In [1]:
```python
# Imports for this Notebook
import pandas as pd
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn import metrics
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.callbacks import EarlyStopping
```

Using TensorFlow backend.

In [2]:
```python
from sklearn import preprocessing
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import shutil
import os
import keras
import matplotlib
import matplotlib.pyplot as plt
%matplotlib inline
```

In [3]:
```python
# Encode text values to dummy variables(i.e. [1,0,0],[0,1,0],[0,0,1] for r
def encode_text_dummy(df, name):
    dummies = pd.get_dummies(df[name])
    for x in dummies.columns:
        dummy_name = "{}-{}".format(name, x)
        df[dummy_name] = dummies[x]
    df.drop(name, axis=1, inplace=True)


# Encode text values to a single dummy variable.  The new columns (which d
# at every location where the original column (name) matches each of the t
# each target value.
def encode_text_single_dummy(df, name, target_values):
    for tv in target_values:
        l = list(df[name].astype(str))
        l = [1 if str(x) == str(tv) else 0 for x in l]
        name2 = "{}-{}".format(name, tv)
        df[name2] = l


# Encode text values to indexes(i.e. [1],[2],[3] for red,green,blue).
def encode_text_index(df, name):
    le = preprocessing.LabelEncoder()
    df[name] = le.fit_transform(df[name])
    return le.classes_
```

```python
# Encode a numeric column as zscores
def encode_numeric_zscore(df, name, mean=None, sd=None):
    if mean is None:
        mean = df[name].mean()

    if sd is None:
        sd = df[name].std()

    df[name] = (df[name] - mean) / sd


# Convert all missing values in the specified column to the median
def missing_median(df, name):
    med = df[name].median()
    df[name] = df[name].fillna(med)


# Convert all missing values in the specified column to the default
def missing_default(df, name, default_value):
    df[name] = df[name].fillna(default_value)


# Convert a Pandas dataframe to the x,y inputs that TensorFlow needs
def to_xy(df, target):
    result = []
    for x in df.columns:
        if x != target:
            result.append(x)
    # find out the type of the target column.  Is it really this hard? :(
    target_type = df[target].dtypes
    target_type = target_type[0] if hasattr(target_type, '__iter__') else
    # Encode to int for classification, float otherwise. TensorFlow likes
    if target_type in (np.int64, np.int32):
        # Classification
        dummies = pd.get_dummies(df[target])
        return df.as_matrix(result).astype(np.float32), dummies.as_matrix(
    else:
        # Regression
        return df.as_matrix(result).astype(np.float32), df.as_matrix([targ

# Nicely formatted time string
def hms_string(sec_elapsed):
    h = int(sec_elapsed / (60 * 60))
    m = int((sec_elapsed % (60 * 60)) / 60)
    s = sec_elapsed % 60
    return "{}:{:>02}:{:>05.2f}".format(h, m, s)


# Regression chart.
```

```python
def chart_regression(pred,y,sort=True):
    t = pd.DataFrame({'pred' : pred, 'y' : y.flatten()})
    if sort:
        t.sort_values(by=['y'],inplace=True)
    a = plt.plot(t['y'].tolist(),label='expected')
    b = plt.plot(t['pred'].tolist(),label='prediction')
    plt.ylabel('output')
    plt.legend()
    plt.show()

# Remove all rows where the specified column is +/- sd standard deviations
def remove_outliers(df, name, sd):
    drop_rows = df.index[(np.abs(df[name] - df[name].mean()) >= (sd * df[n
    df.drop(drop_rows, axis=0, inplace=True)


# Encode a column to a range between normalized_low and normalized_high.
def encode_numeric_range(df, name, normalized_low=-1, normalized_high=1,
                         data_low=None, data_high=None):
    if data_low is None:
        data_low = min(df[name])
        data_high = max(df[name])

    df[name] = ((df[name] - data_low) / (data_high - data_low)) \
               * (normalized_high - normalized_low) + normalized_low
```

