**Ayesha Ashraf**

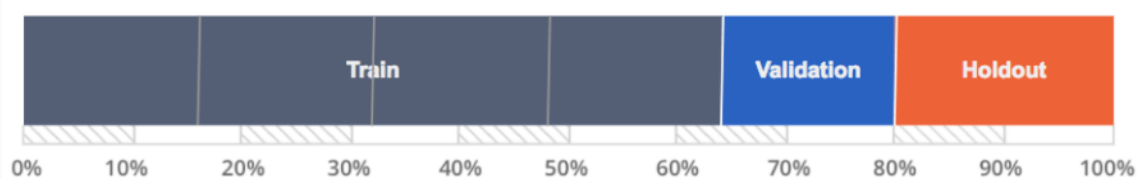**Task 23 (CheatSheet)**

**Training Sets, Validation Sets, and Holdout Sets**

**What are Training, Validation, and Holdout?**

**Partitioning Data**

The first step in developing a machine learning model is training and validation. In order to train and validate a model, you must first **partition** your dataset, which involves choosing what percentage of your data to use for the training, validation, and holdout sets. The following example shows a dataset with 64% training data, 16% validation data, and 20% holdout data.



**What is a Training Set?**

A **training** set is the subsection of a dataset from which the machine learning algorithm uncovers, or "learns," relationships between the features and the target variable. In supervised machine learning, training data is labeled with known outcomes.

**What is a Validation Set?**

A **validation** set is another subset of the input data to which we apply the machine learning algorithm to see how accurately it identifies relationships between the known outcomes for the target variable and the dataset's other features.

## What is a Holdout Set?

Sometimes referred to as "testing" data, a **holdout** subset provides a final estimate of the machine learning model's performance after it has been trained and validated. Holdout sets should never be used to make decisions about which algorithms to use or for improving or <u>tuning</u> algorithms.

## Why are Training, Validation, and Holdout Sets Important?

Partitioning data into training, validation, and holdout sets allows you to develop highly accurate models that are relevant to data that you collect in the future, not just the data the model was trained on. By training your data, validating it, and testing it on the holdout set, you get a real sense of how accurate the model's outcomes will be, leading to better decisions and greater confidence in your model's accuracy.

## Train vs. Validation vs. Test set

For training and testing purposes of our model, we should have our data broken down into three distinct dataset splits.

## The Training Set

It is the set of data that is used to train and make the model learn the hidden features/patterns in the data.

In each epoch, the same training data is fed to the neural network architecture repeatedly, and the model continues to learn the features of the data.

The training set should have a diversified set of inputs so that the model is trained in all scenarios and can predict any unseen data sample that may appear in the future.

**The Validation Set**

The validation set is a set of data, separate from the training set, that is used to validate our model performance during training.

This validation process gives information that helps us tune the model's hyperparameters and configurations accordingly. It is like a critic telling us whether the training is moving in the right direction or not.

The model is trained on the training set, and, simultaneously, the model evaluation is performed on the validation set after every epoch.

