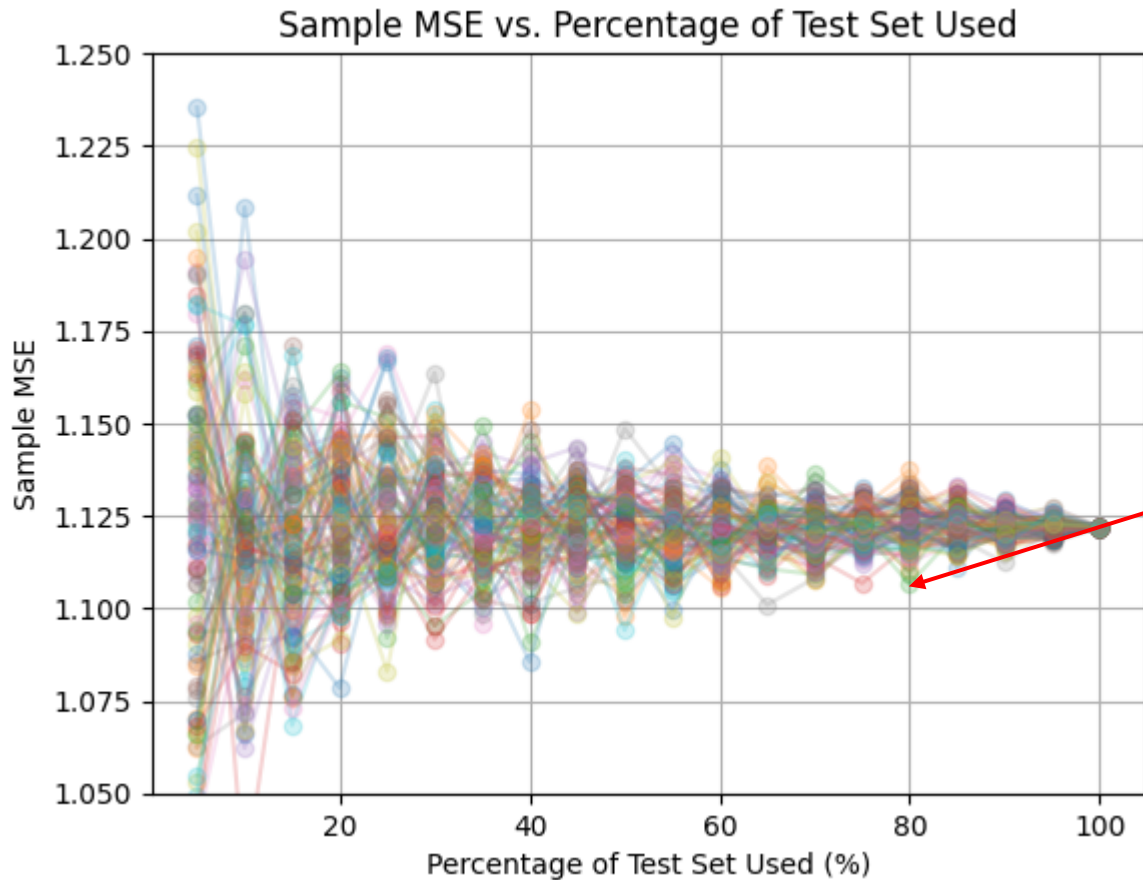


Plot: Sample MSE using different amounts of the test data.



Notice: There was significant variation across runs!

Overlay of 100 runs



Each entry in the tables from before showed *one single point* from this plot.

	Model	MSE	RMSE	MAE	R ²
0	k-NN k=1 sigma=None	1.152084	1.073352	0.823743	-0.687769
1	k-NN k=100 sigma=None	0.579404	0.761186	0.596919	0.151190
2	k-NN k=100 sigma=100	0.579572	0.761297	0.596952	0.150943
3	k-NN k=200 sigma=100	0.577554	0.759970	0.596220	0.153901
4	k-NN k=300 sigma=100	0.577443	0.759897	0.596408	0.154062
5	k-NN k=400 sigma=100	0.577620	0.760013	0.596670	0.153804
6	k-NN k=500 sigma=100	0.578077	0.760314	0.597044	0.153135

Note: It isn't precisely the indicated point, since this old table was for a model trained with more data.

Variance

- The previous plots show the *variance* of the sample statistics (sample metrics).
- When we compute sample statistics from data, they are not exactly equal to the population statistic (parameter).
- If we had many data sets and computed the sample statistics many times, we would see this variance.
- With only one data set, we get our estimate (e.g., the sample MSE)
- How much should we trust this estimate?
 - We can use *standard error* to get an idea.