In [4]:
```python
# This file is a CSV, just no CSV extension or headers
# Download from: http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html
df = pd.read_csv("/Users/tsaiflow/Downloads/kddcup.data_10_percent_correct
# df1 = pd.read_csv("/Users/tsaiflow/Downloads/kddcup.data.corrected", hea
print("Read {} rows.".format(len(df)))
#print("Read {} rows.".format(len(df1)))
# df = df.sample(frac=0.1, replace=False) # Uncomment this line to sample
df.dropna(inplace=True,axis=1) # For now, just drop NA's (rows with missi
# df1.dropna(inplace=True,axis=1)

# The CSV file has no column heads, so add them
df.columns = [
    'duration',
    'protocol_type',
    'service',#
    'flag',#
    'src_bytes',#
    'dst_bytes',#
    'land',#
    'wrong_fragment',
    'urgent',

    'hot',
    'num_failed_logins',
```

```
    'logged_in',
    'num_compromised',
    'root_shell',
    'su_attempted',
    'num_root',#
    'num_file_creations',
    'num_shells',
    'num_access_files',
    'num_outbound_cmds',
    'is_host_login',
    'is_guest_login',

    'count',
    'srv_count',
    'serror_rate',
    'srv_serror_rate',
    'rerror_rate',#
    'srv_rerror_rate',
    'same_srv_rate',
    'diff_srv_rate',#
    'srv_diff_host_rate',
    'dst_host_count',
    'dst_host_srv_count',
    'dst_host_same_srv_rate',
    'dst_host_diff_srv_rate',#
    'dst_host_same_src_port_rate',
    'dst_host_srv_diff_host_rate',
    'dst_host_serror_rate',
    'dst_host_srv_serror_rate',#
    'dst_host_rerror_rate',
    'dst_host_srv_rerror_rate',


    'outcome'
]
```

Read 494021 rows.

In [5]:  `df.head(5)`

Out[5]:

| | duration | protocol_type | service | flag | src_bytes | dst_bytes | land | wrong_fragment | urgent | ho |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | tcp | http | SF | 181 | 5450 | 0 | 0 | 0 | ( |
| **1** | 0 | tcp | http | SF | 239 | 486 | 0 | 0 | 0 | ( |
| **2** | 0 | tcp | http | SF | 235 | 1337 | 0 | 0 | 0 | ( |
| **3** | 0 | tcp | http | SF | 219 | 1337 | 0 | 0 | 0 | ( |
| **4** | 0 | tcp | http | SF | 217 | 2032 | 0 | 0 | 0 | ( |

5 rows × 42 columns

In [6]:
```
df_new = df[['service',
    'flag',
    'src_bytes',
    'dst_bytes',
    'land',
    'num_root',
    'rerror_rate',
    'diff_srv_rate',
    'dst_host_diff_srv_rate',
    'dst_host_srv_serror_rate',
    'outcome']]
```

