

Morse Code Datasets for Machine Learning

Sourya Dey, Keith Chugg, Peter Beerel

9th International Conference on Computing,
Communication and Networking Technologies

July 2018



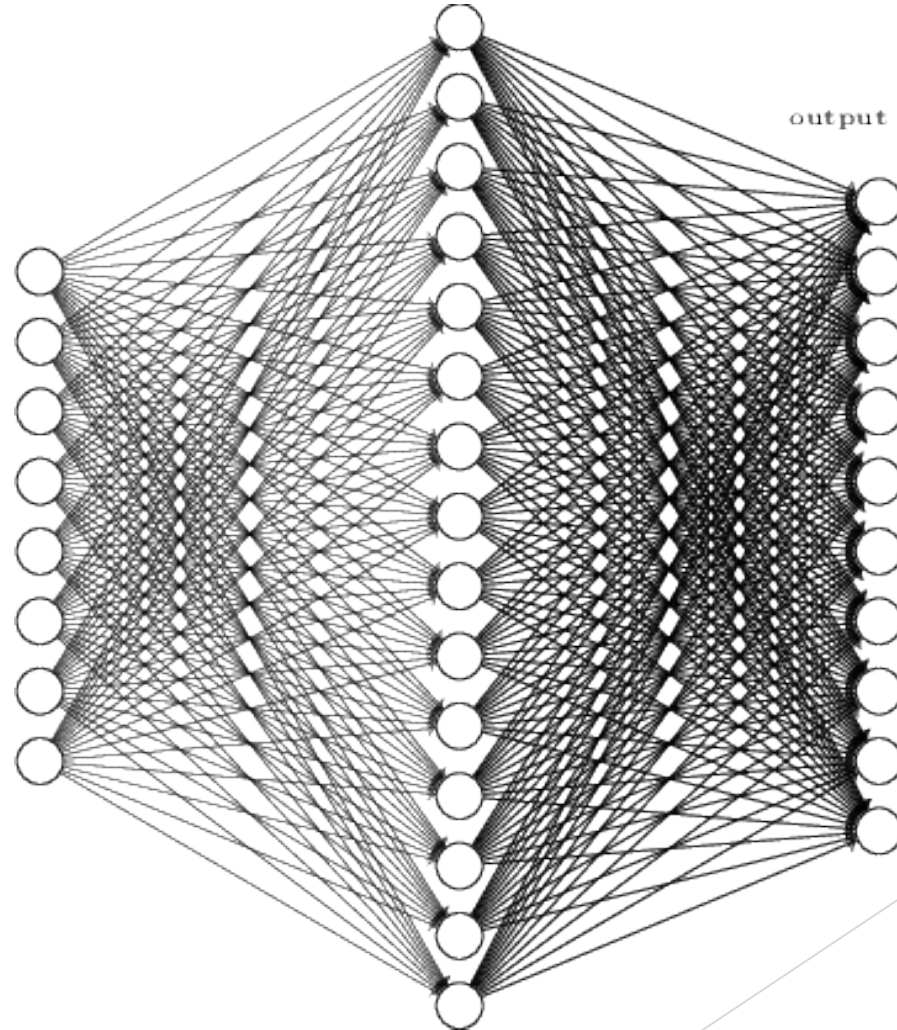
USC University of
Southern California

Machine Learning and Neural Networks

*An algorithm to
learn from data
and classify it*

Machine Learning and Neural Networks

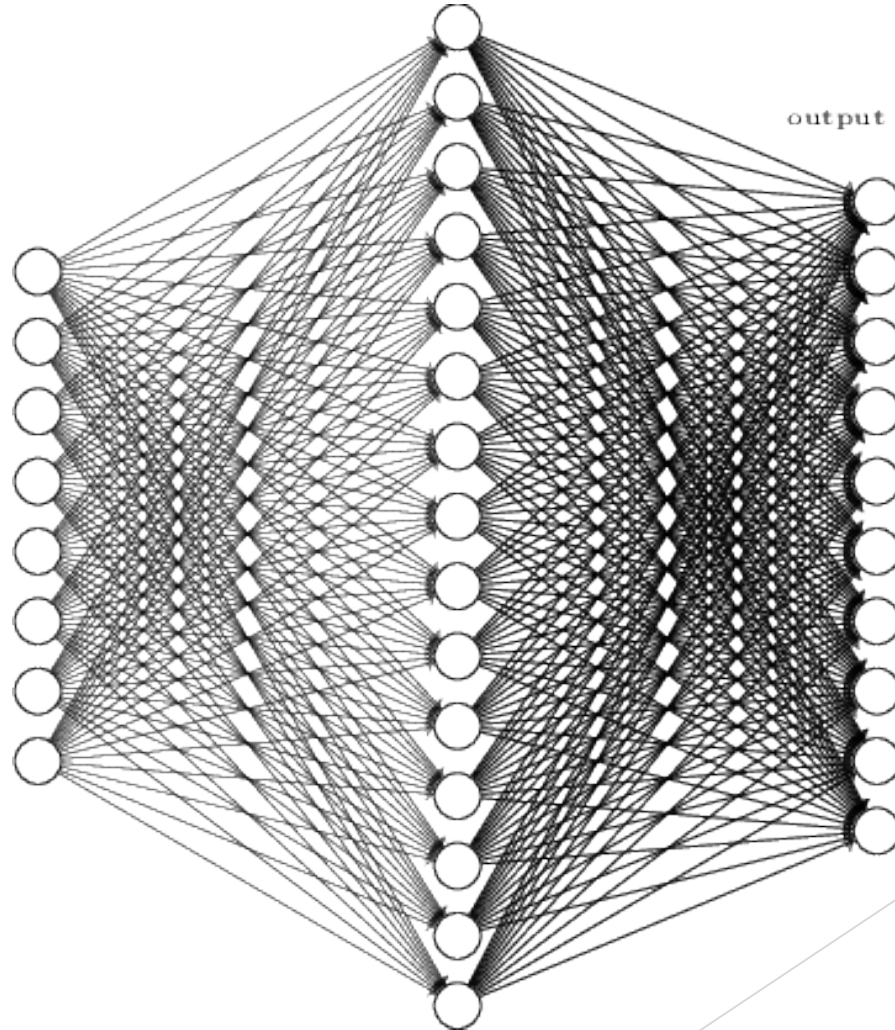
*An algorithm to
learn from data
and classify it*



