

Modelo de deep learning fast.ai para a detecção de neoplasias cutâneas

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

```
1 !ls
```

drive sample_data

Bibliotecas e Banco de Dados

```
1 #Carregar toda vez que acionar a biblioteca FASTAI
2 #Recarregar notebook quaisquer alterações feitas em qualquer biblioteca usada.
3 %reload_ext autoreload
4 %autoreload 2
5 #Garantir que todos os gráficos plotados sejam mostrados
6 %matplotlib inline
```

```
1 #Carga bibliotecas
2 from fastai.vision.all import *
3 from fastai.metrics import *
4 import pandas as pd
5 from pathlib import Path
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import os
9 from glob import glob
10 import seaborn as sns
```

```
1 base_skin_dir = os.path.join('', 'drive/MyDrive/archive')
2
3 #Banco de metadados de imagens
4 csv_path = "drive/MyDrive/archive/HAM10000_metadata.csv"
5 skin_df = pd.read_csv(csv_path)
```

Análise Exporatória de Dados

```
1 # Amostra descritiva do banco de imagens
2 skin_df.sort_values(by="image_id")
```

	lesion_id	image_id	dx	dx_type	age	sex	localization
4349	HAM_0000550	ISIC_0024306	nv	follow_up	45.0	male	trunk
4263	HAM_0003577	ISIC_0024307	nv	follow_up	50.0	male	lower extremity
4217	HAM_0001477	ISIC_0024308	nv	follow_up	55.0	female	trunk
3587	HAM_0000484	ISIC_0024309	nv	follow_up	40.0	male	trunk
1451	HAM_0003350	ISIC_0024310	mel	histo	60.0	male	chest
...
1721	HAM_0004304	ISIC_0034316	mel	histo	85.0	male	upper extremity
1888	HAM_0006376	ISIC_0034317	mel	histo	70.0	female	lower extremity
121	HAM_0000344	ISIC_0034318	bkl	histo	55.0	male	trunk
7440	HAM_0000747	ISIC_0034319	nv	histo	30.0	male	trunk
7363	HAM_0002244	ISIC_0034320	nv	histo	25.0	female	chest

10015 rows × 7 columns

```
1 # Contagem de imagens por neoplasia
2 dx = skin_df['dx'].value_counts().sort_index()
3 print(dx)
```

akiec 327
bcc 514

```
bkl      1099
df        115
mel      1113
nv       6705
vasc      142
Name: dx, dtype: int64
```

```
1 # Contagem de imagens por categorias de neoplasia
2 categories = dx.index.values
3 print(categories)
4
5 counts = dx.values
6 print(counts)
```

```
['akiec' 'bcc' 'bkl' 'df' 'mel' 'nv' 'vasc']
[ 327  514 1099  115 1113 6705  142]
```

```
1 #Estatísticas básicas dos dados
2 skin_df.describe(exclude=[np.number])
```

	lesion_id	image_id	dx	dx_type	sex	localization
count	10015	10015	10015	10015	10015	10015
unique	7470	10015	7	4	3	15
top	HAM_0003789	ISIC_0027419	nv	histo	male	back
freq	6	1	6705	5340	5406	2192

```
1 #Identificacao da estrutura do banco de dados
2 skin_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10015 entries, 0 to 10014
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  -
0   lesion_id       10015 non-null  object
1   image_id        10015 non-null  object
2   dx              10015 non-null  object
3   dx_type         10015 non-null  object
4   age             9958 non-null   float64
5   sex             10015 non-null  object
6   localization    10015 non-null  object
dtypes: float64(1), object(6)
memory usage: 547.8+ KB
```

Podemos observar que não há valores nulos

```
1 imageid_path_dict = {os.path.splitext(os.path.basename(x))[0]: x
2                       for x in glob(os.path.join(base_skin_dir, '*', '*.jpg'))}
3
4 lesion_type_dict = {
5     'nv': 'nevos melanocíticos',
6     'mel': 'dermatofibroma',
7     'bkl': 'Lesões semelhantes a ceratose benigna',
8     'bcc': 'Carcinoma basocelular',
9     'akiec': 'Queratose actínica',
10    'vasc': 'lesões vasculares',
11    'df': 'Dermatofibroma'
12 }
```