In [7]:
```
# Now encode the feature vector

#encode_numeric_zscore(df, 'duration')
#encode_text_dummy(df, 'protocol_type')
encode_text_dummy(df_new, 'service')
encode_text_dummy(df_new, 'flag')
encode_numeric_zscore(df_new, 'src_bytes')
encode_numeric_zscore(df_new, 'dst_bytes')
encode_text_dummy(df_new, 'land')
#encode_numeric_zscore(df, 'wrong_fragment')
# encode_numeric_zscore(df, 'urgent')
# encode_numeric_zscore(df, 'hot')
# encode_numeric_zscore(df, 'num_failed_logins')
# encode_text_dummy(df, 'logged_in')
# encode_numeric_zscore(df, 'num_compromised')
# encode_numeric_zscore(df, 'root_shell')
# encode_numeric_zscore(df, 'su_attempted')
encode_numeric_zscore(df_new, 'num_root')
# encode_numeric_zscore(df, 'num_file_creations')
# encode_numeric_zscore(df, 'num_shells')
# encode_numeric_zscore(df, 'num_access_files')
# encode_numeric_zscore(df, 'num_outbound_cmds')
```

```python
# encode_numeric_zscore(df, 'num_outbound_cmds')
# encode_text_dummy(df, 'is_host_login')
# encode_text_dummy(df, 'is_guest_login')
# encode_numeric_zscore(df, 'count')
# encode_numeric_zscore(df, 'srv_count')
# encode_numeric_zscore(df, 'serror_rate')
# encode_numeric_zscore(df, 'srv_serror_rate')
encode_numeric_zscore(df_new, 'rerror_rate')
# encode_numeric_zscore(df, 'srv_rerror_rate')
# encode_numeric_zscore(df, 'same_srv_rate')
encode_numeric_zscore(df_new, 'diff_srv_rate')
# encode_numeric_zscore(df, 'srv_diff_host_rate')
# encode_numeric_zscore(df, 'dst_host_count')
# encode_numeric_zscore(df, 'dst_host_srv_count')
# encode_numeric_zscore(df, 'dst_host_same_srv_rate')
encode_numeric_zscore(df_new, 'dst_host_diff_srv_rate')
# encode_numeric_zscore(df, 'dst_host_same_src_port_rate')
# encode_numeric_zscore(df, 'dst_host_srv_diff_host_rate')
# encode_numeric_zscore(df, 'dst_host_serror_rate')
encode_numeric_zscore(df_new, 'dst_host_srv_serror_rate')
# encode_numeric_zscore(df, 'dst_host_rerror_rate')
# encode_numeric_zscore(df, 'dst_host_srv_rerror_rate')
outcomes = encode_text_index(df_new, 'outcome')
num_classes = len(outcomes)


# display 5 rows


df_new.dropna(inplace=True,axis=1)
df_new[0:5]
# This is the numeric feature vector, as it goes to the neural net
```

/anaconda/lib/python3.6/site-packages/ipykernel_launcher.py:6: Setting
WithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-
docs/stable/indexing.html#indexing-view-versus-copy
(http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-vi
ew-versus-copy)

Out[7]:

|   | src_bytes | dst_bytes | num_root | rerror_rate | diff_srv_rate | dst_host_diff_srv_rate | dst_host_srv_s |
|---|-----------|-----------|----------|-------------|---------------|------------------------|----------------|
| 0 | -0.002879 | 0.138664  | -0.00564 | -0.24796    | -0.255243     | -0.282866              |                |
| 1 | -0.002820 | -0.011578 | -0.00564 | -0.24796    | -0.255243     | -0.282866              |                |
| 2 | -0.002824 | 0.014179  | -0.00564 | -0.24796    | -0.255243     | -0.282866              |                |
| 3 | -0.002840 | 0.014179  | -0.00564 | -0.24796    | -0.255243     | -0.282866              |                |

| | | | | | | |
|---|---|---|---|---|---|---|
| **4** | -0.002842 | 0.035214 | -0.00564 | -0.24796 | -0.255243 | -0.282866 |

