# Lab 5: Pitch Detection in Audio

In this lab, we will use numerical optimization to find the pitch and harmonics in a simple audio signal. In addition to the concepts in the gradient descent demo (./grad_descent.ipynb), you will learn to:

- Load, visualize and play audio recordings
- Divide audio data into frames
- Perform nested minimization

The ML method presented here for pitch detection is actually not a very good one. As we will see, it is highly suseptible to local minima and quite slow. There are several better pitch detection algorithms (https://en.wikipedia.org/wiki/Pitch_detection_algorithm), mostly using frequency-domain techniques. But, the method here will illustrate non-linear estimation well.

## Reading the Audio File

Python provides a very simple method to read a `wav` file in the `scipy.io.wavefile` package. We first load that along with the other packages.

```
In [66]:  from scipy.io.wavfile import read
          import numpy as np
          import matplotlib.pyplot as plt
          %matplotlib inline
```

In the github repository, you should find a file, viola.wav (./viola.wav). Download this file to your local directory. You can then read the file with the `read` command. Print the sample rate in Hz, the number of samples in the file and the file length in seconds.

```
In [2]:  # Read the file
         sr, y = read('viola.wav')

         # TODO:  Print sample rate and file length
         nsamp = len(y)
         print("Sample rate   = %f Hz" % sr)
         print("Number sample = %d" % nsamp)
         print("File length   = %f sec" % (nsamp/sr))
```

```
Sample rate   = 44100.000000 Hz
Number sample = 299350
File length   = 6.787982 sec
```

You can then play the file with the following command. You should hear the viola play a sequence of simple notes.

In [3]:
```
import IPython.display as ipd
ipd.Audio(y, rate=sr) # Load a NumPy array
```

Out[3]:
▶  0:00 / 0:06  ●——————  🔊  —————●  ⬇

For the analysis below, it will be easier to re-scale the samples so that they have an average squared value of 1. Find the `scale` value in the code below to do this.

In [67]:
```
# TODO
scale = np.sqrt(np.mean(y**2))
y = y / scale
```

# Dividing the Audio File into Frames

In audio processing, it is common to divide audio streams into short frames (typically between 10 to 40 ms long). Since frames are often processed with an FFT, the frames are typically a power of two. Analysis is then performed in the frames separately. Given the vector y, create a `nfft x nframe` matrix `yframe` where

```
yframe[:,0] = samples y[k], k=0,...,nfft-1
yframe[:,1] = samples y[k], k=nfft,...,2*nfft-1,
yframe[:,2] = samples y[k], k=2*nfft,...,3*nfft-1,
...
```

You can do this with the `reshape` command with `order=F`. Zero pad y if the number of samples of y is not divisible by `nfft`. Print the total number of frames as well as the length (in milliseconds) of each frame.

Note that in actual audio processing, the frames are typically overlapping and use careful windowing. But, we will ignore that here for simplicity.

In [68]:
```python
# Frame size
nfft = 1024

# TODO:
#   nframe = ...
#   yframe = ...
ny = len(y)
if (ny % nfft > 0):
    npad = nfft - (ny%nfft)
    y1 = np.hstack((y, np.zeros(npad)))
else:
    y1 = y

# Divide into frames
nframe = len(y1) // nfft
yframe = y1.reshape((nfft,nframe),order='F')

# Print dimensions
print("Frame length  = %.2f ms" % (1000*nfft/sr))
print("Number frames = %d" % nframe)
```
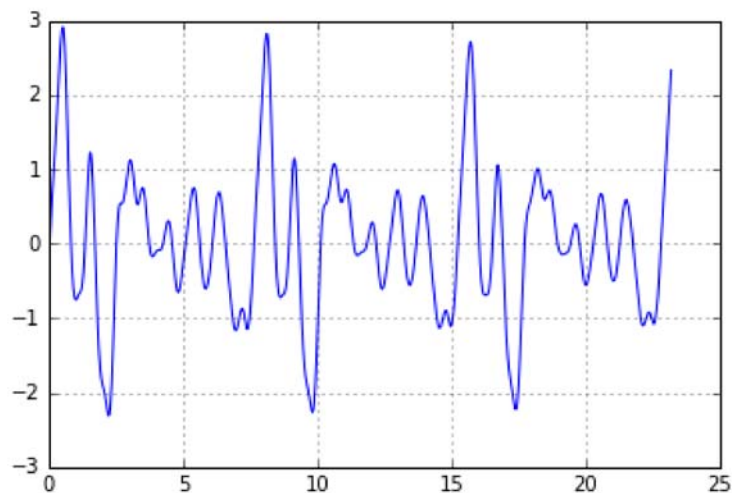
```
Frame length  = 23.22 ms
Number frames = 293
```