Machine Learning and Neural Networks

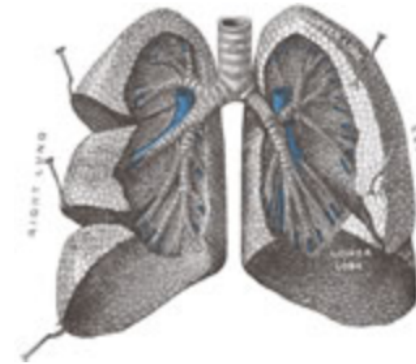
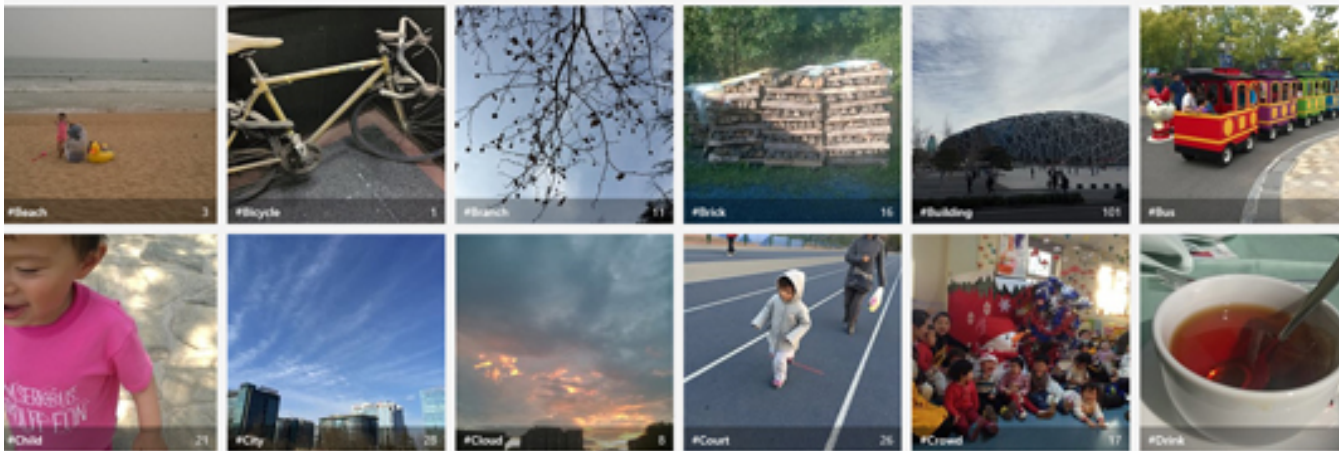
*An algorithm to
learn from data
and classify it*

*Need a lot of
data for good
performance*



Issues with Natural Data

- ▶ Most data is naturally collected and labeled by humans
- ▶ Labeling is **time-consuming** (e.g. Imagenet¹)
- ▶ Data can have **missing features** (e.g. Lung cancer dataset²)



Synthetic data as a Solution

- ▶ **Synthetic data** generated and labeled using algorithms
- ▶ Can be mass-produced cheaply without missing features
- ▶ Algorithm can be tuned to:
 - ▶ *Adjust difficulty*
 - ▶ Get any distribution

Overview of our Work

- ▶ Algorithm to generate Morse code classification datasets of varying difficulty
- ▶ Metrics to evaluate difficulty of a dataset

Overview of our Work

- ▶ Algorithm to generate Morse code classification datasets of varying difficulty
- ▶ Metrics to evaluate difficulty of a dataset

Morse code is a system of communication to encode characters as dots and dashes

+ • — • — •

Overview of our Work

- ▶ Algorithm to generate Morse code classification datasets of varying difficulty
- ▶ Metrics to evaluate difficulty of a dataset

Morse code is a system of communication to encode characters as dots and dashes

+

• — • — •

64 character
classes

The Algorithm

Step 1:

Frame length: 64

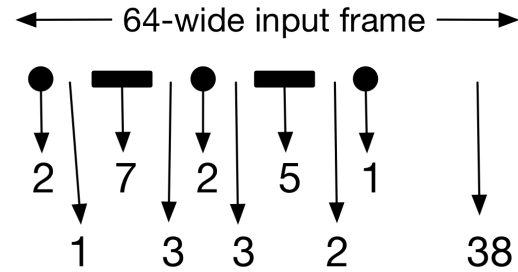
Dot: 1-3

Dash: 4-9

Intermediate space: 1-3

Leading spaces: None

Trailing spaces: Remaining at end



Codeword Length = 26. Remaining spaces = 38

The Algorithm

Step 1:

Frame length: 64

Dot: 1-3

Dash: 4-9

Intermediate space: 1-3

Leading spaces: None

Trailing spaces: Remaining at end



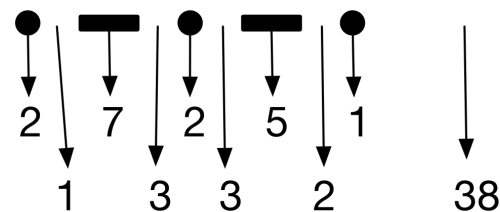
Step 2:

Expected value range = [0, 16]

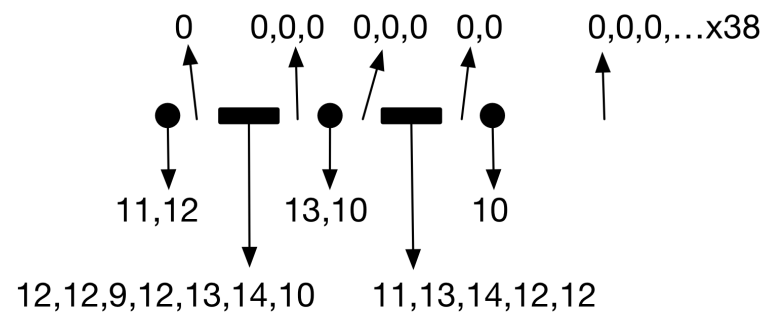
Dot, dash = $\text{Normal}(12, 4/3)$

Space = 0

← 64-wide input frame →



Codeword Length = 26. Remaining spaces = 38



The Algorithm

Step 1:

Frame length: 64

Dot: 1-3

Dash: 4-9

Intermediate space: 1-3

Leading spaces: None

Trailing spaces: Remaining at end

Step 2:

Expected value range = $[0, 16]$

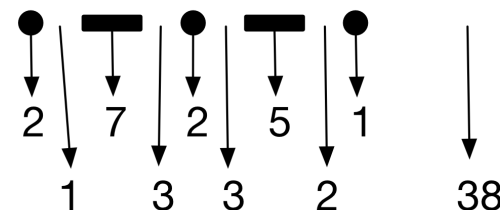
Dot, dash = $\text{Normal}(12, 4/3)$

Space = 0

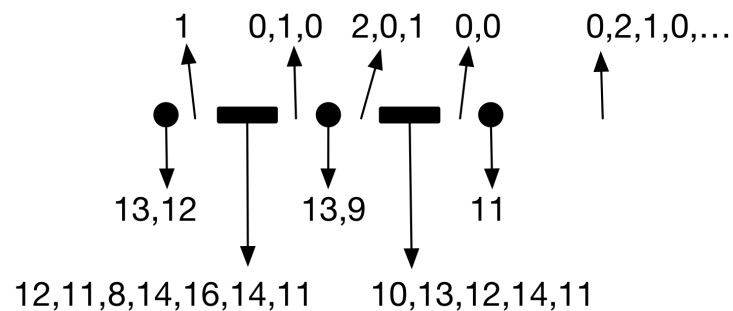
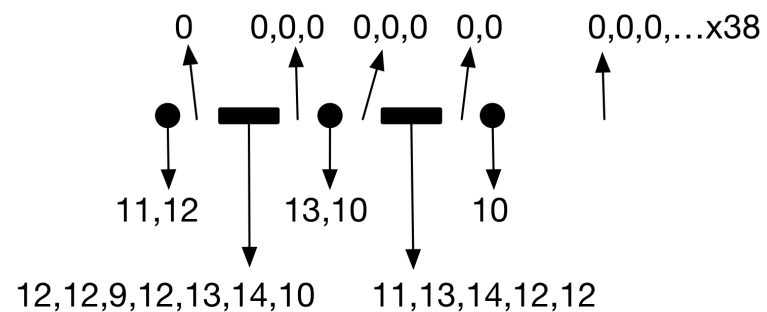
Step 3:

Additive Noise = $\text{Normal}(0, \sigma)$
(For this case, $\sigma=1$)