5 rows × 87 columns

```python
In [8]: df1 = pd.read_csv("/Users/tsaiflow/Downloads/kddcup.data.corrected", head
        print("Read {} rows.".format(len(df1)))
        df1.dropna(inplace=True,axis=1) # For now, just drop NA's (rows with miss
        df1.columns = [
            'duration',
            'protocol_type',
            'service',#
            'flag',#
            'src_bytes',#
            'dst_bytes',#
            'land',#
            'wrong_fragment',
            'urgent',

            'hot',
            'num_failed_logins',
            'logged_in',
            'num_compromised',
            'root_shell',
            'su_attempted',
            'num_root',#
            'num_file_creations',
            'num_shells',
            'num_access_files',
            'num_outbound_cmds',
            'is_host_login',
            'is_guest_login',

            'count',
            'srv_count',
            'serror_rate',
            'srv_serror_rate',
            'rerror_rate',#
            'srv_rerror_rate',
            'same_srv_rate',
            'diff_srv_rate',#
            'srv_diff_host_rate',
            'dst_host_count',
            'dst_host_srv_count',
            'dst_host_same_srv_rate',
            'dst_host_diff_srv_rate',#
            'dst_host_same_src_port_rate',
            'dst_host_srv_diff_host_rate',
            'dst_host_serror_rate',
            'dst_host_srv_serror_rate',#
            'dst_host_rerror_rate'
```

```
    dst_host_rerror_rate',
    'dst_host_srv_rerror_rate',


    'outcome'
]

df1_new = df1[['service',
    'flag',
    'src_bytes',
    'dst_bytes',
    'land',
    'num_root',
    'rerror_rate',
    'diff_srv_rate',
    'dst_host_diff_srv_rate',
    'dst_host_srv_serror_rate',
    'outcome']]

encode_text_dummy(df1_new, 'service')
encode_text_dummy(df1_new, 'flag')
encode_numeric_zscore(df1_new, 'src_bytes')
encode_numeric_zscore(df1_new, 'dst_bytes')
encode_text_dummy(df1_new, 'land')
encode_numeric_zscore(df1_new, 'num_root')
encode_numeric_zscore(df1_new, 'rerror_rate')
encode_numeric_zscore(df1_new, 'diff_srv_rate')
encode_numeric_zscore(df1_new, 'dst_host_diff_srv_rate')
encode_numeric_zscore(df1_new, 'dst_host_srv_serror_rate')
outcomes = encode_text_index(df1_new, 'outcome')
# num_classes = len(outcomes)

# display 5 rows

df1_new.dropna(inplace=True,axis=1)

del df1_new['service-aol']
del df1_new['service-harvest']
del df1_new['service-http_2784']
del df1_new['service-http_8001']
#del df1_new['is_host_login-1']
df1_new[0:5]
```

Read 4898431 rows.

/anaconda/lib/python3.6/site-packages/ipykernel_launcher.py:6: Setting
WithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-

docs/stable/indexing.html#indexing-view-versus-copy
(http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-vi
ew-versus-copy)

Out[8]:

|   | src_bytes | dst_bytes | num_root | rerror_rate | diff_srv_rate | dst_host_diff_srv_rate | dst_host_srv_s |
|---|-----------|-----------|----------|-------------|---------------|------------------------|----------------|
| 0 | -0.001720 | 0.068188  | -0.003285 | -0.248286  | -0.256056     | -0.282939              |                |
| 1 | -0.001777 | 0.005325  | -0.003285 | -0.248286  | -0.256056     | -0.282939              |                |
| 2 | -0.001698 | 0.000208  | -0.003285 | -0.248286  | -0.256056     | -0.282939              |                |
| 3 | -0.001701 | 0.001455  | -0.003285 | -0.248286  | -0.256056     | -0.282939              |                |
| 4 | -0.001695 | -0.000942 | -0.003285 | -0.248286  | -0.256056     | -0.282939              |                |

5 rows × 87 columns

In [9]:
```python
import numpy
from pandas import read_csv
from sklearn.decomposition import PCA
# load data

array_outcome = df_new[['outcome']].values
array_outcome[array_outcome!=11] = 1
array_outcome[array_outcome==11] = 0

array_outcome_whole = df1_new[['outcome']].values
array_outcome_whole[array_outcome_whole!=11] = 1
array_outcome_whole[array_outcome_whole==11] = 0

del df1_new['outcome']

X_whole = df1_new.values
Y_whole = array_outcome_whole

del df_new['outcome']

X = df_new.values
Y = array_outcome

print(np.shape(X),np.shape(Y))
```

(494021, 86) (494021, 1)

In [10]:
```python
# #####################
# #########PCA##########
# #####################


# # feature extraction
# pca = PCA(n_components=20)
# fit = pca.fit(X.T)
# # summarize components
# print(("Explained Variance: %s") % fit.explained_variance_ratio_)
# print(np.shape(fit.components_.T))
# X = fit.components_.T
# # print(fit.components_)
```

In [11]:
```python
import keras.backend as K
K.clear_session()
```

In [12]:
```python
# Break into X (predictors) & y (prediction)
# x, y = to_xy(df,'outcome')
# x_train, y_train = to_xy(df,'outcome')
# x_test, y_test = to_xy(df1,'outcome')

# Create a test/train split.  25% test
# Split into train/test
# x_train, x_test, y_train, y_test = train_test_split(
#     X, Y, test_size=0.25, random_state=42)
x_train, x_test, y_train, y_test = train_test_split(
    X, Y, test_size=0.25)
```