The main idea of splitting the dataset into a validation set is to prevent our model from overfitting i.e., the model becomes really good at classifying the samples in the training set but cannot generalize and make accurate classifications on the data it has not seen before.
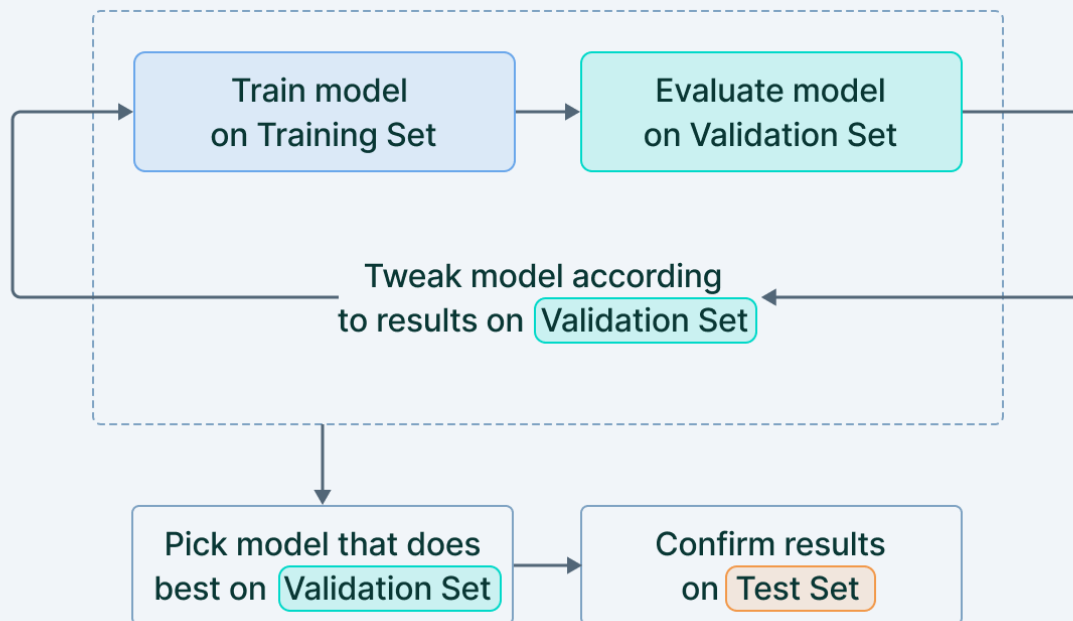
**The Test Set**

The test set is a separate set of data used to test the model *after* completing the training.

It provides an unbiased final model performance metric in terms of accuracy, precision, etc. To put it simply, it answers the question of "*How well does the model perform?*"

**Training data/validation/test**

How to split your Machine Learning data?

The creation of different samples and splits in the dataset helps us judge the true model performance.

The dataset split ratio depends on the number of samples present in the dataset and the model.

Some common inferences that can be derived on dataset split include:

- If there are several hyperparameters to tune, the machine learning model requires a larger validation set to optimize the model performance. Similarly, if the model has

fewer or no hyperparameters, it would be easy to validate the model using a small set of data.

- If a model use case is such that a false prediction can drastically hamper the model performance—like falsely predicting cancer—it's better to validate the model after each epoch to make the model learn varied scenarios.

- With the increase in the dimension/features of the data, the hyperparameters of the neural network functions also increase making the model more complex. In these scenarios, a large split of data should be kept in training set with a validation set.

The truth is—

There is no optimal split percentage.

One has to come to a split percentage that suits the requirements and meets the model's needs.

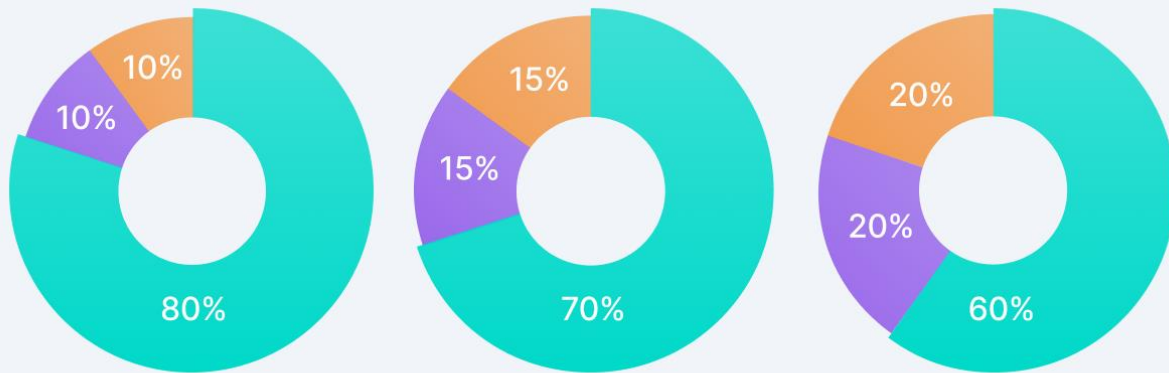However, there are two major concerns while deciding on the optimum split:

- If there is less training data, the machine learning model will show high variance in training.

- With less testing data/validation data, your model evaluation/model performance statistic will have greater variance.

Essentially, you need to come up with an optimum split that suits the need of the dataset/model.

But here's the rough standard split that you might encounter.