← 64-wide input frame →

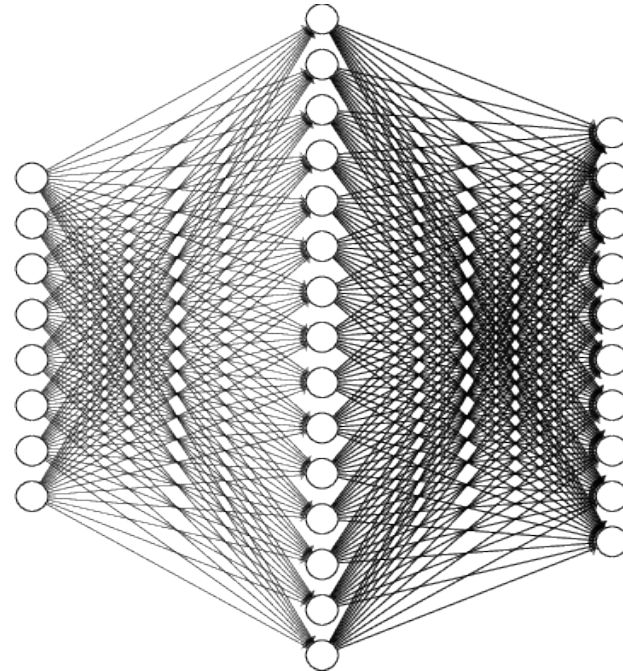


Codeword Length = 26. Remaining spaces = 38



The Neural Network

64 input neurons =
Frame length of each
Morse codeword

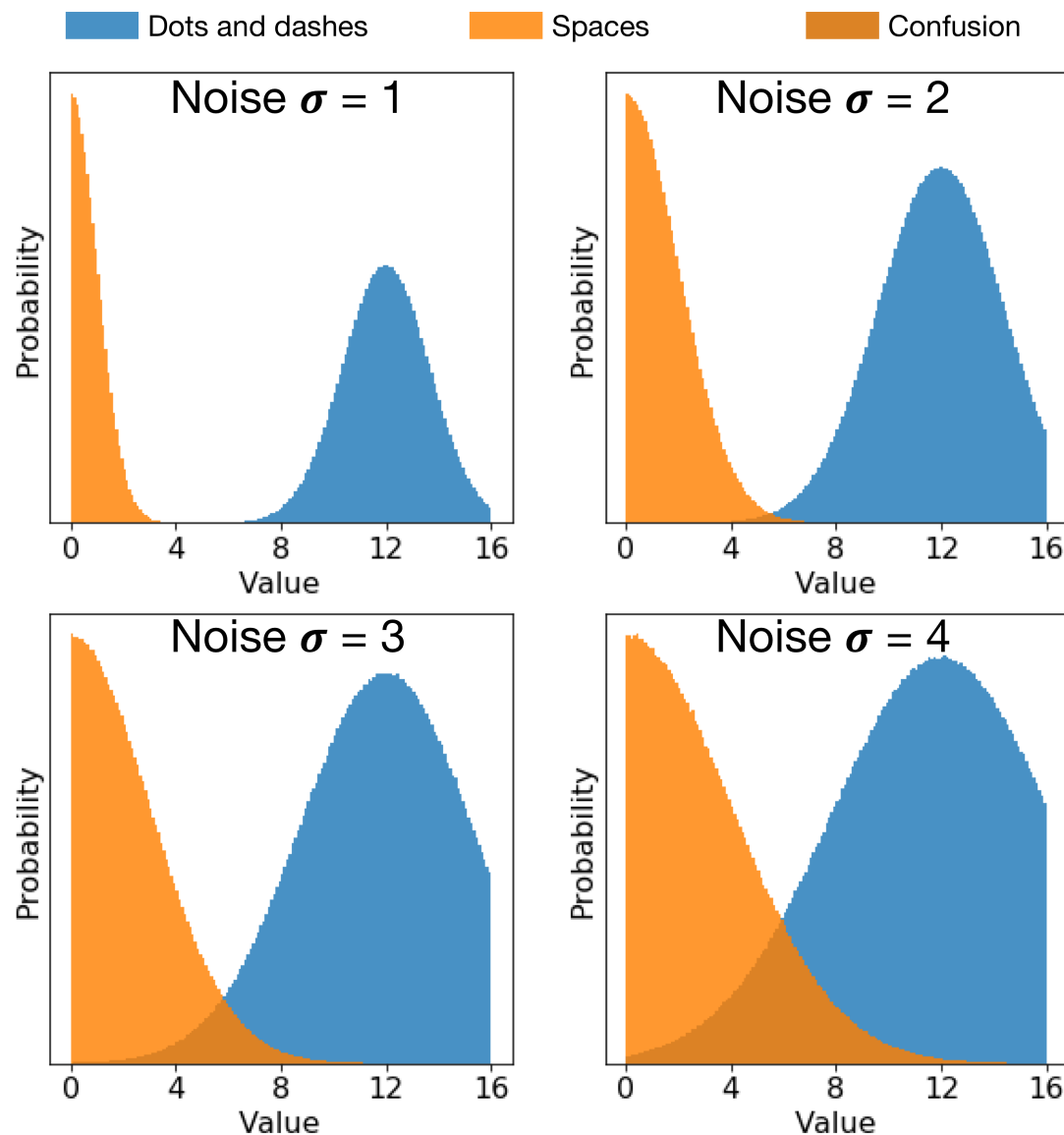


1024 hidden neurons

64 output neurons =
Number of character
classes

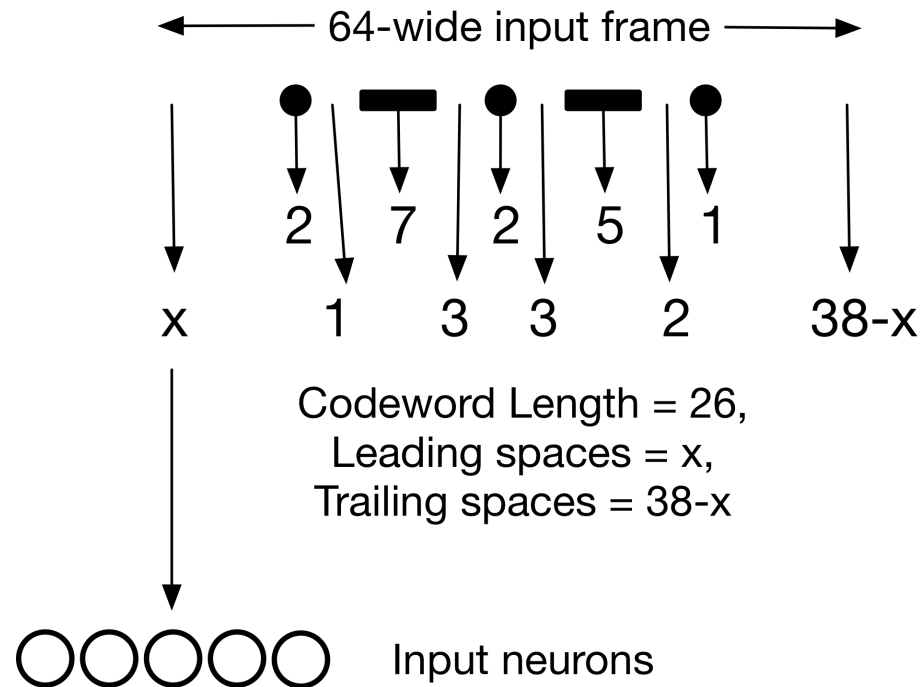
Variations and Difficulty Scaling - 1