In [13]:
```python
# Create neural net
model = Sequential()

# Used relu for activation function
model.add(Dense(10, input_dim=X.shape[1], kernel_initializer='normal', ac
model.add(Dense(50, input_dim=X.shape[1], kernel_initializer='normal', ac
model.add(Dense(10, input_dim=X.shape[1], kernel_initializer='normal', ac

# output layer
model.add(Dense(1, kernel_initializer='normal'))
model.add(Dense(1,activation = 'sigmoid'))
```

In [14]:
```python
# TODO:  Print the model summary
model.summary()
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| dense_1 (Dense) | (None, 10) | 870 |
| dense_2 (Dense) | (None, 50) | 550 |
| dense_3 (Dense) | (None, 10) | 510 |
| dense_4 (Dense) | (None, 1) | 11 |
| dense_5 (Dense) | (None, 1) | 2 |

Total params: 1,943
Trainable params: 1,943
Non-trainable params: 0

In [15]:
```python
class LossHistory(keras.callbacks.Callback):
    def on_train_begin(self, logs={}):
        # TODO:  Create two empty lists, self.loss and self.val_acc
        self.losses = []
        self.accs = []
        self.val_acc = []
        self.val_loss = []
    def on_batch_end(self, batch, logs={}):
        # TODO:  This is called at the end of each batch.
        # Add the loss in logs.get('loss') to the loss list
        loss = logs.get('loss')
        acc = logs.get('acc')
        self.losses.append(loss)
        self.accs.append(acc)
    def on_epoch_end(self, epoch, logs):
        # TODO:  This is called at the end of each epoch.
        # Add the test accuracy in logs.get('loss') to the val_acc list
        val_acc = logs.get('val_acc')
        val_loss = logs.get('val_loss')
        self.val_acc.append(val_acc)
        self.val_loss.append(val_loss)

# Create an instance of the history callback
history_cb = LossHistory()
```

In [16]:
```python
from keras import optimizers

# optimizer defined
opt = optimizers.Adam(lr=0.001, beta_1=0.9, beta_2=0.999, epsilon=1e-08,
# opt = optimizers.Adam(lr=0.001)
```

In [17]:
```python
## compile
model.compile(optimizer=opt,
              loss='binary_crossentropy',
              metrics=['accuracy'])
```

In [18]:
```python
# model.compile(loss='categorical_crossentropy', optimizer='adam', metric
# monitor = EarlyStopping(monitor='val_loss', min_delta=1e-3, patience=5,
```

In [19]:
```python
# model.fit(x,y,validation_data=(x_test,y_test),callbacks=[monitor],verbo
batch_size = 100
epochs = 5
model.fit(x_train,y_train,callbacks=[history_cb],verbose=1,epochs=epochs,
          validation_data=(x_test,y_test))
```