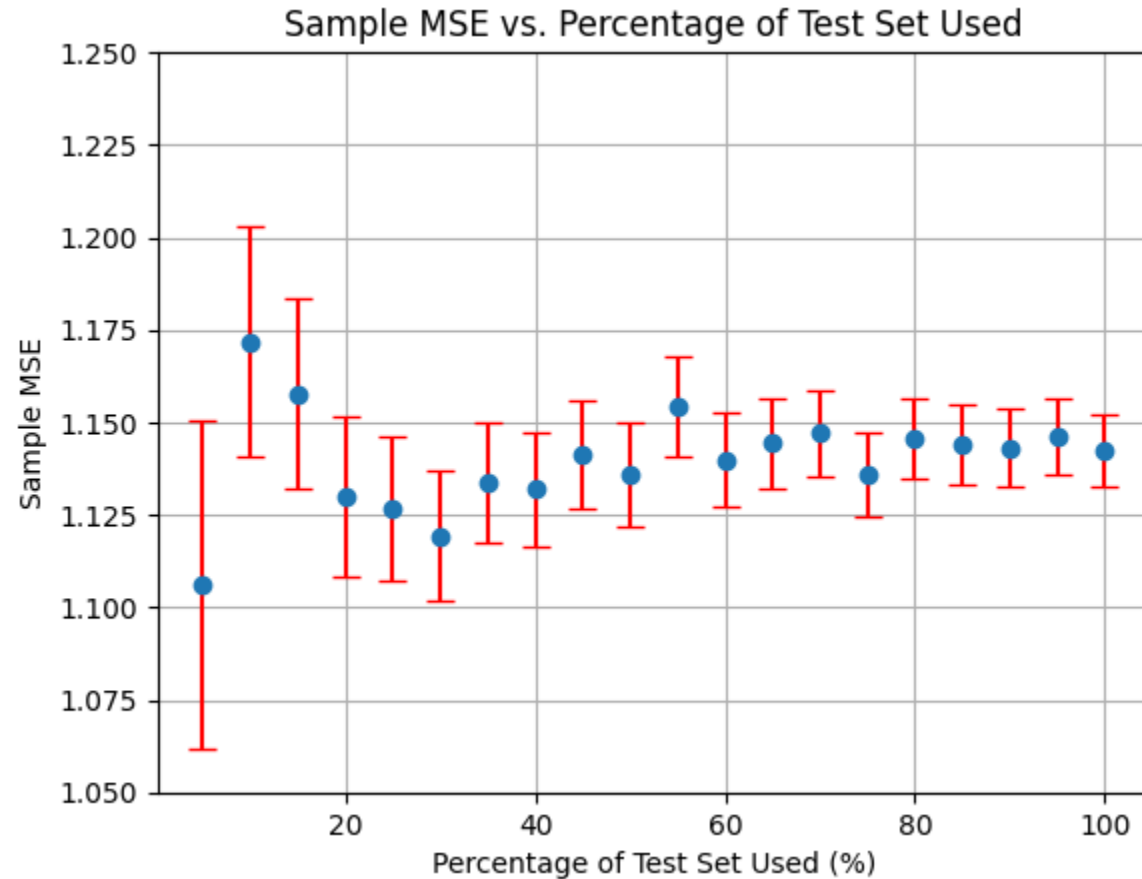
Error Bars: Standard Error

- We can plot error bars on each point
- The width of each error bar is the *standard error*
 - The point is the average of the `squared_error` samples
 - The error bar width is the standard error of the `squared_error` samples
 - Recall:

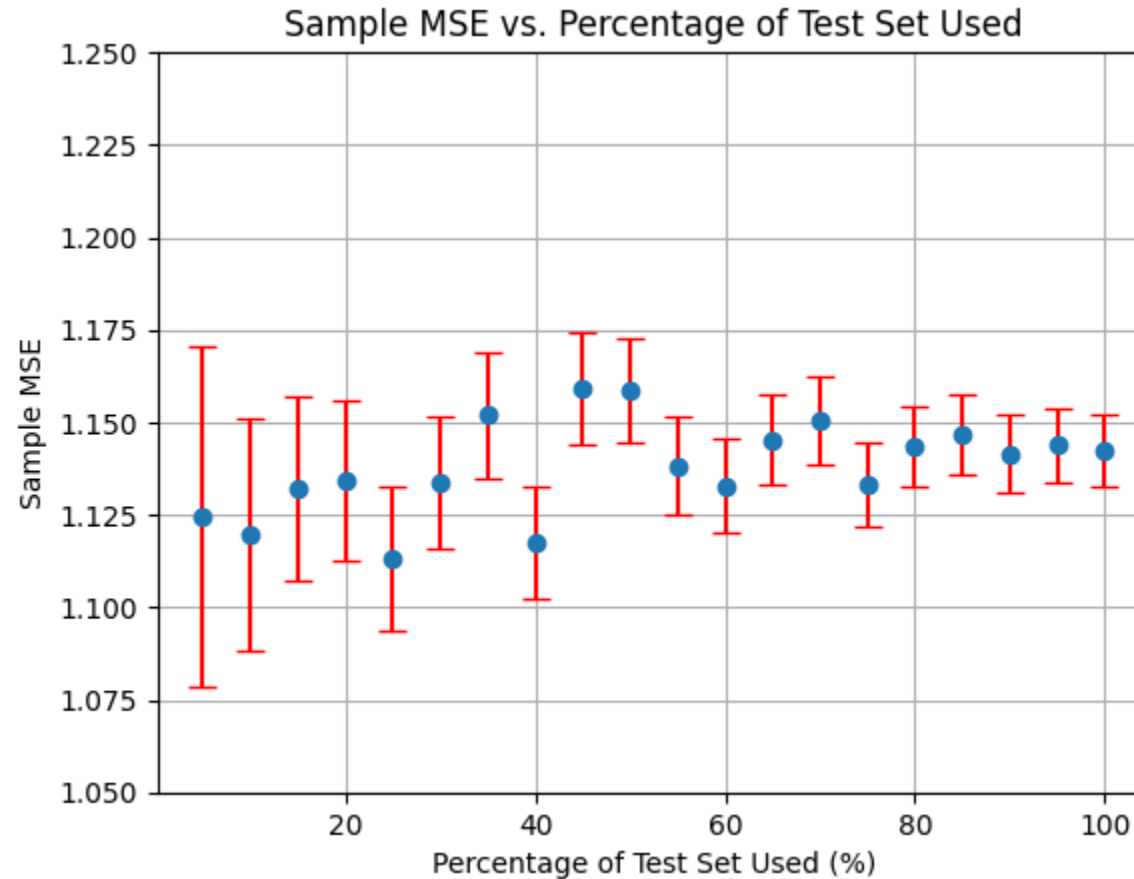
$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad SE = \frac{\sigma}{\sqrt{n}}$$

- Here, x_i is one of the `squared_error` values.
- The standard deviation σ captures how much the squared errors vary.
- $\pm 1.96 \times SE$ gives a 95% confidence interval (under normality assumptions)