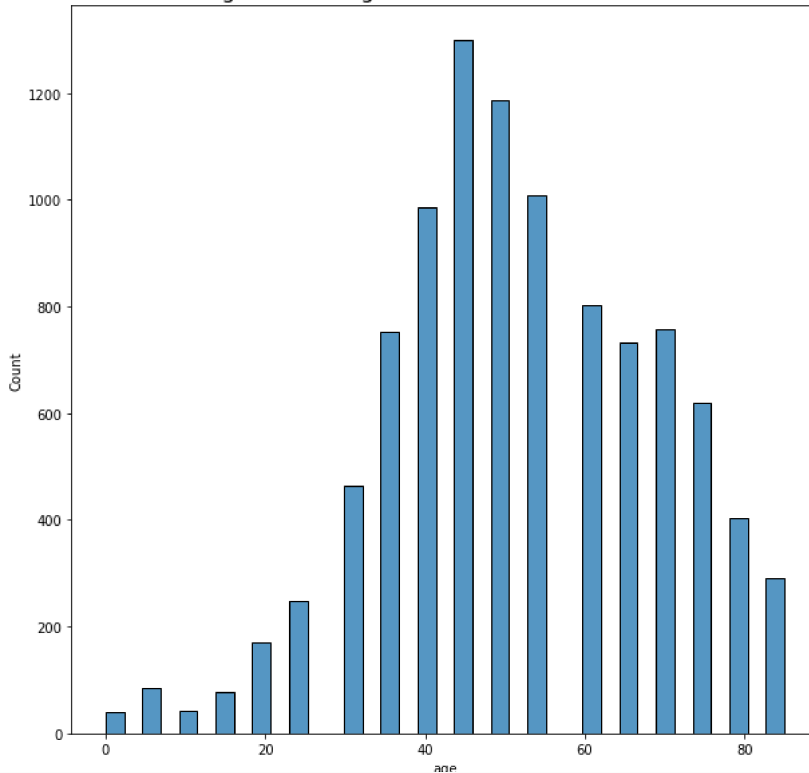
```
1 tile_df = pd.read_csv(os.path.join(base_skin_dir, 'HAM10000_metadata.csv'))
2 tile_df['path'] = tile_df['image_id'].map(imageid_path_dict.get)
3 tile_df['cell_type'] = tile_df['dx'].map(lesion_type_dict.get)
4 tile_df['cell_type_idx'] = pd.Categorical(tile_df['cell_type']).codes
5 tile_df.sample(3)
```

	lesion_id	image_id	dx	dx_type	age	sex	localization	path
80	HAM_0000959	ISIC_0030189	bkl	histo	75.0	female	face	drive/MyDrive/archive/HAM10000_images_part_2/ISIC_0030189.jpg

```
1 bar, ax = plt.subplots(figsize=(10,10))
2 sns.histplot(skin_df['age'])
3 plt.title('Figura 5: Histograma de Idade dos Pacientes', size=16)
```

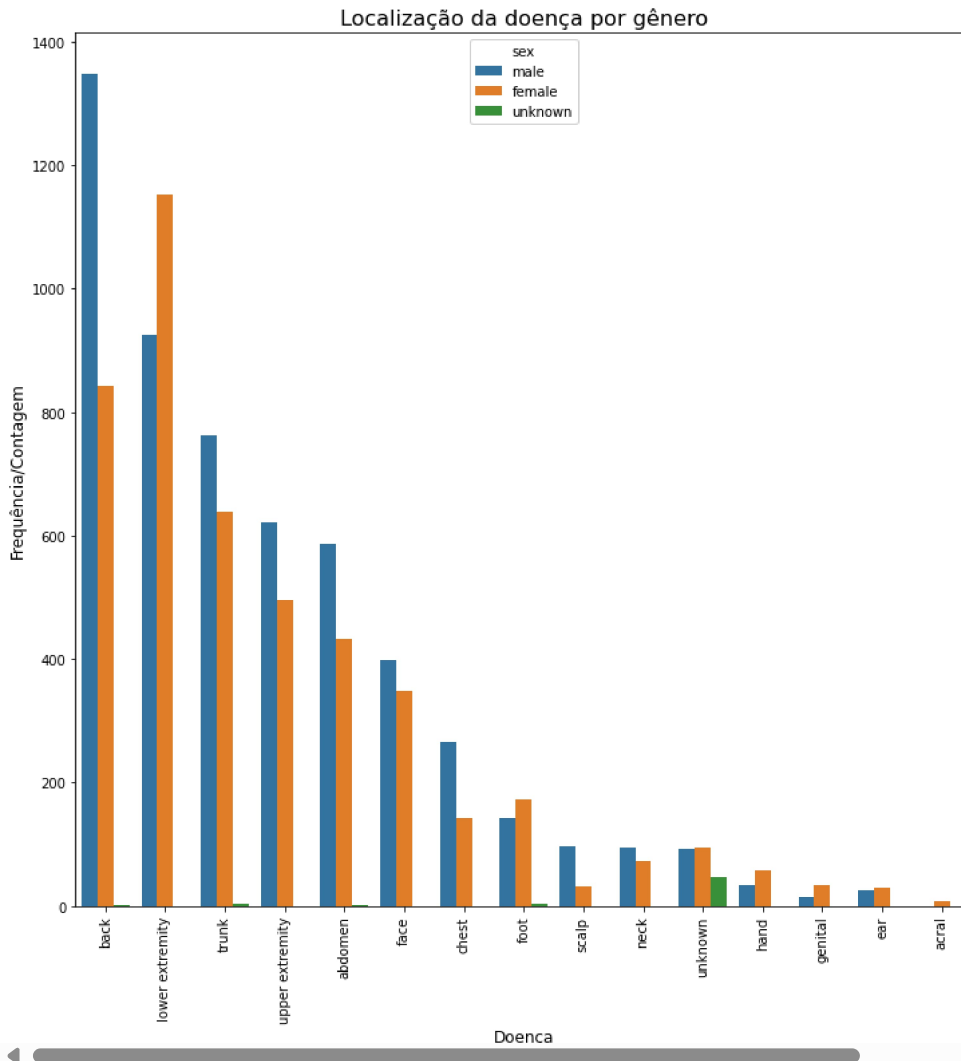
Text(0.5, 1.0, 'Figura 5: Histograma de Idade dos Pacientes')