```
Train on 370515 samples, validate on 123506 samples
Epoch 1/5
370515/370515 [==============================] - 15s - loss: 0.0580 -
acc: 0.9841 - val_loss: 0.0239 - val_acc: 0.9913
Epoch 2/5
370515/370515 [==============================] - 15s - loss: 0.0219 -
acc: 0.9946 - val_loss: 0.0195 - val_acc: 0.9949
Epoch 3/5
370515/370515 [==============================] - 15s - loss: 0.0198 -
acc: 0.9950 - val_loss: 0.0189 - val_acc: 0.9949
Epoch 4/5
370515/370515 [==============================] - 15s - loss: 0.0188 -
acc: 0.9950 - val_loss: 0.0174 - val_acc: 0.9958
Epoch 5/5
370515/370515 [==============================] - 14s - loss: 0.0179 -
acc: 0.9954 - val_loss: 0.0180 - val_acc: 0.9947
```
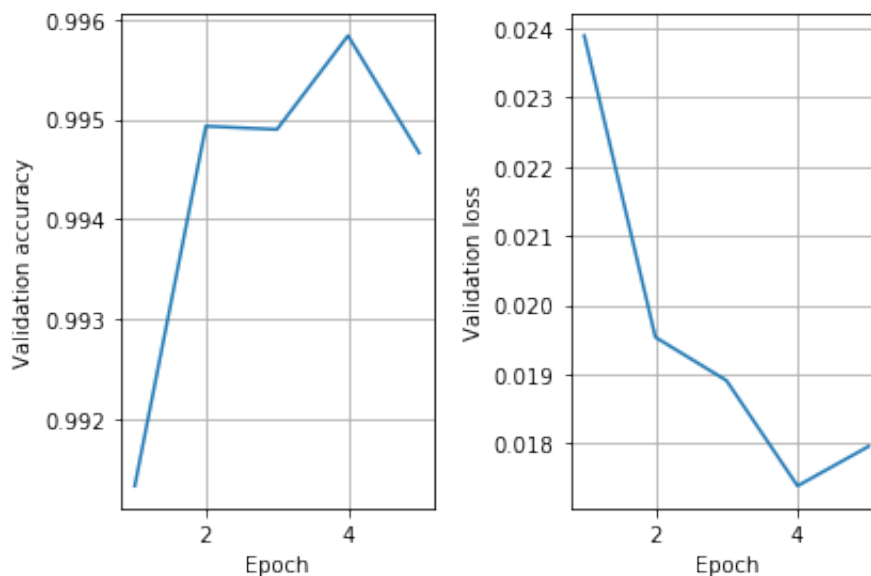
Out[19]: <keras.callbacks.History at 0x19735cdd8>

In [20]:
```python
yhat = model.predict(x_test)
yhat[yhat>0.5] = 1
yhat[yhat<=0.5] = 0
count = 0
yhat_size = np.shape(yhat)[0]

for i in range (yhat_size):
    if yhat[i] != y_test[i]:
        count = count + 1
print("Accuracy is: %s" %(1-count/yhat_size))
print("# of wrongly classified: %s" %(count))
```
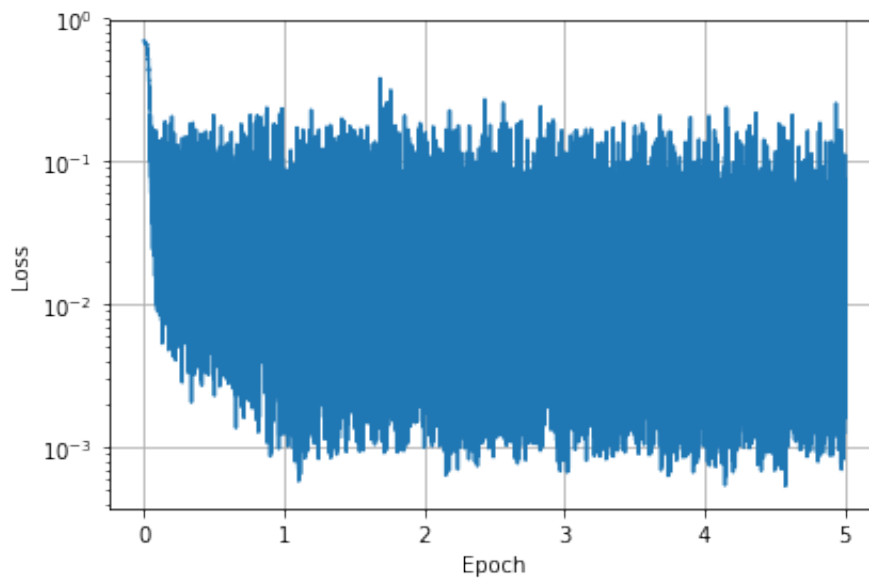
```
Accuracy is: 0.9946642268391819
# of wrongly classified: 659
```

In [21]:
```python
valAcc = history_cb.val_acc
valLoss = history_cb.val_loss
epoch_it = np.arange(1,6)
plt.subplot(121)
plt.plot(epoch_it, valAcc)
plt.grid()
plt.xlabel('Epoch')
plt.ylabel('Validation accuracy')
plt.tight_layout()
plt.subplot(122)
plt.plot(epoch_it, valLoss)
plt.grid()
plt.xlabel('Epoch')
plt.ylabel('Validation loss')
plt.tight_layout()
plt.show()
```
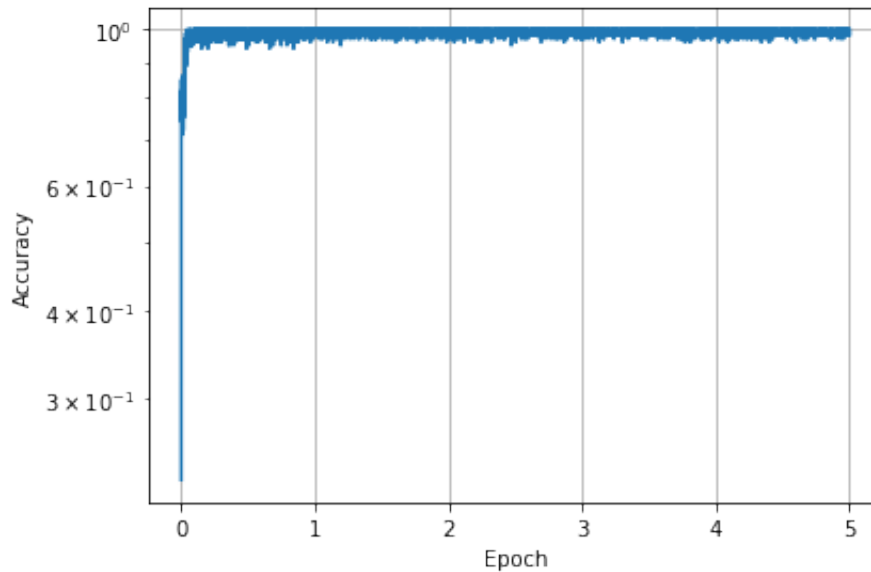
In [22]:
```python
Ls = history_cb.losses
Acc = history_cb.accs
ntr = np.shape(x_train)[0]
epochNum = []
for i in range(18530):
    epochNum.append(i*100/ntr)