Increasing σ of noise leads to confusion between dots, dashes and spaces



Variations and Difficulty Scaling - 2

Distribute remaining spaces randomly between leading and trailing



Variations and Difficulty Scaling - 3, 4

Dash length is 3-9, can be confused with dots and spaces

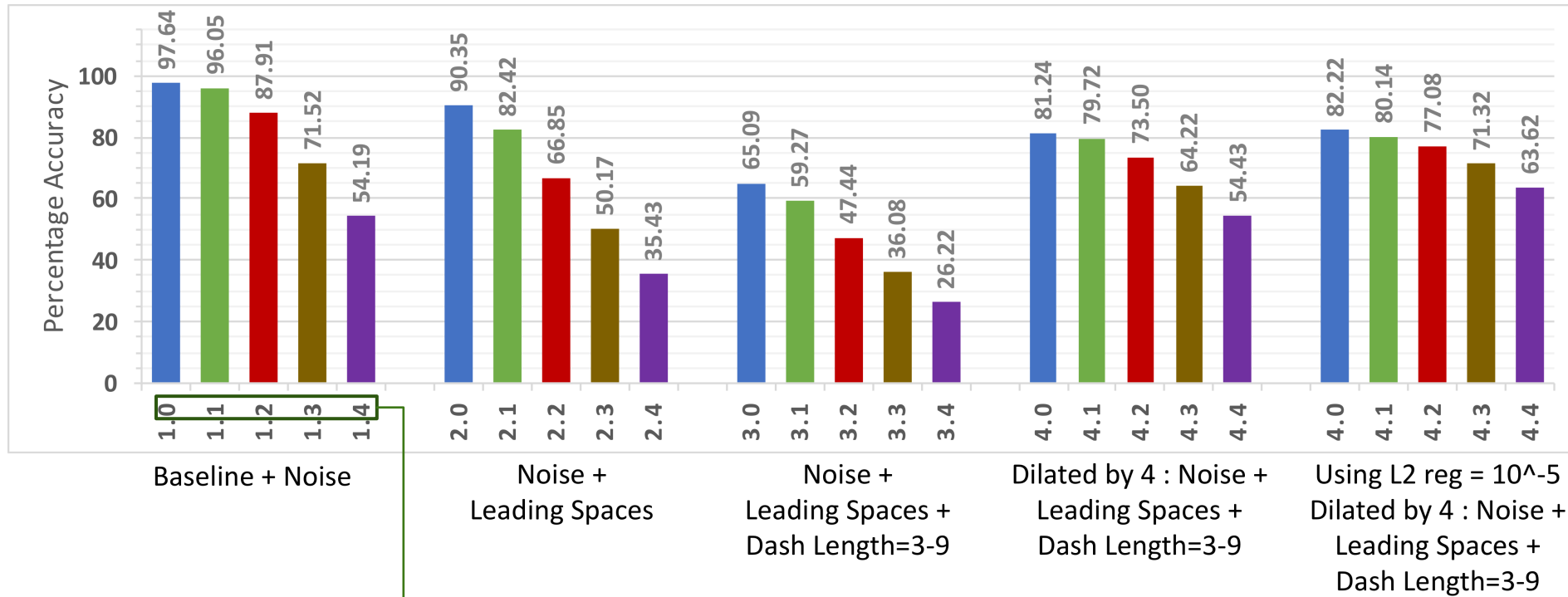
Variations and Difficulty Scaling - 3, 4