**Data Training Needs**

- Training data
- Validation data
- Test data

Chart 1: 80% Training data, 10% Test data, 10% Validation data

Chart 2: 70% Training data, 15% Test data, 15% Validation data

Chart 3: 60% Training data, 20% Test data, 20% Validation data

**Advanced techniques for data splitting**

Various data splitting techniques have been implemented in the Computer Vision literature to ensure a robust and fair way of testing machine learning models. Some of the most popular ones are explained below.

*Random*

Random sampling is the oldest and most popular method for dividing a dataset. As the name suggests, the dataset is shuffled, and samples are picked randomly and put in the train, validation, or the test set based on what percentage ratio is given by the user.

However, this method has a significant drawback. Random sampling works optimally on class-balanced datasets, i.e., datasets with the more or less the same number of samples in every dataset category. In the case of class-imbalanced datasets, such a data splitting method may create a bias.
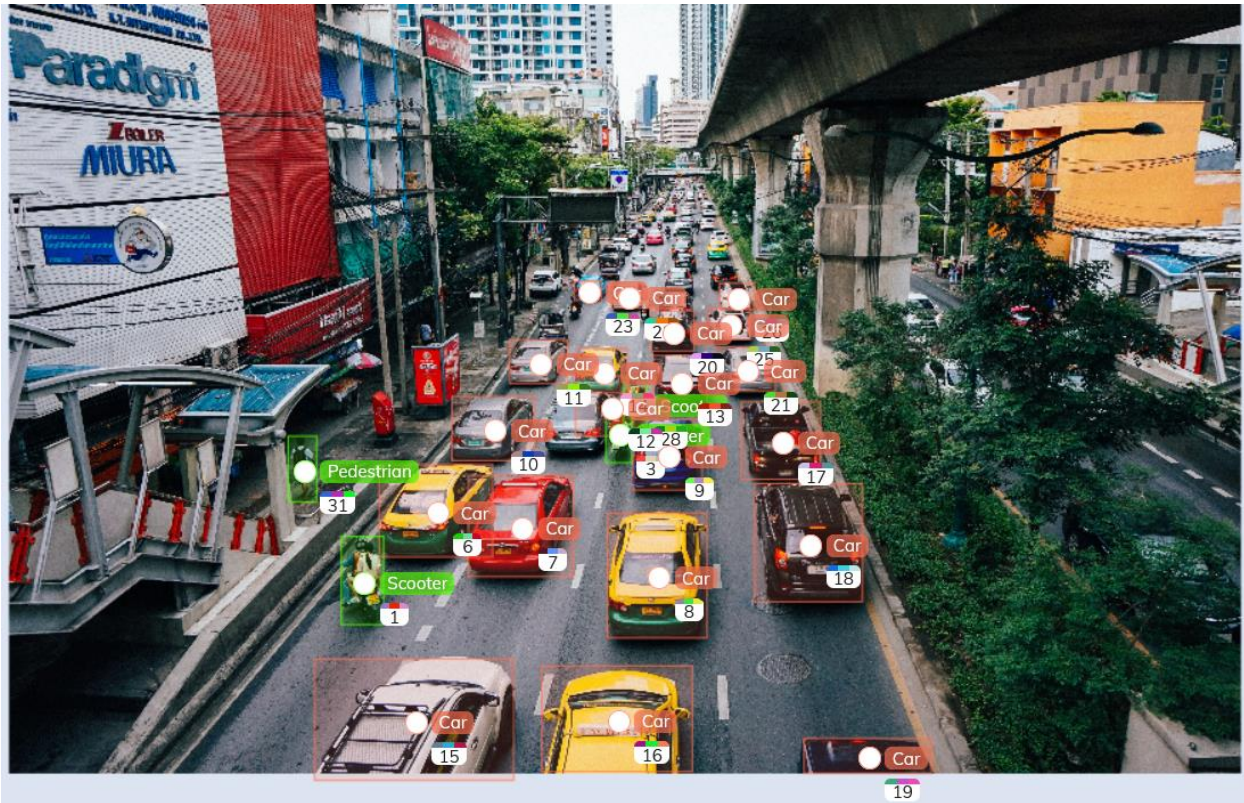
For example, if a dataset has 1000 images, of which 800 belong to the category: "hen" and 200 belong to the category: "cat," random sampling is employed to split the data into training and validation sets in an 80%-20% ratio (respectively), it might so happen that the training set consists only hen images, while the validation set consists of only cat images. Even if it's not so extreme, there will still be an imbalance distribution among the train and validation sets.