matplotlib.pyplot.semilogy(epochNum, Ls)
plt.grid()
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.tight_layout()
plt.show()
```

```
In [23]:  matplotlib.pyplot.semilogy(epochNum, Acc)
          # plt.plot(epochNum, Acc)
          plt.grid()
          plt.xlabel('Epoch')
          plt.ylabel('Accuracy')
          plt.tight_layout()
          # plt.axis([0.3,3,0.96,1])
          plt.show()
```



```
In [24]:  ##################
          #####details######
          ##################
          ##The following is the trained model validated on the whole dataset
```

```
In [25]:  from sklearn.metrics import classification_report,confusion_matrix
```

```
In [26]:  y_pred_whole = model.predict_classes(X_whole,verbose=1)
          #print(y_pred[:5])
```

```
4896480/4898431 [==============================>.] - ETA: 0s
```

```
In [27]: count = 0
         yhat_size = np.shape(y_pred_whole)[0]

         for i in range (yhat_size):
             if y_pred_whole[i] != Y_whole[i]:
                 count = count + 1
         print("Accuracy is: %s" %(1-count/yhat_size))
         print("# of wrongly classified: %s" %(count))
```

```
Accuracy is: 0.9977035503817446
# of wrongly classified: 11249
```

```
In [28]: print(classification_report(Y_whole, y_pred_whole))
```

```
             precision    recall  f1-score   support

          0       1.00      0.99      0.99    972781
          1       1.00      1.00      1.00   3925650

avg / total       1.00      1.00      1.00   4898431
```

```
In [29]: Cm = confusion_matrix(Y_whole,y_pred_whole)
         C = np.sum(Cm)
         Cm = Cm/C
         print('Confusion Matrix:')
         print(np.array_str(Cm, precision=4, suppress_small=True))
```

```
Confusion Matrix:
[[ 0.1969  0.0017]
 [ 0.0006  0.8008]]
```