Dash length is 3-9, can be confused with dots and spaces

Dilate inputs by 4x

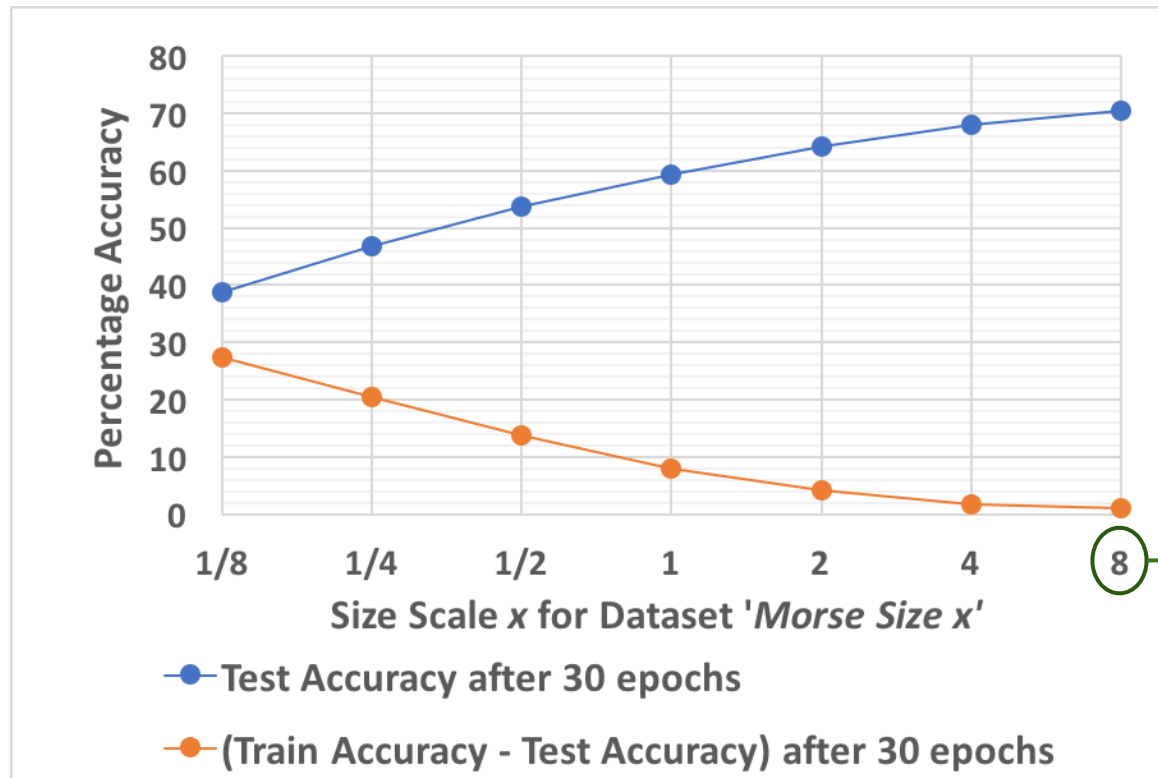
Property	Before Dilation	After Dilation
Frame length (= Number of inputs)	64	256
Space	1-3	4-12
Dot	1-3	4-12
Dash	3-9	12-36

Classification Accuracy on Test Data



Increasing Dataset Size

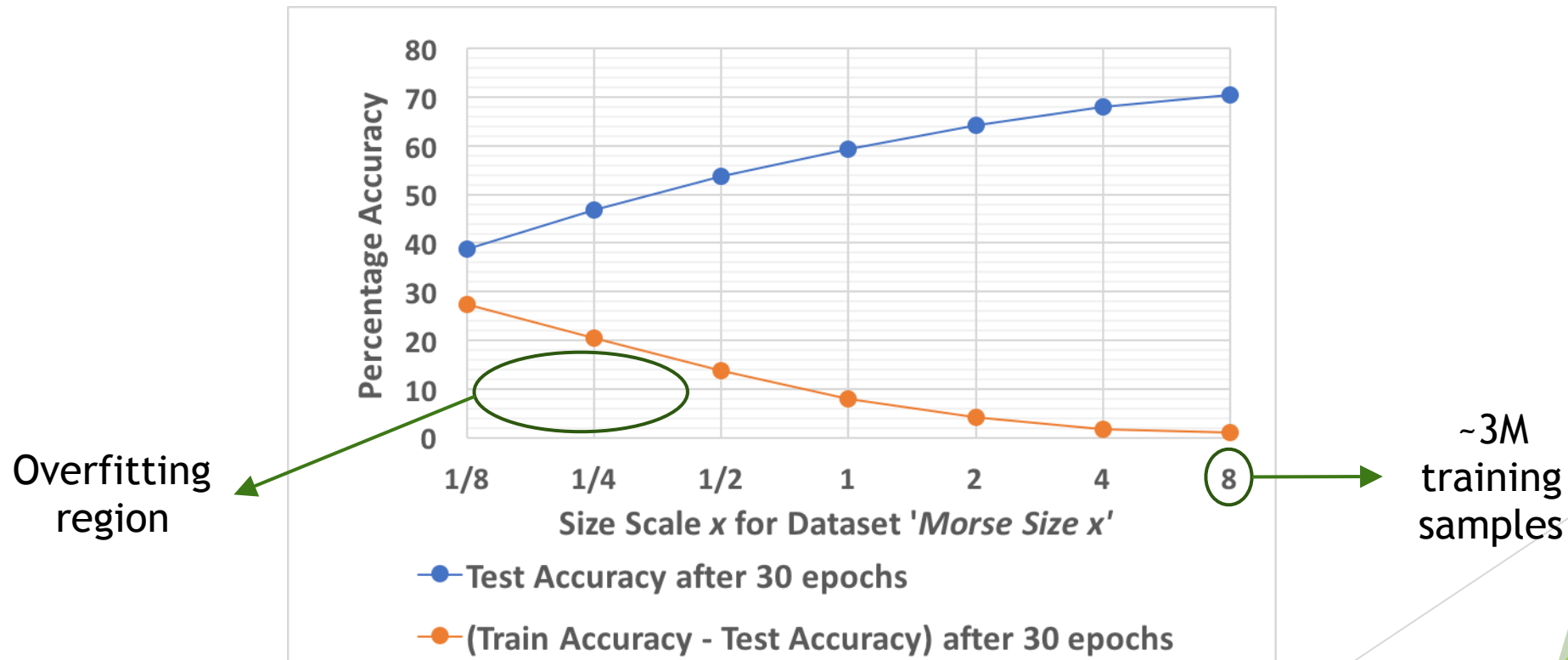
Unlimited amounts of data can be easily generated using computer algorithms