### Stratified

Stratified sampling for splitting a dataset alleviates the problem of Random Sampling in datasets with an imbalanced-class distribution. Here, the distribution of classes in each of the train, validation, and test sets is preserved.

Suppose a dataset consists of 1000 images, of which 600 are hen images and 400 are cat images. In that case, stratified sampling ensures that 60% of the images are of category "hen" and 40% are of category "cat" in the training and validation sets. That is, if a train-validation split of 80%-20% is desired, out of the 800 images in the training set, 480 images (60%) will be of hens, and the rest 320 (40%) will be of cats.

Let us consider another example. In object detection tasks, single images may contain several different objects belonging to disparate categories. In a dataset, some images may contain 10 examples of a hen but only 1 example of a person, while others may contain 10 people and 2 hens, and another one might contain 5 cats and 5 hens, with no people. In such cases, a random split of images may skew the object category-wise distribution. Stratified sampling on the other hand can partition the data in a way that the resulting object category distribution is balanced.

Stratified sampling is, thus, a more fair way of data splitting, such that the machine learning model will be trained and validated on the same data distribution.

### *Cross-Validation*

Cross-Validation or K-Fold Cross-Validation is a more robust technique for data splitting, where a model is trained and evaluated "K" times on different samples. Let us understand this with an example.

Suppose we have a balanced, 2-class dataset consisting of 1000 images of raccoons and ringtails (to be used for training and validation only). Now, we want to perform a 5-Fold cross-validation. We first split the datasets into 5 equal and non-overlapping parts: each consisting of 200 images; label them as Parts 1, 2, 3, 4, and 5. Each of these subsets of 200 images consists of mutually different samples.

Now, we will create 5 complete datasets (labeled as Datasets 1-5) using Parts 1-5 in the following manner: For Dataset-1, use Part-1 as the validation set, and consolidate Parts 2-5 to create the training set; for Dataset-2, use Part-2 as the validation set, and consolidate Parts 1, 3, 4 and 5 to create the training set, and so on. Notice that since each part consists of 20% of the data of the original dataset, each of Datasets 1-5 has an 80%-20% train-validation split ratio. Generalizing, each K-Fold cross-validation dataset has (100/K)% data in its validation set (here, 100/5 = 20% was in validation set).

Using K-Fold cross-validation exposes the machine learning model (trained independently of each other) to different distributions of data ("K" times to be exact). This alleviates any bias that may occur while selecting data in the training and validation sets. Average and standard deviation values are usually reported while using K-Fold cross-validation schemes.

A simple analysis can suggest that K-Fold cross-validation also suffers from the same problem as random sampling. Data distribution may get skewed when the "K" parts or subsets of the datasets are created. In the example above, it may so happen that Part-2 of the dataset consists of 200 images of raccoons and no images of ringtails.

Thus, the "Stratified K-Fold Cross-Validation" technique avoids such inconsistencies. Similar to stratified sampling, the class-ratio of the data is maintained while generating the "K" subsets or parts of the data. Thus, the same class distribution is carried forward when these "K" parts are coalesced to form the final complete datasets.

**3 common pitfalls in the training data split**

Finally, let's briefly discuss common mistakes that data scientists make when building their models.

Low-quality training data

The quality of the training data is crucial for the model performance to improve.

If the training data is "garbage," one cannot expect the model to perform well.

Moreover, since the machine learning algorithms are sensitive to the training data, even small variations/errors in the training set can lead to significant errors in the model performance.

**Overfitting**

Overfitting happens when the machine learning model memorizes the pattern in the training data to such an extent that it fails to classify unseen data.

The noise or fluctuations in the training data is seen as features and learned by the model. This leads to the model outperforming in the training set but poor performance in the validation and testing sets.

**Overemphasis on Validation and Test Set metrics**

The validation set metric is the one that decides the path of the training of the model.