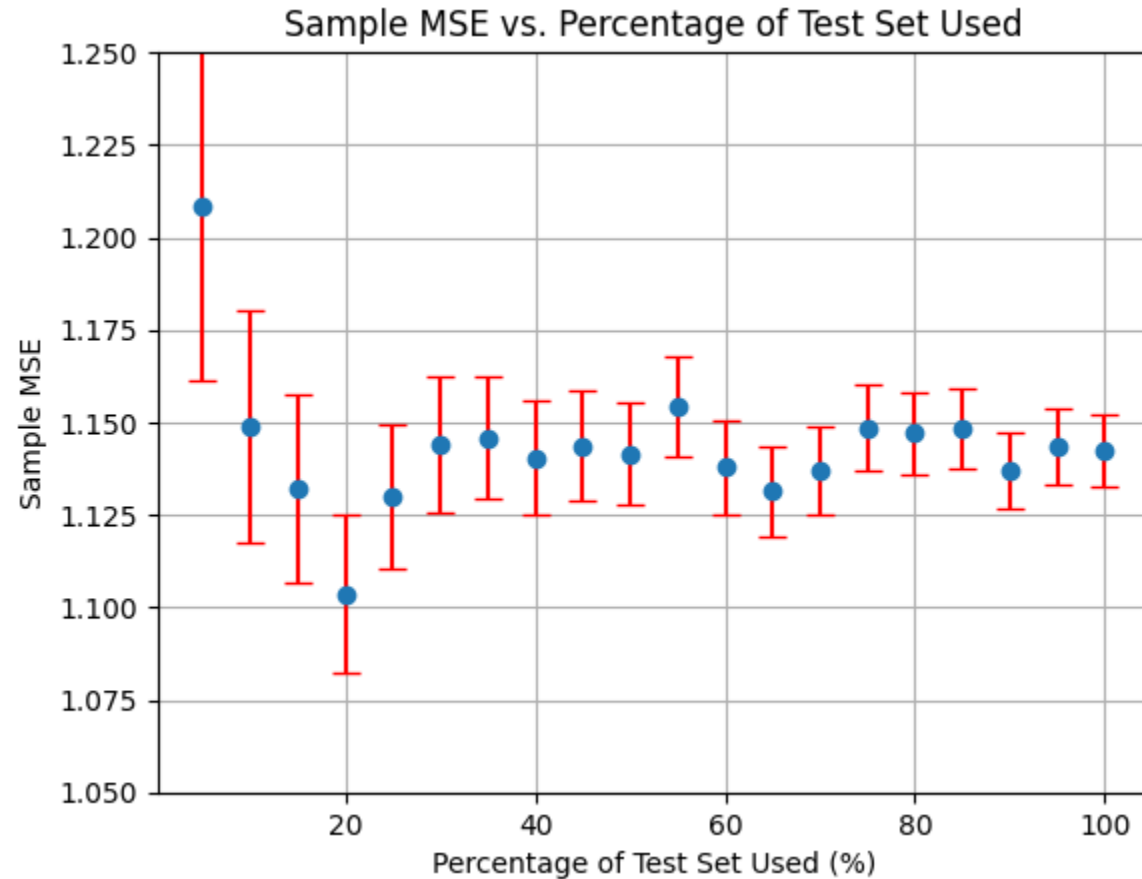
Plot: Sample MSE using different amounts of the test data, **with error bars (standard error)**.



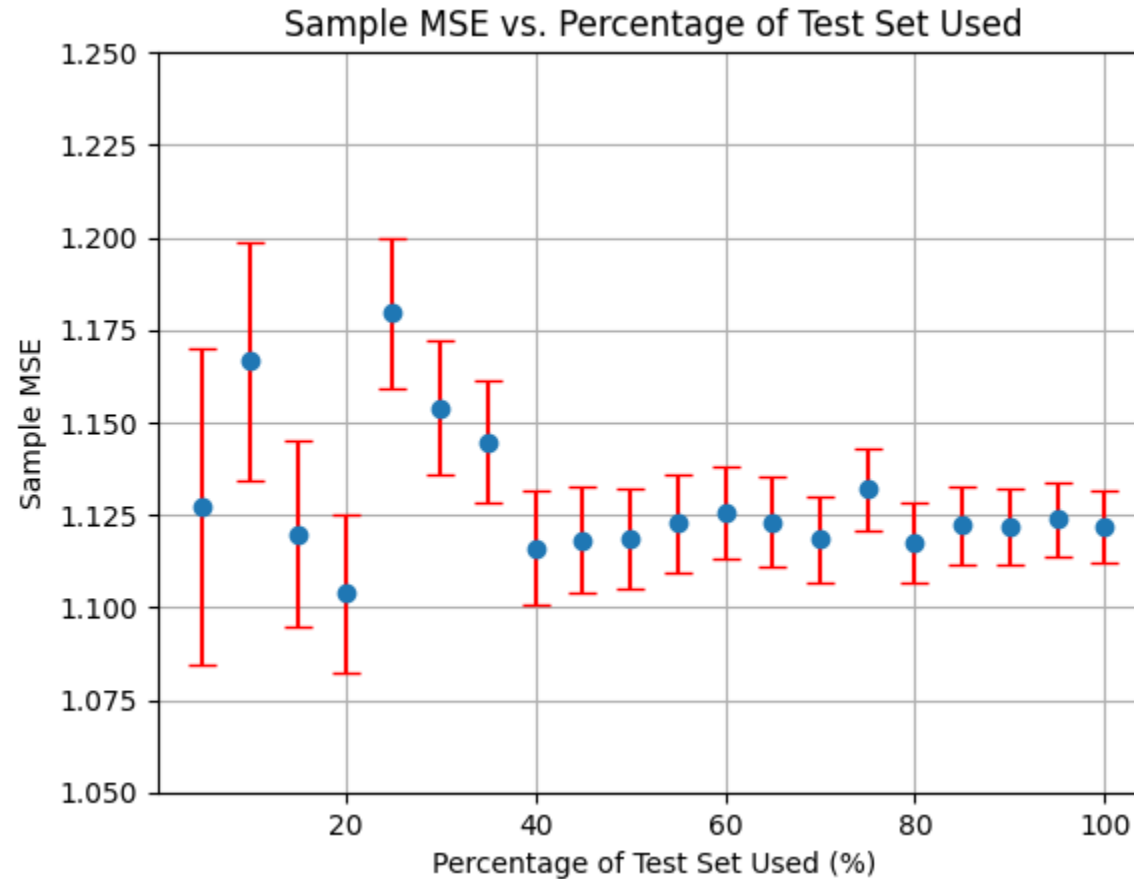
Plot: Sample MSE using different amounts of the test data, **with error bars (standard error)**.