Figura 5: Histograma de Idade dos Pacientes



```
1 value = skin_df[['localization', 'sex']].value_counts().to_frame()
2 value.reset_index(level=[1,0 ], inplace=True)
3 temp = value.rename(columns = {'localization':'location', 0: 'count'})
4
5 bar, ax = plt.subplots(figsize = (12, 12))
6 sns.barplot(x = 'location', y='count', hue = 'sex', data = temp)
7 plt.title('Localização da doença por gênero', size = 16)
8 plt.xlabel('Doença', size=12)
9 plt.ylabel('Frequência/Contagem', size=12)
10 plt.xticks(rotation = 90)
```

(array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14]),
<a list of 15 Text major ticklabel objects>)



Amostra estatística de 1000 registros

```

1 # Definindo o gráfico de coluna por distribuição
2 def plotPerColumnDistribution(df, nGraphShown, nGraphPerRow):
3     nunique = df.nunique()
4     df = df[[col for col in df if nunique[col] > 1 and nunique[col] < 50]] # Para fins de exibição, escolha colunas que tenham entre
5     nRow, nCol = df.shape
6     columnNames = list(df)
7     nGraphRow = (nCol + nGraphPerRow - 1) // nGraphPerRow
8     plt.figure(num = None, figsize = (6 * nGraphPerRow, 8 * nGraphRow), dpi = 80, facecolor = 'w', edgecolor = 'k')
9     for i in range(min(nCol, nGraphShown)):
10         plt.subplot(nGraphRow, nGraphPerRow, i + 1)
11         columnDf = df.iloc[:, i]
12         if (not np.issubdtype(type(columnDf.iloc[0]), np.number)):
13             valueCounts = columnDf.value_counts()
14             valueCounts.plot.bar()
15         else:
16             columnDf.hist()
17             plt.ylabel('qde')
18             plt.xticks(rotation = 90)
19             plt.title(f'{columnNames[i]} (column {i})')
20     plt.tight_layout(pad = 1.0, w_pad = 1.0, h_pad = 1.0)
21     plt.show()