After each epoch, the machine learning model is evaluated on the validation set. Based on the validation set metrics, the corresponding loss terms are calculated, and the hyperparameters are modified.

Metrics should be chosen so that they bring a positive effect on the overall trajectory of the model performance.

**Pull, split, and load data for training a model**

Using V7, data can be uploaded to a dataset, new versions of a collaborative dataset can be downloaded, and split into training, testing, and validation sets. Using simple PyTorch scripts, you can then use the data to train a deep learning model in Darwin. PyTorch is a Machine
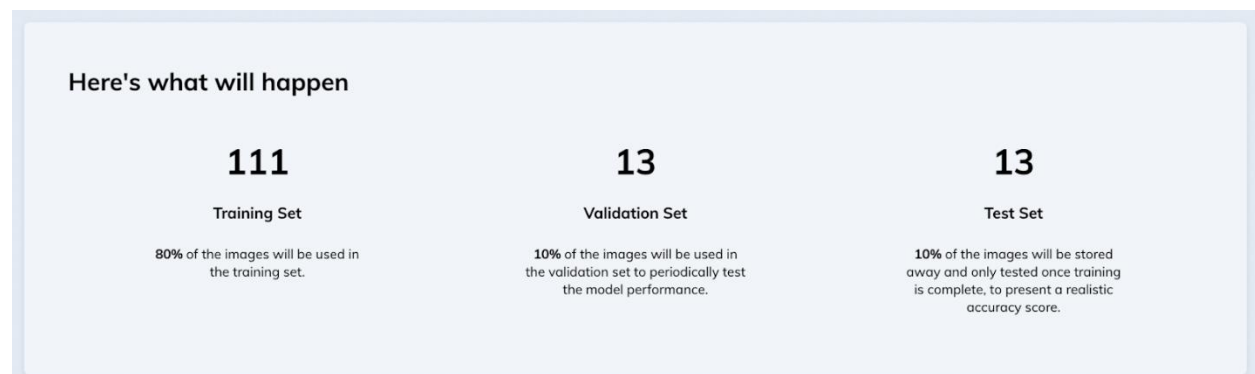
Learning framework that allows you to train Neural Networks. To take advantage of this, the Darwin SDK allows some integrations with PyTorch, making training models in Darwin much simpler for programmers who are used to PyTorch.

 Let us see how we can do this.

1. *Pulling Data*- Using v7, you can download a remote dataset (images and annotations) to the datasets directory. A simple python script can help accomplish this, as shown below:

2. *Loading a Dataset*- Darwin datasets can be directly used in your PyTorch-based code using the "get_dataset" function, which works for classification, instance segmentation, and semantic segmentation. The simple Python code for this is as shown below:

3. *Splitting Data*- You can split the data into training, testing, and validation sets using the "darwin.dataset.split_manager" command in the Darwin SDK. All you need is the dataset path for this. You can specify the percentage of data in the validation and testing sets or let them be the default values of 10% and 20%, respectively. The dataset splitting will be a process unless a list of stratified types is explicitly provided.

Here's what will happen

| 111 | 13 | 13 |
|---|---|---|
| **Training Set** | **Validation Set** | **Test Set** |
| **80%** of the images will be used in the training set. | **10%** of the images will be used in the validation set to periodically test the model performance. | **10%** of the images will be stored away and only tested once training is complete, to present a realistic accuracy score. |

Now you are all set to train a Neural Network in V7, a detailed guide to which can be found here.

Train, Validation, and Test Set: Key takeaways

Finally, here's a recap of everything we've learned:

- Training data is the set of the data on which the actual training takes place. Validation split helps to improve the model performance by fine-tuning the model after each epoch. The test set informs us about the final accuracy of the model after completing the training phase.

- The training set should not be too small; else, the model will not have enough data to learn. On the other hand, if the validation set is too small, then the evaluation metrics like accuracy, precision, recall, and F1 score will have large variance and will not lead to the proper tuning of the model.

- In general, putting 80% of the data in the training set, 10% in the validation set, and 10% in the test set is a good split to start with.

- The optimum split of the test, validation, and train set depends upon factors such as the use case, the structure of the model, dimension of the data, etc.