Plot: Sample MSE using different amounts of the test data, **with error bars (standard error)**.

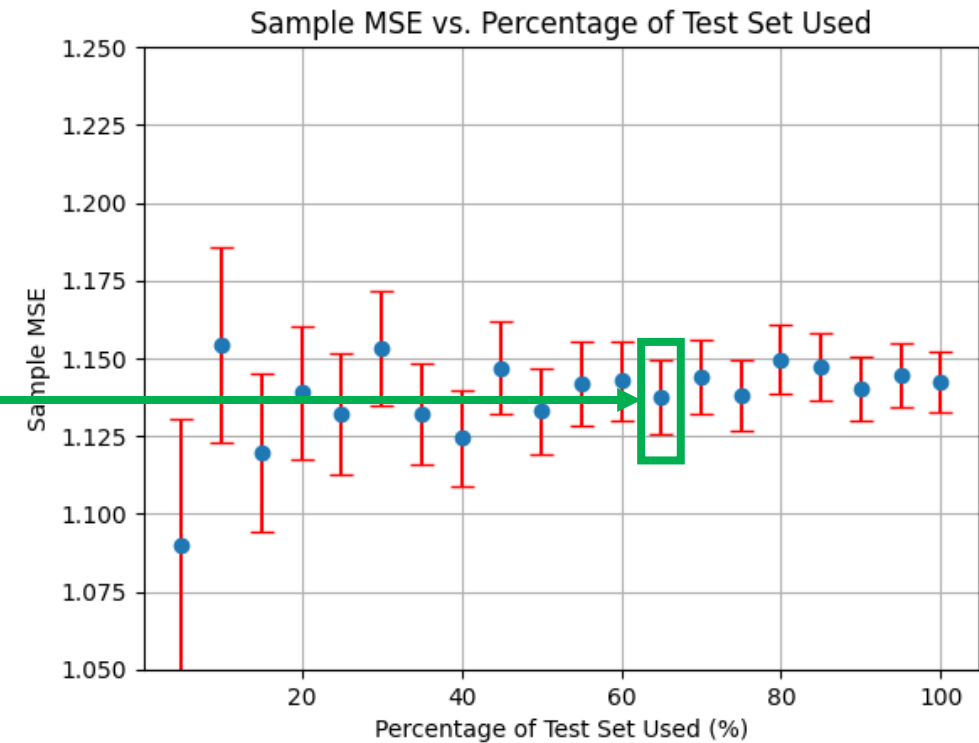


Plot: Sample MSE using different amounts of the test data, **with error bars (standard error)**.



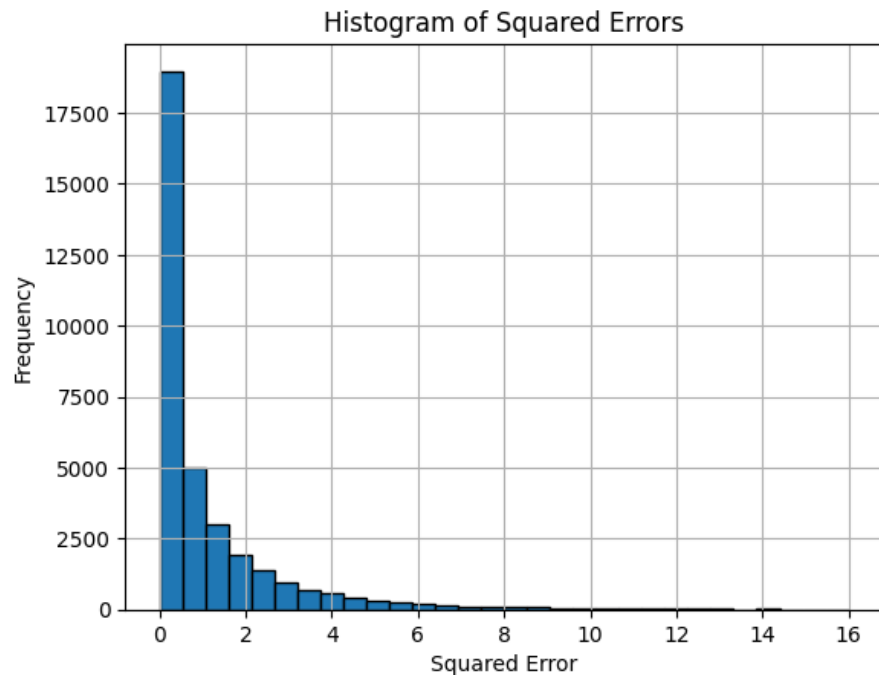
Error Bars: Standard Error

- Notice that given a single sample MSE (one point on this plot) we can compute the standard error.
- If we computed the sample MSE for our test set and the standard error and observed **this**  can we conclude that the actual MSE is within this interval with 95% confidence?