```

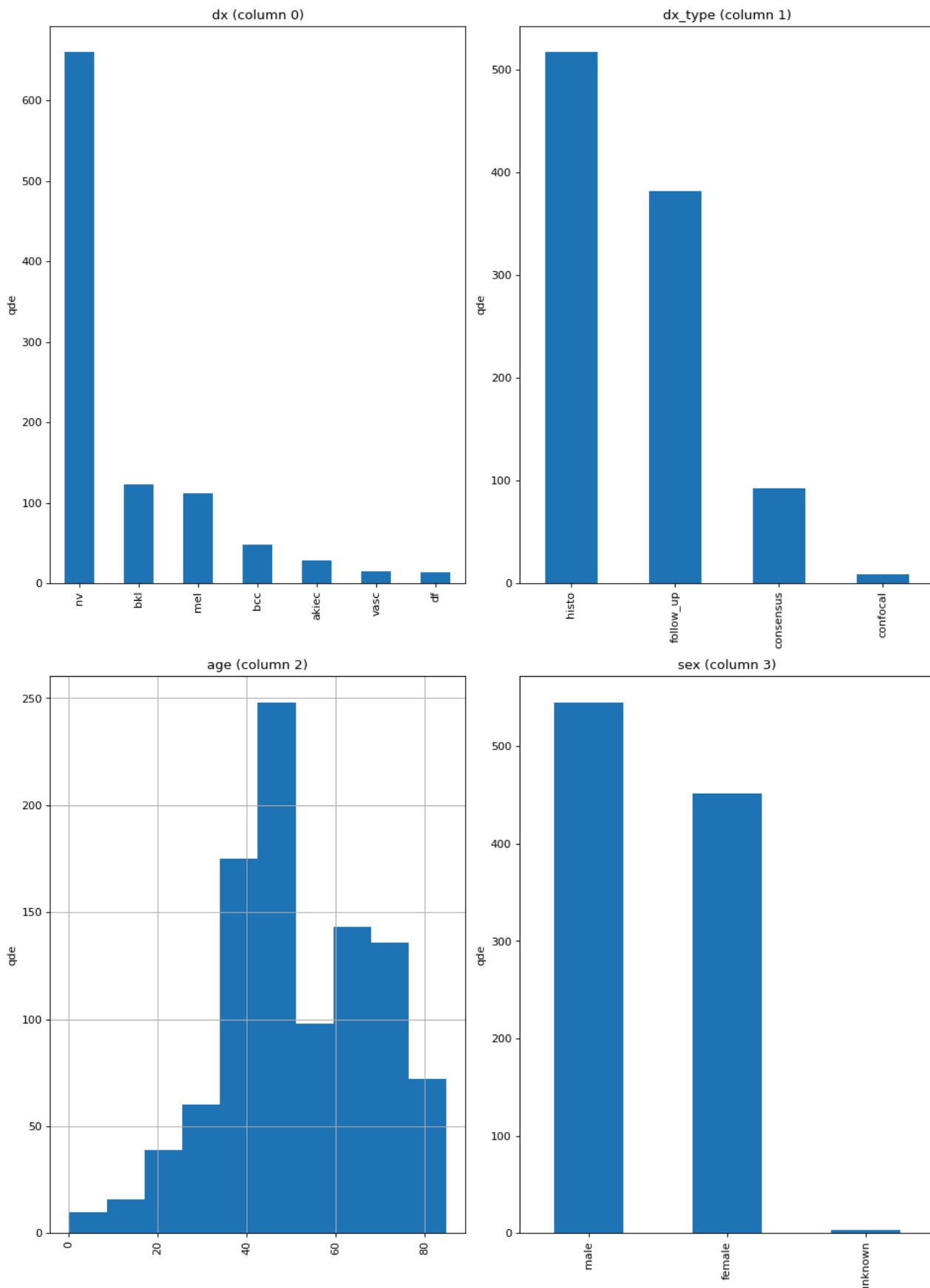
```

1 # Gerando a amostra
2 df1 = skin_df.sample(n=1000, random_state=1)
3
4 df1.dataframeName = 'HAM10000_metadata.csv'
5 nRow, nCol = df1.shape
6 print(f'Há {nRow} linhas e {nCol} colunas')

```

Há 1000 linhas e 7 colunas

```
1 plotPerColumnDistribution(df1, 4, 2)
```



✓ Verificando o equilíbrio dos dados

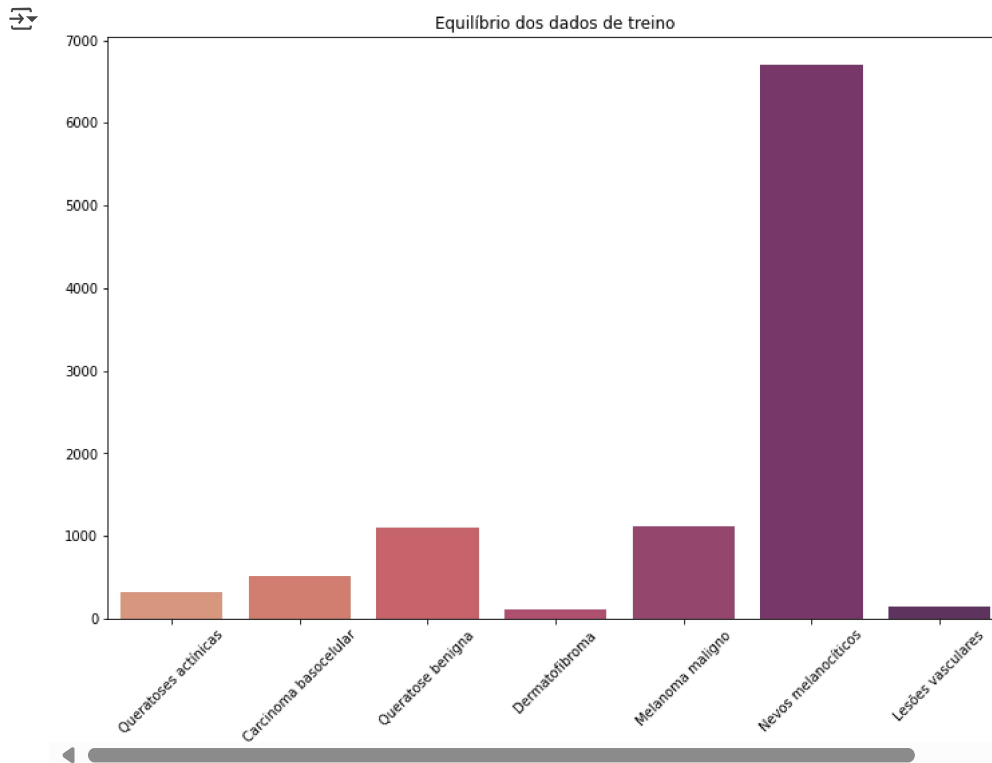
```
1 categories = ['Queratoses actínicas', 'Carcinoma basocelular', 'Queratose benigna', 'Dermatofibroma', 'Melanoma maligno', 'Nevos melâ  
2  
3 num_classes = len(categories)
```