~3M
training
samples

Increasing Dataset Size

Unlimited amounts of data can be easily generated using computer algorithms



Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors

Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors

$$\begin{aligned} \sum_{m=1}^M P(m) \left[\max_{\substack{j \in \{1, 2, \dots, M\} \\ j \neq m}} P_{PW}(j|m) \right] &\leq P(E) \\ &\leq \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M P_{PW}(j|m) \end{aligned}$$

Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors

$$L = \sum_{m=1}^M P(m) Q \left(\sqrt{\frac{d_{\min}(m)^2}{4\sigma_m^2}} \right) \leftarrow \sum_{m=1}^M P(m) \left[\max_{\substack{j \in \{1,2,\dots,M\} \\ j \neq m}} P_{PW}(j|m) \right] \leq P(E)$$
$$\leq \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M P_{PW}(j|m)$$

Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors

$$\begin{aligned}
 L = \sum_{m=1}^M P(m) Q \left(\sqrt{\frac{d_{\min}(m)^2}{4\sigma_m^2}} \right) &\leftarrow \sum_{m=1}^M P(m) \left[\max_{\substack{j \in \{1, 2, \dots, M\} \\ j \neq m}} P_{PW}(j|m) \right] \leq P(E) \\
 &\leq \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M P_{PW}(j|m) \\
 &\downarrow \\
 U = \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M Q \left(\sqrt{\frac{d(m, j)^2}{4\sigma_m^2}} \right)
 \end{aligned}$$

Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors

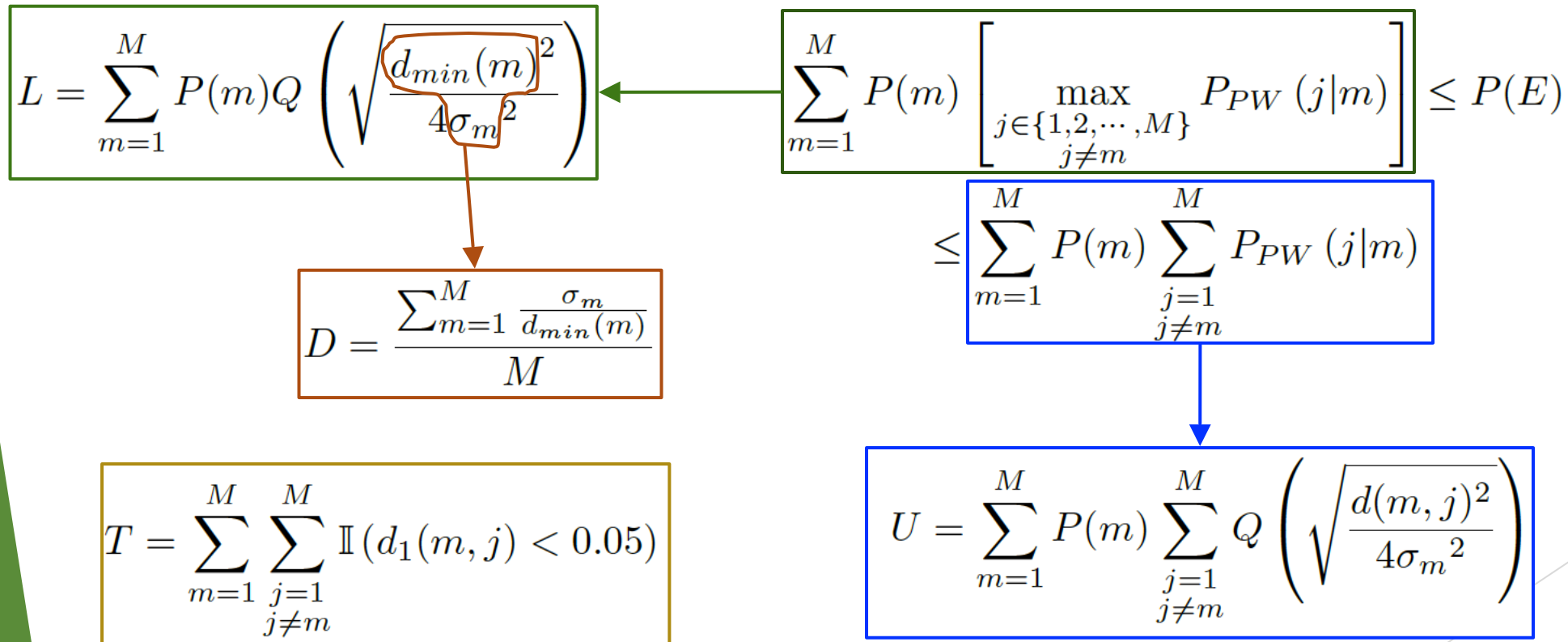
$$\begin{aligned}
 L &= \sum_{m=1}^M P(m) Q \left(\sqrt{\frac{d_{\min}(m)^2}{4\sigma_m^2}} \right) \\
 &\leftarrow \sum_{m=1}^M P(m) \left[\max_{\substack{j \in \{1, 2, \dots, M\} \\ j \neq m}} P_{PW}(j|m) \right] \leq P(E) \\
 &\leq \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M P_{PW}(j|m) \\
 &\downarrow \\
 U &= \sum_{m=1}^M P(m) \sum_{\substack{j=1 \\ j \neq m}}^M Q \left(\sqrt{\frac{d(m, j)^2}{4\sigma_m^2}} \right)
 \end{aligned}$$