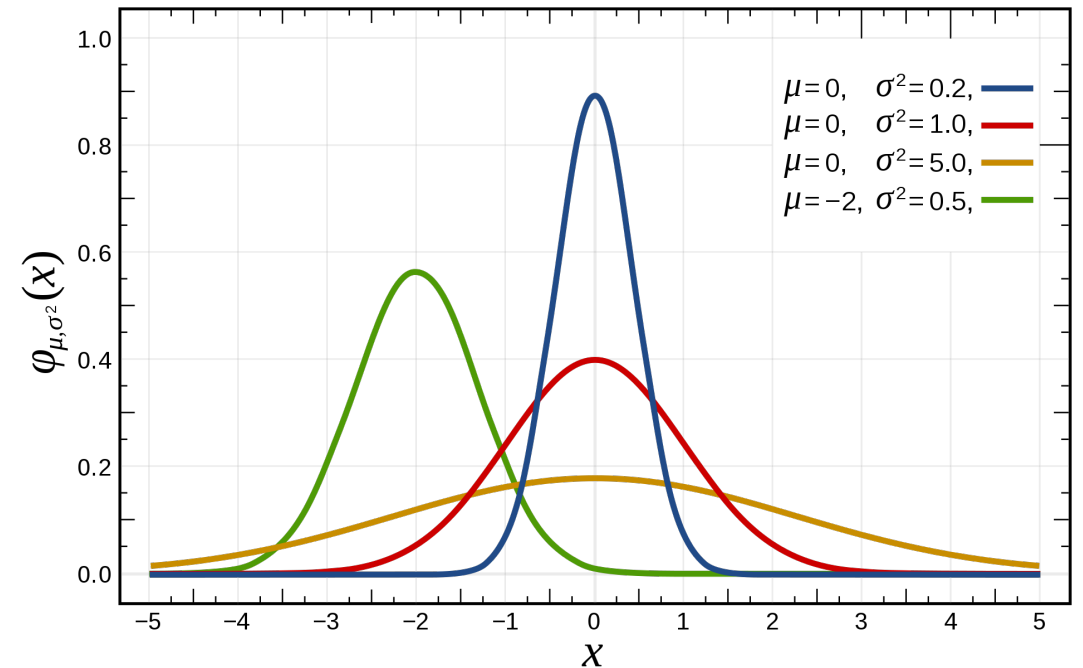


Error Bars: Standard Error

- Yes, if
 - The squared errors are normally distributed
 - The error bars were inflated by a factor of 1.96