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 def plot_equilibrio(categories, counts):
5
6     plt.figure(figsize=(12, 8))
7
8     sns_bar = sns.barplot(x=categories, y=counts, saturation=0.75, palette = "flare")
9     sns_bar.set_xticklabels(categories, rotation=45)
10
11     plt.title('Equilíbrio dos dados de treino')
12     plt.show()

```

```
1 plot_equilibrio(categories, counts)
```



```

1 #Definindo a localização das imagens
2 path = Path('drive/MyDrive/archive/')
3 Path.BASE_PATH = path
4 path.ls()

```

```

(#7)
[Path('HAM10000_metadata.csv'), Path('hmnist_28_28_L.csv'), Path('hmnist_28_28_RGB.csv'), Path('hmnist_8_8_L.csv'), Path('hmnist_8_8_RGB.csv')]

```

✓ Renomeando variáveis

```

1 short_to_full_name_dict = {
2     "akiec" : "Doença de Bowen", # forma muito precoce de câncer de pele
3     "bcc" : "carcinoma basocelular", # câncer basocelular ou câncer de pele branca
4     "bkl" : "lesões semelhantes a ceratose benigna", # tumor de pele não canceroso
5     "df" : "dermatofibroma", # saliências arredondadas não cancerosas
6     "mel" : "melanoma", # câncer de pele
7     "nv" : "nevos melanocíticos", # mole não canceroso
8     "vasc" : "lesões vasculares", # condição da pele
9 }

```

✓ Obtendo imagens do arquivo

```

1 # retorna apenas a coluna dx e id da imagem
2 img_to_class_dict = skin_df.loc[:, ["image_id", "dx"]]
3 # retorna colunas como listas em um dict
4 img_to_class_dict = img_to_class_dict.to_dict('list')
5 # retorna um id de imagem de mapeamento dict para o nome da doença
6 img_to_class_dict = {img_id : short_to_full_name_dict[disease] for img_id,disease in zip(img_to_class_dict['image_id'], img_to_class_dict['dx'])}
7 [x for x in img_to_class_dict.items()][:5]

```

```

[('ISIC_0027419', 'benign keratosis-like lesions'),
 ('ISIC_0025030', 'benign keratosis-like lesions'),
 ('ISIC_0026769', 'benign keratosis-like lesions'),
 ('ISIC_0025661', 'benign keratosis-like lesions'),
 ('ISIC_0031633', 'benign keratosis-like lesions')]

```

