

# Tumor-CIFAR: Toy Datasets to Simulate Longitudinal Lung Cancer Prediction -- Supplementary Materials

**Abstract.** In this material, we describe how the demonstration dataset, Tumor-CIFAR, is created. This is the supplementary document of the DLSTM paper. The references and mathematic computing are provided. The code and datasets are publicly available at <https://github.com/MASILab/tumor-cifar>.

## 1 Introduction

Medical images are a precious and rare resources. It is hard to share large amounts of lung CT scans, especially for longitudinal data. CIFAR10 [1] dataset consists of 60K color natural images with the size of  $32 \times 32$ , which is convenient to verify the effectiveness of algorithms without excessive computing cost. Compared with MNIST [2], CIFAR10 consists more diverse image patterns. We manually add synthetic “nodules” to CIFAR10 with approximate clinical knowledge of lung nodules. Motivated by [3], we model the growth rate of malignant nodules at three times than the benignant one. This information guides us to simulate the Tumor-CIFAR dataset.

## 2 The Implementing of Tumor-CIFAR

**Step 1.** Create 5 time-intervals  $X = [x_0, x_1, x_2, x_3, x_4]$  with absolute of Gaussian distribution:

$$X \sim |N(\mu, \sigma^2)|$$

where  $\mu = 1.67$  and  $\sigma^2 = 1$  for version 1,  $\mu = 5$  and  $\sigma^2 = 3$  for version 2, and  $N(\mu, \sigma^2)$  represents is the Gaussian distribution with  $\mu$  as mean and  $\sigma^2$  as variance. The operation  $|\cdot|$  is the absolute operation.

**Step 2.** Create 5 time-points  $T = [t_0, t_1, t_2, t_3, t_4]$  by boosting the 5 time-intervals, where

$$t_i = \sum_{k=0}^i x_k$$

**Step 3.** Create the nodule growth rate  $g$ , where  $g \sim |N(1, 0.2)|$ .

**Step 4.** Randomly choose threshold  $th$ , where  $th$  belongs to uniform distribution  $th \sim U(0,1)$ . If  $th < 0.5$  we regard this subject as BENIGN, otherwise, this subject is MALIGNANT.

We create two different versions in step 5 according the rule learned from [1] that the growth rate of malignant nodules is about 3 times than the benign one (we call this NODULE RULE in the following).

**Step 5.** For version 1, the growth rate  $g$  of MALIGNANT subjects would increase to  $3 \times g$  to match the NODULE RULE. For version 2, the time points  $T$  of MALIGNANT subjects would reduce to  $T/3$  to match the NODULE RULE.

**Step 6.** Randomly choose two nodule locations, and generate the nodules. The nodule location  $(x_l, y_l)$  is chosen by

$$(x_l, y_l) \sim (U(7,25), U(7,25))$$

Note that the image size is  $32 \times 32$ . Then the nodule size  $s_i$  is computed by

$$s_i = t_i \times g$$

And the intensity of the nodule pixel  $p_i^{nodule}$  is calculated by

$$p_i^{nodule} = (1 + \alpha \cdot s_i) \cdot p_i^{origin}$$

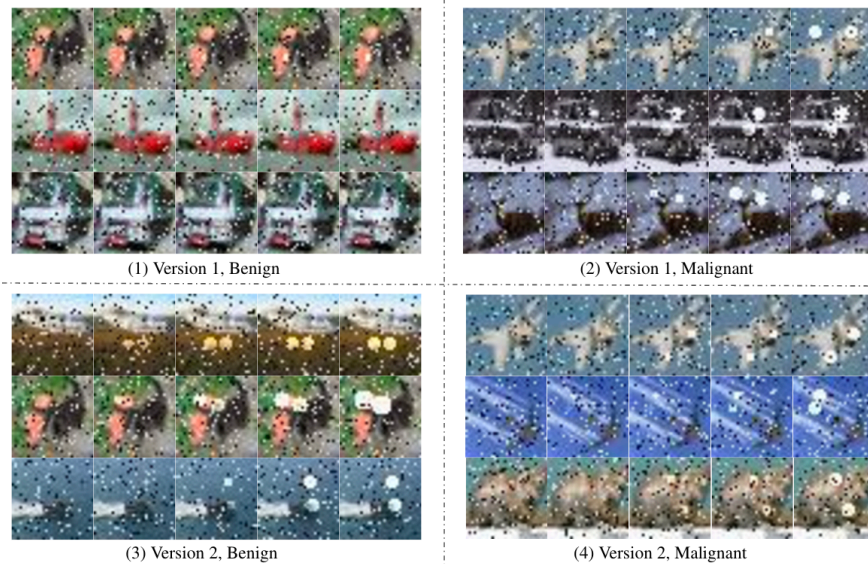
where  $i$  is the image index, and  $p_i^{origin}$  is the pixel value before adding nodule,  $\alpha = 0.1$ .

**Step 7.** Finally, we add 10% salt-and-pepper noise among all pixels.

### 3 Example images of Tumor-CIFAR

Fig. 1 presents image examples. We can see both salt-and-pepper noise and “nodules” on the images. In version 1 Tumor-CIFAR, images are modeled with same time interval distribution, but with different nodule sizes between benign and malignant. While in version 2, the data with same nodule size distribution, different time intervals between benign and malignant.

The source code of generating these datasets, the raw datasets (some examples), and related raw files (including time intervals and nodule sizes) can be found at the link: <https://github.com/MASILab/tumor-cifar>.



**Fig. 1.** Example images from two version Tumor-CIFAR datasets. In version 1, the nodule size of malignant image is apparently bigger than the benign ones because the time interval distributions are the same between Benign and Malignant. Since nodule size distribution are the same in version 2, we barely can see the nodule differences according only image data.

## 4 References

1. Krizhevsky, A. and G. Hinton, *Learning multiple layers of features from tiny images*. 2009, Citeseer.
2. LeCun, Y., C. Cortes, and C.J.J.U.h.y.l.c.e.m. Burges, *The MNIST database of handwritten digits, 1998*. 1998. **10**: p. 34.
3. Duhaylongsod, F.G., et al., *Lung tumor growth correlates with glucose metabolism measured by fluoride-18 fluorodeoxyglucose positron emission tomography*. 1995. **60**(5): p. 1348-1352.