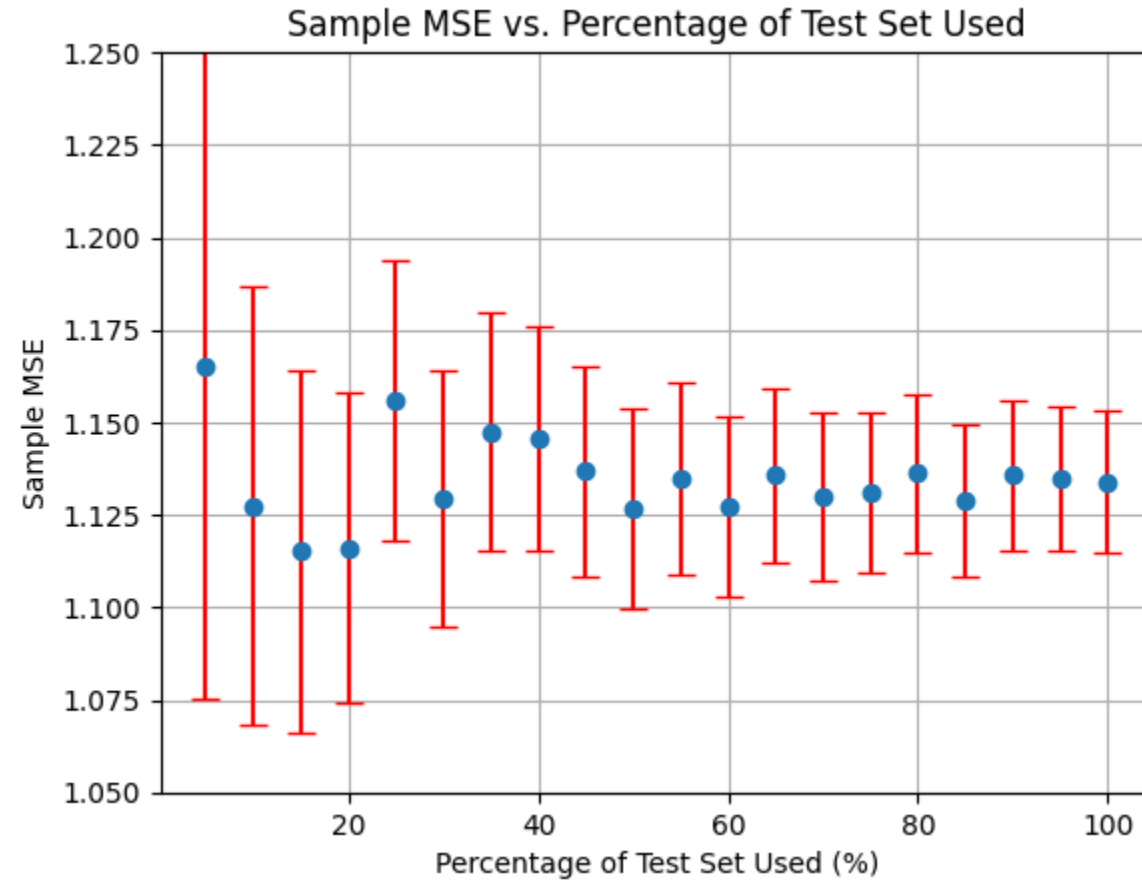


\neq

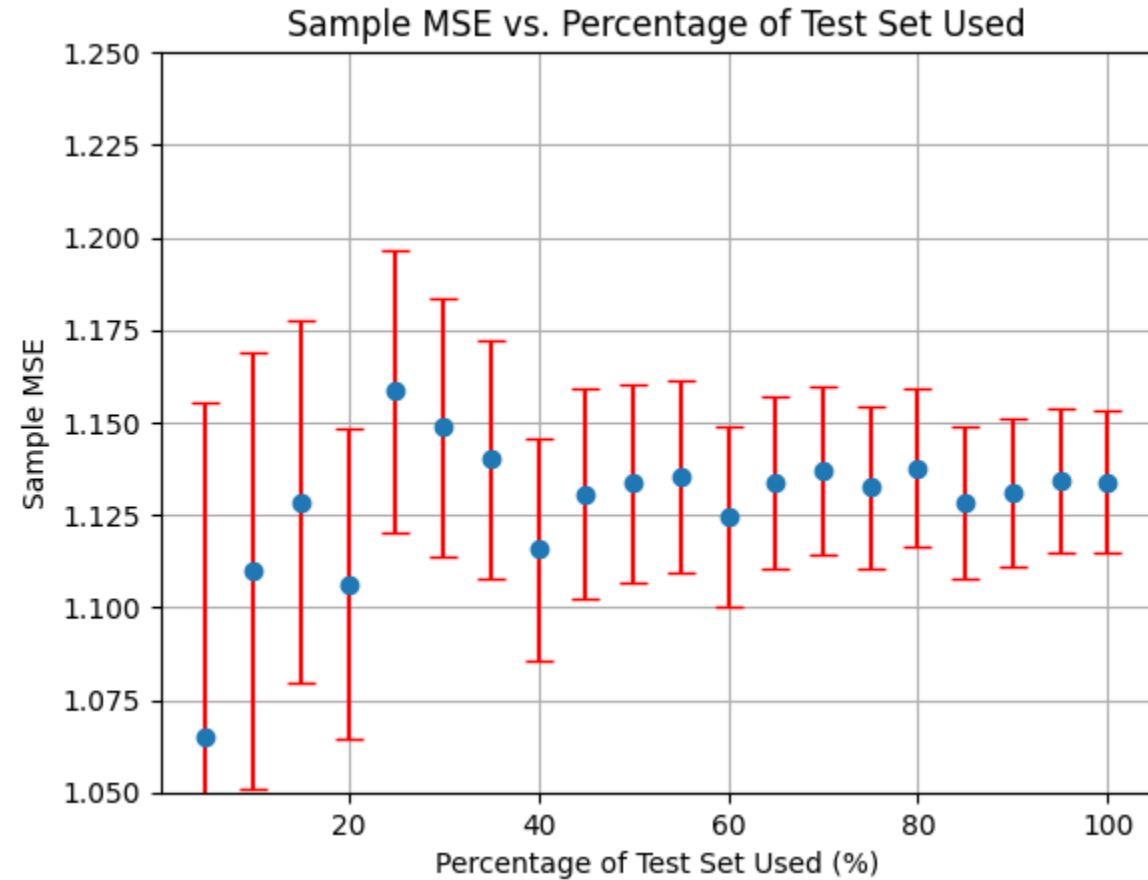


(the squared errors are far from normally distributed)