```

1 # path.stem retorna o nome do arquivo sem sufixo
2 def get_label_from_dict(path):
3     return img_to_class_dict[path.stem]

```

✓ Construindo um bloco de imagens DataBlock

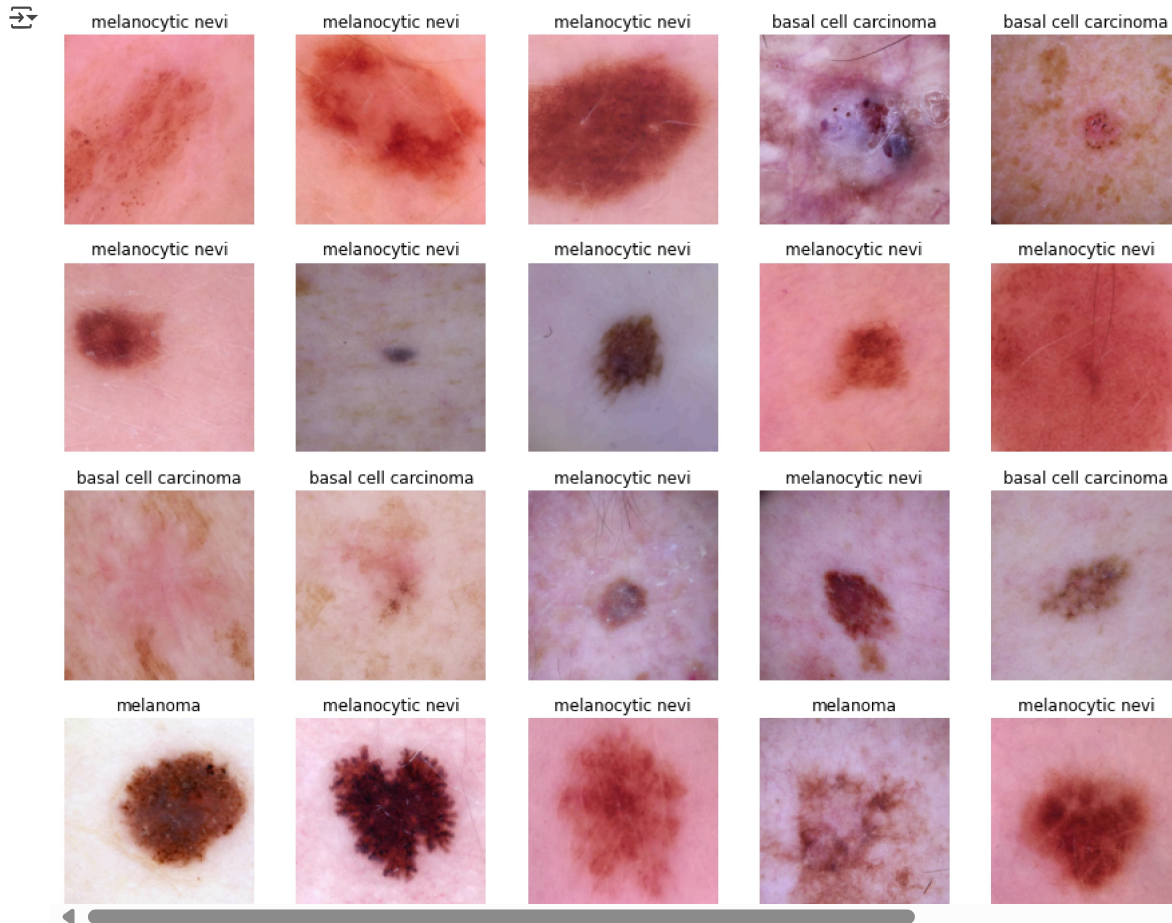
```

1 dblock = DataBlock(
2     # Designação das variáveis independentes e dependentes
3     blocks = (ImageBlock, CategoryBlock),
4     # Para obter uma lista desses arquivos e retornar uma lista de todas as imagens nesse caminho
5     get_items = get_image_files,
6     # Divida nossos conjuntos de treinamento e validação aleatoriamente
7     splitter = RandomSplitter(valid_pct=0.2, seed=42),
8     # Estamos dizendo ao fastai qual função chamar para criar os rótulos em nosso conjunto de dados, no nosso caso é uma variável ir
9     get_y = get_label_from_dict,
10    # DihedralItem todas as 4 rotações de 90 graus e para cada:
11    # 2 flips horizontais -> 8 orientações
12    item_tfms=[Resize(448), DihedralItem()],
13    # Escolhe um recorte escalado aleatório de uma imagem e a redimensiona para o tamanho
14    batch_tfms=RandomResizedCrop(size=224, min_scale=0.75, max_scale=1.0))
15
16 img_path = "drive/MyDrive/archive/HAM10000_images_part_1"
17 # cria dataloader usando img_path
18 dls = dblock.dataloaders(img_path, bs=64) # bs = tamanho do lote

```

✓ Visualizando as imagens

```
1 dls.show_batch(max_n=20)
```



Observações dessas imagens serão anotadas abaixo. Primeiro, farei mais algumas verificações para confirmar que nossas categorias são apenas "doença de Bowen", 'carcinoma basocelular', 'lesões semelhantes a ceratose benigna', 'dermatofibroma', 'nevus melanocíticos',

'melanoma', 'lesões vasculares':

```
1 print(dls.vocab)
```

```
["Bowen's disease", 'basal cell carcinoma', 'benign keratosis-like lesions', 'dermatofibroma', 'melanocytic nevi', 'melanoma', 'vasc
```

Vamos visualizar o tamanho dos nossos conjuntos de dados:

```
1 len(dls.train_ds), len(dls.valid_ds)
```

```
(4000, 1000)
```

✓ Treine um modelo simples

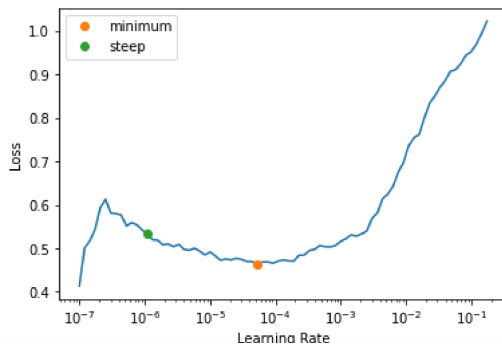
```
1 learn = vision_learner(dls,
2                         resnet18,
3                         metrics=accuracy)
4 weights='IMAGENET1K_V1'
5 learn.fine_tune(4)
```



epoch	train_loss	valid_loss	accuracy	time
0	2.107257	1.014550	0.733000	01:19

epoch	train_loss	valid_loss	accuracy	time
0	1.162495	0.803819	0.793000	01:23
1	0.910802	0.594211	0.835000	01:22
2	0.679382	0.481674	0.837000	01:23
3	0.553778	0.479971	0.841000	01:22

```
1 # Definindo a taxa de aprendizagem em relação as perdas
2 lr_min,lr_steep = learn.lr_find(suggest_funcs=(minimum, steep))
```



```
1 print(f"Minimum/10: {lr_min:.2e}, steepest point: {lr_steep:.2e}")
```



```
Minimum/10: 5.25e-06, steepest point: 1.10e-06
```

treinando fit_one_cycle por 3 ciclos para ter uma ideia de quão preciso seria o modelo com resnet34.

```
1 learn = vision_learner(dls,resnet34, metrics = accuracy)
2 weights='IMAGENET1K_V1'
3 learn.fit_one_cycle(3,1e-2)
```



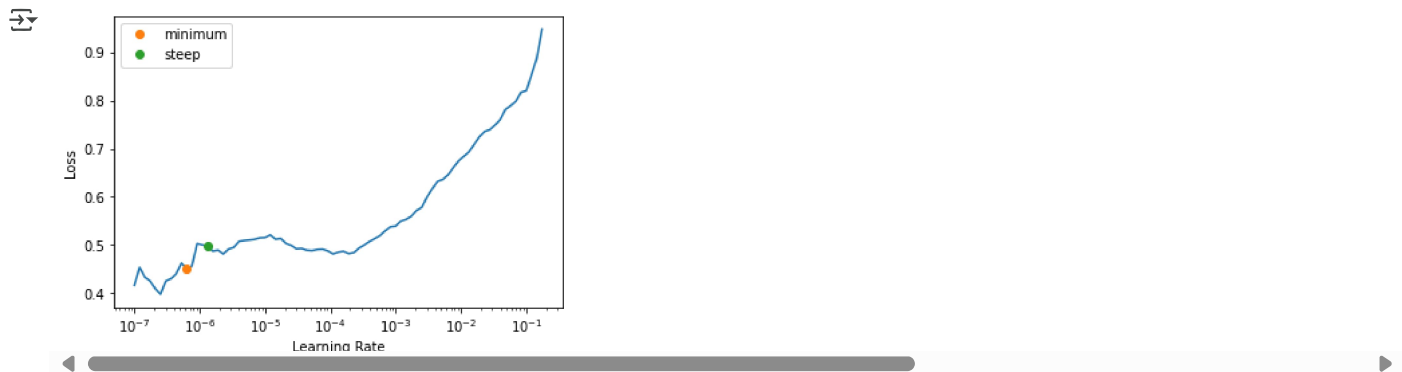
epoch	train_loss	valid_loss	accuracy	time
0	1.484736	1.193570	0.713000	01:26
1	0.915688	0.601399	0.812000	01:39
2	0.632683	0.461788	0.835000	01:28

Aplicando o algoritmo Resnet34 observamos uma taxa de acurácia de 83,5 %

✓ Descongelamento e Transferência de Aprendizagem