Diagram illustrating the relationship between classification error probability and dataset metrics:

- The top-left box shows the loss function L as a function of the minimum distance $d_{\min}(m)$ and standard deviation σ_m . A red box highlights $d_{\min}(m)$ and σ_m , with an arrow pointing to the bottom-left box.
- The bottom-left box shows the average distance metric $D = \frac{\sum_{m=1}^M \frac{\sigma_m}{d_{\min}(m)}}{M}$.
- The top-right box shows the probability of error $P(E)$ as a function of the maximum pairwise probability $P_{PW}(j|m)$.
- The middle-right box shows the pairwise probability $P_{PW}(j|m)$ as a function of the distance $d(m, j)$ and standard deviation σ_m .
- The bottom-right box shows the loss function U as a function of the distance $d(m, j)$ and standard deviation σ_m .

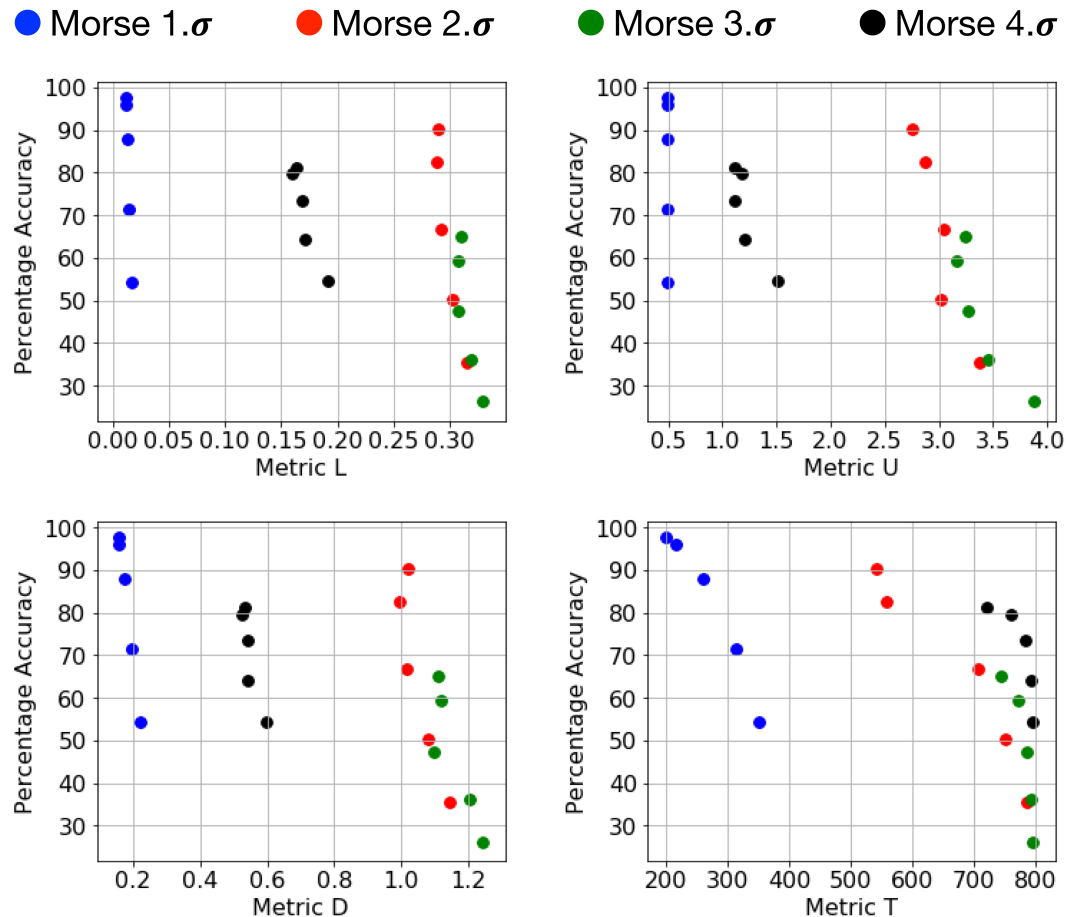
Dataset Evaluating Metrics

Difficult datasets have increased probability of classification errors



Performance of the Metrics

Harder datasets have lower accuracy and higher metric values



Metric	$-\rho$
L	0.59
U	0.64
D	0.63
T	0.64

Conclusion

- ▶ Algorithm to generate machine learning datasets of tunable difficulty
- ▶ Synthetic data to solve challenges associated with natural data
- ▶ Metrics to evaluate dataset difficulty prior to training

Thank you!

Questions?