Error Bars: Standard Error $\times 1.96$



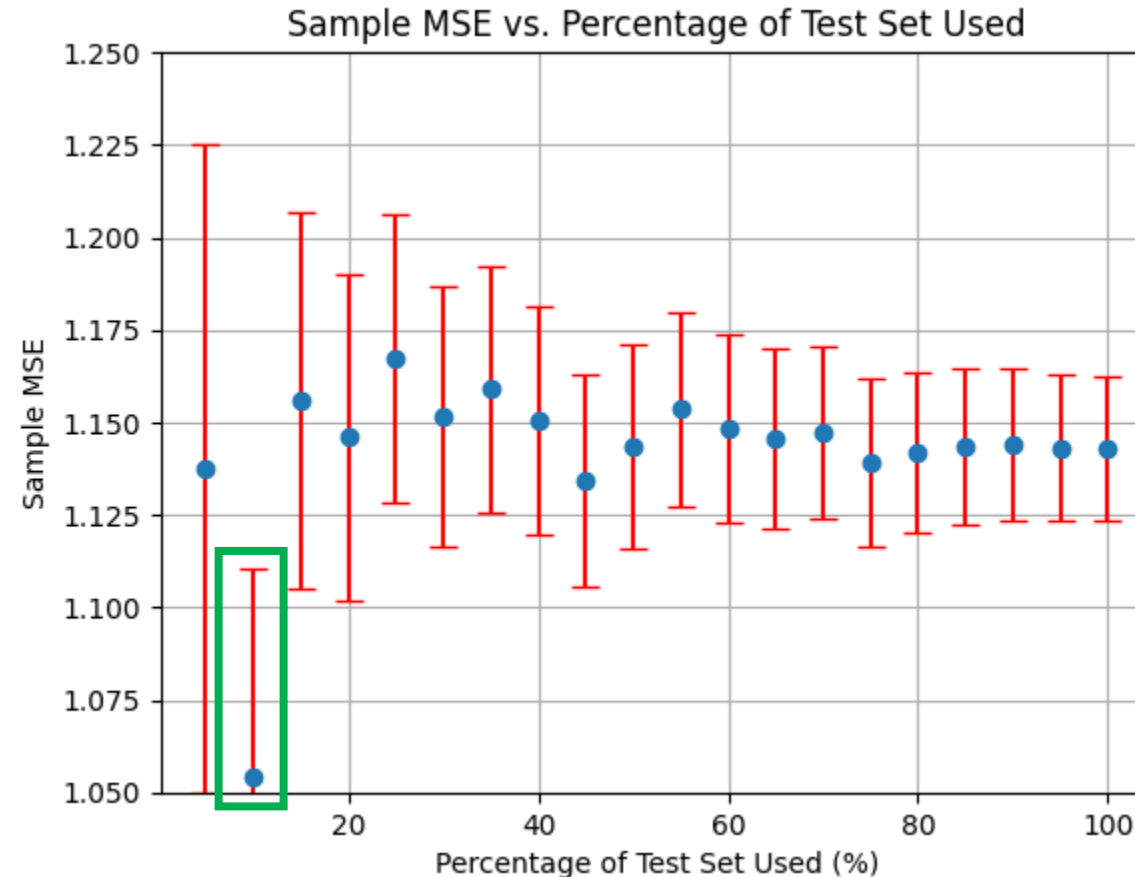
Error Bars: Standard Error x 1.96



Error Bars: Standard Error x 1.96

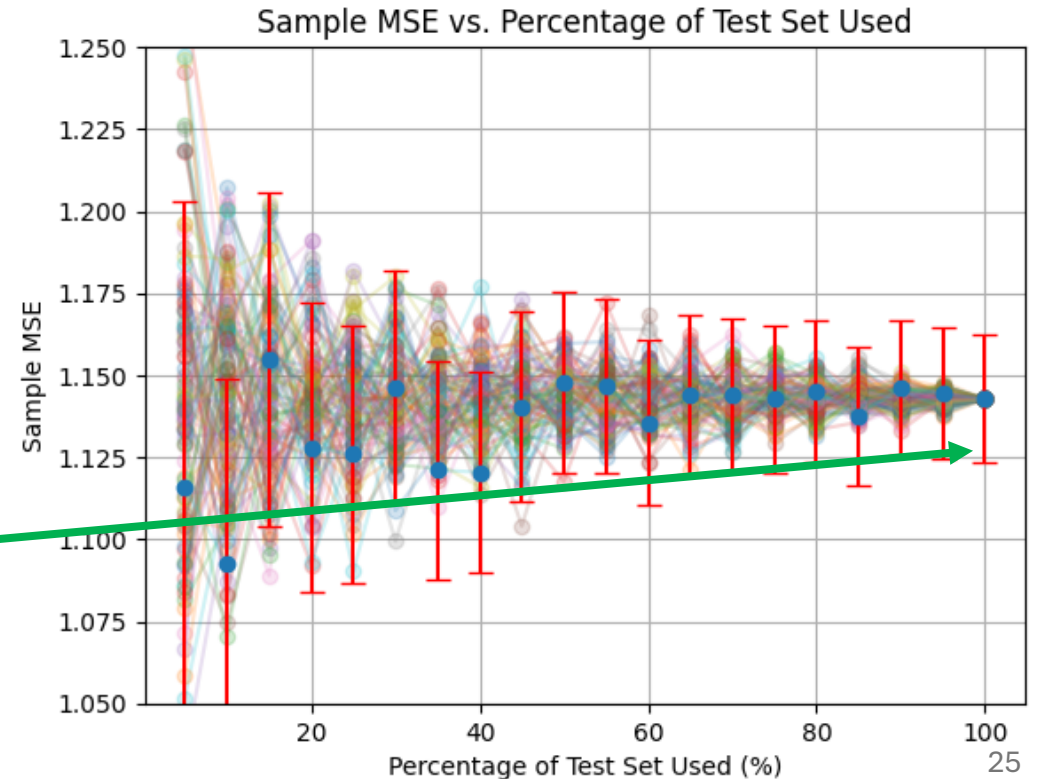
- Remember, each error bar has a 5% error chance.
- With 20 points and error bars, we expect *roughly* one (on average) to not contain the true MSE.

(If data was normally distributed)



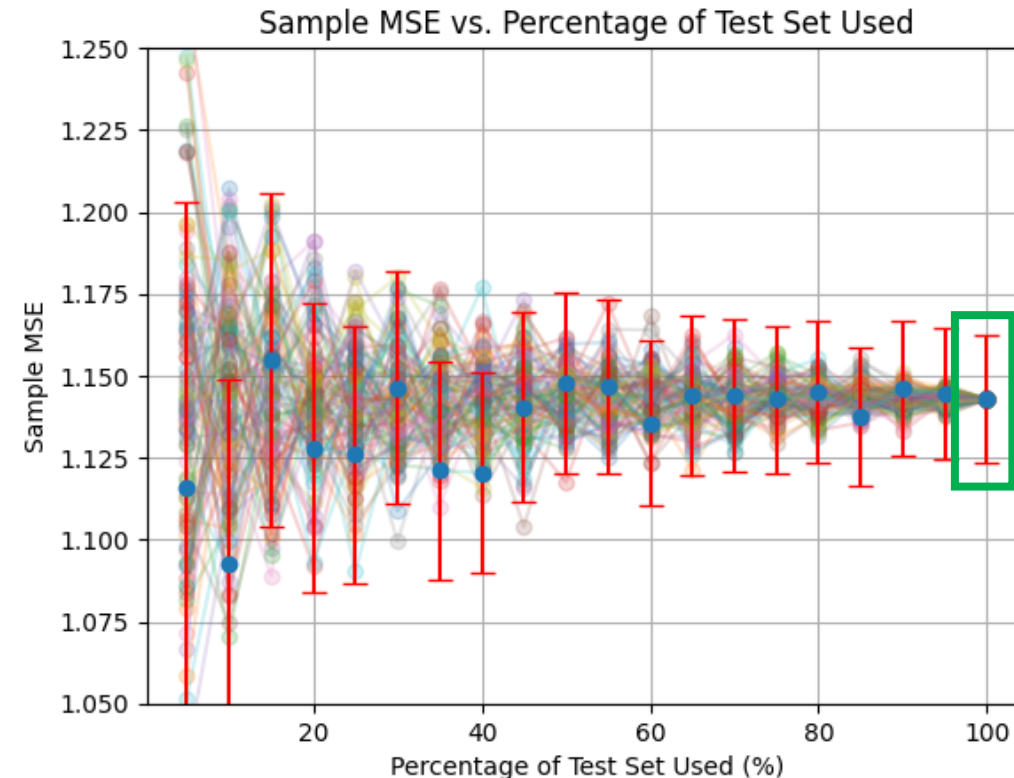
Error Bars: Standard Error x 1.96

- Although not entirely reliable (due to lack of normality and evaluating many estimates), error bars give an *idea* of how much results can be trusted.
- **Note:** ML texts/papers almost never inflate standard-error error-bars by a factor of 1.96.
- **Note:** The blue points and error bars **can** be computed from a single run
- **Question:** Why are the error bars still so wide when using 100% of the testing data?