```
1 learn.unfreeze()
```

```
1 lr_min,lr_steep = learn.lr_find(suggest_funcs=(minimum, steep))
```



```
1 print(f"Minimum/10: {lr_min:.2e}, steepest point: {lr_steep:.2e}")
```

```
↕ Minimum/10: 6.31e-08, steepest point: 1.32e-06
```

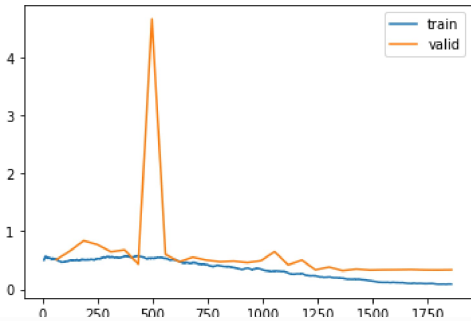
```
1 learn.fit_one_cycle(30 ,lr_max=slice(1e-4, 1e-2))
```



epoch	train_loss	valid_loss	accuracy	time
0	0.511097	0.513052	0.824000	01:39
1	0.503994	0.662045	0.796000	01:44
2	0.514480	0.839405	0.729000	01:41
3	0.530039	0.771729	0.793000	01:36
4	0.558221	0.647117	0.794000	01:45
5	0.567121	0.679792	0.792000	01:35
6	0.575276	0.434001	0.852000	01:36
7	0.539968	4.668460	0.733000	01:35
8	0.536520	0.603031	0.791000	01:34
9	0.484933	0.473026	0.845000	01:36
10	0.464415	0.553452	0.788000	01:37
11	0.420540	0.503143	0.831000	01:34
12	0.406005	0.475321	0.829000	01:36
13	0.376242	0.486368	0.830000	01:37
14	0.363214	0.460963	0.843000	01:37
15	0.349533	0.492556	0.830000	01:36
16	0.311974	0.649859	0.853000	01:36
17	0.282202	0.420920	0.866000	01:36
18	0.272512	0.503506	0.837000	01:35
19	0.235566	0.334372	0.886000	01:37
20	0.209646	0.382954	0.875000	01:36
21	0.189831	0.320385	0.896000	01:36
22	0.171440	0.348050	0.897000	01:35
23	0.149904	0.329796	0.898000	01:36
24	0.120866	0.334872	0.901000	01:36
25	0.114491	0.336159	0.906000	01:35
26	0.099070	0.340293	0.907000	01:37
27	0.098978	0.332594	0.907000	01:36
28	0.087987	0.331241	0.911000	01:35
29	0.087777	0.334828	0.911000	01:35

1 O modelo com fast.ai chegou a taxa de 91,1%

1 learn.recorder.plot_loss()



▼ salvando o modelo

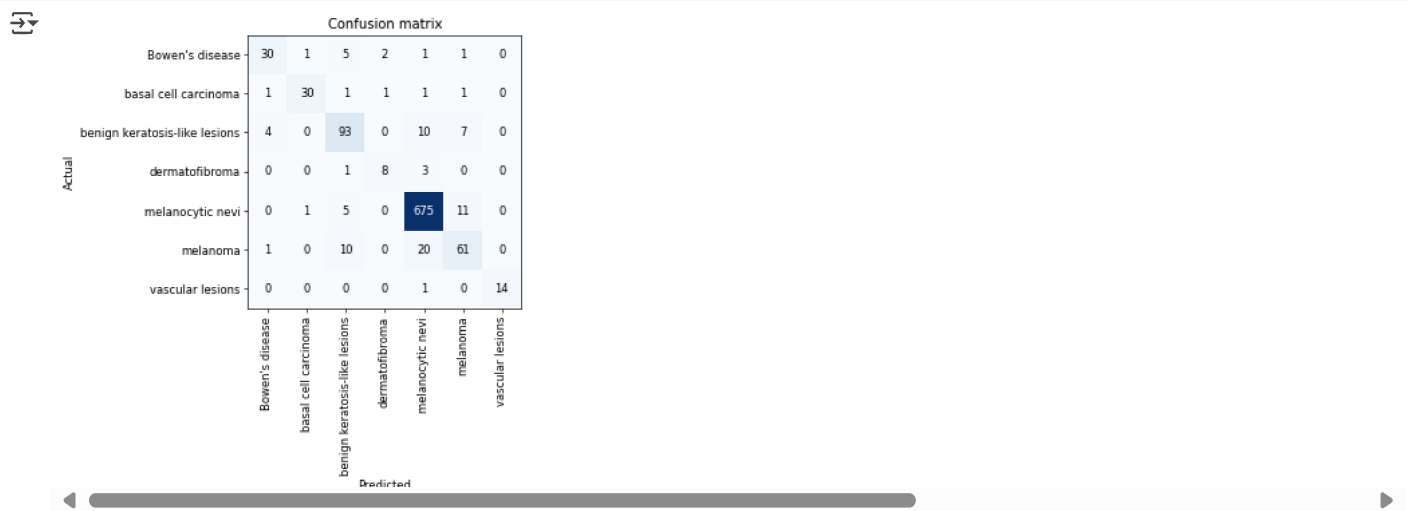
1 learn.save('model1')



Path('models/model1.pth')

Interpretação do modelo

```
1 interp = ClassificationInterpretation.from_learner(learn)
2 interp.plot_confusion_matrix(figsize=(6,6), dpi=60)
```



As 6 principais perdas

```
1 interp.plot_top_losses(6, nrows=6)
```