Answer

- Recall that the squared errors are computed from a portion of a test set.
- At 100% of the test set, the plotted MSE is always the MSE on the full test set.
- This test set is still finite, and the sample MSE is still an approximation of the true MSE!
- This plot *underestimates* how much the estimates vary, particularly for large %.
- The red error bars do not have this bias.
 - They are a more accurate depiction of uncertainty.



Common Evaluation

- Often ML texts evaluate models by doing the following:
 - Partition the data into train/test.
 - Train the model on the training data.
 - Evaluate the model on the testing data.
 - Report a performance metric and a number representing the *uncertainty* in this performance metric.
 - Format: performance \pm uncertainty
 - The uncertainty value can be standard deviation, standard error (SE), or a confidence interval (e.g., $1.96 \times \text{SE}$).
 - Example (MSE of NN on GPA):

$$0.149729 \pm 0.0013.$$

	Model	MSE	RMSE	MAE	R ²
0	k-NN k=1 sigma=None	1.066188	1.032564	0.793455	-0.635682
1	k-NN k=100 sigma=None	0.556796	0.746187	0.587380	0.145797
2	k-NN k=110 sigma=90	0.555601	0.745386	0.586671	0.147631



$\pm 1.96 \times SE$

	Model	MSE	RMSE	MAE
0	k-NN k=1 sigma=None	1.104 \pm 0.075	1.051 \pm 0.029	0.803 \pm 0.029
1	k-NN k=100 sigma=None	0.565 \pm 0.041	0.752 \pm 0.020	0.586 \pm 0.020
2	k-NN k=110 sigma=90	0.565 \pm 0.041	0.752 \pm 0.020	0.586 \pm 0.020

We can be *somewhat* confident that the model learned by NN is worse than the model learned by k-NN ($k = 100$) and weighted k-NN ($k = 110, \sigma = 90$).

We are not confident about k-NN vs weighted k-NN.

Note: Always check for the *meaning* of the \pm value! Standard error, standard deviation, and confidence intervals all have very different meanings!

Does this table give us confidence that k -NN with $k = 100$ will outperform NN on the GPA data set?

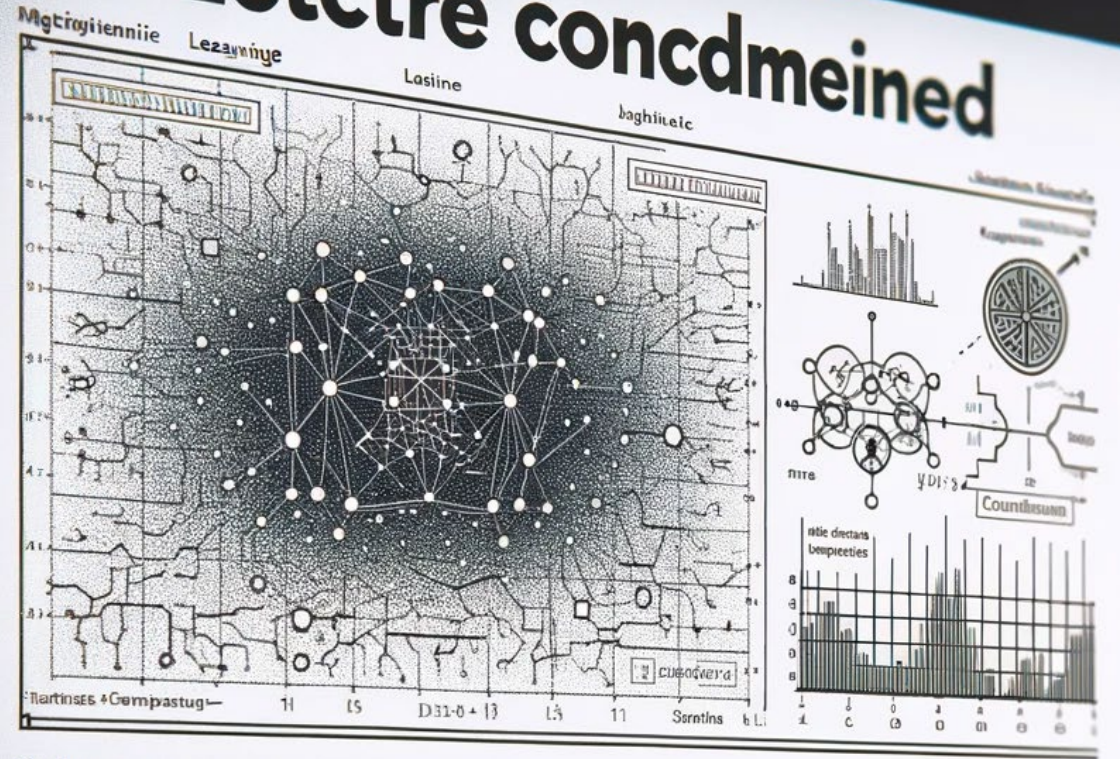
	Model	MSE	RMSE	MAE
0	k-NN k=1 sigma=None	1.104 ± 0.075	1.051 ± 0.029	0.803 ± 0.029
1	k-NN k=100 sigma=None	0.565 ± 0.041	0.752 ± 0.020	0.586 ± 0.020
2	k-NN k=110 sigma=90	0.565 ± 0.041	0.752 ± 0.020	0.586 ± 0.020

No!

- This evaluates a single model learned by each of these algorithms.
- Different training sets will result in different models.
 - Perhaps this training set was “lucky” for the latter two methods and “unlucky” for the first!

End

Letctre concdmeined



Dgoinnric



Mbcine Learning

Thank you.