Let i0=10 and set yi=yframe[:,i0] be the samples of frame i0. We will use this frame for most of the rest of the lab. Plot the samples of yi. Label the time axis in milliseconds (ms).

In [69]:
```python
# Get samples from frame 10
i0 = 10
yi = yframe[:,i0]

# TODO:  Plot yi vs. time (in ms)
t = np.arange(nfft)/sr*1000
plt.plot(t,yi)
plt.grid()
```

# Fitting a Multi-Sinusoid

A common model for audio samples, `yi[k]`, containing an instrument playing a single note is the multi-sinusoid model:

```
yi[k] \approx yhati[k] = c + \sum_{j=0}^{nterms-1} a[j]*cos(2*np.pi*k*freq0*(j+1)/s
r) +
                                      b[j]*sin(2*np.pi*k*freq0*(j+1)/s
r)
```

where `sr` is the sample rate. The parameter `freq0` is called the fundamental frequency and the audio signal is modeled as being composed of sinusoids and cosinusoids with frequencies equal to integer multiples of the fundamental. In audio processing, these terms are called *harmonics*. In analyzing audio signals, a common goal is to determine both the fundamental frequency `freq0` (the pitch of the audio) as well as the coefficients of the harmonics,

```
beta = (c, a[0], ..., a[nterms-1], b[0], ..., b[nterms-1]).
```

To find the parameters, we will fit the mean squared error loss function:

```
mse(freq0,beta) := 1/N * \sum_k (yi[k] - yhati[k])**2,   N = len(yi).
```

In practice, a separate model would be fit for each audio frame. But, in this lab, we will mostly look at a single frame.

## Nested Minimization

We will perform the minimization of `mse` in a nested manner: First, given a fundamental frequency `freq0`, we minimize over the coefficients `beta`. Call this minimum `mse1`:

```
mse1(freq0) := min_beta mse(freq0,beta)
```

Importantly, this minimizaiton can be performed by least-squares. Then, we find the fundamental frequency `freq0` by minimizing `mse1`:

```
min_{freq0} mse1(freq0)
```

We will use gradient-descent minimization with `mse1(freq0)` as the objective function. This form of *nested* minimization is commonly used whenever we can minimize over one set of parameters easily given the other.

# Setting Up the Objective Function

We will use the class `AudioFitFn` below to perform the two-part minimization. Complete the `feval` method in the class. The method should take the argument `freq0` and perform the minimization of the MSE over `beta`. Specifically, fill the code in `feval` to perform the following:

- Construct a matrix, `A` such that `yhati = A*beta`.
- Find `betahat` with the `np.linalg.lstsq()` method using the matrix `A` and the samples `self.yi`. This is simpler than constructing a linear regression object.
- Compute and store the estimate `self.yhati = A.dot(betahat)`.
- Compute the `mse1`, the minimum MSE, by comparing `self.yhati` and `self.yi`.
- For now, set the gradient to `mse1_grad=0`. We will fill this part in later.
- Return `mse1` and `mse1_grad`.

In [33]:
```python
class AudioFitFn(object):
    def __init__(self,yi,sr=44100,nterms=8):
        """
        A class for fitting

        yi:   One frame of audio
        sr:   Sample rate (in Hz)
        nterms:  Number of harmonics used in the model (default=8)
        """
        self.yi = yi
        self.sr = sr
        self.nterms = nterms

    def feval(self,freq0):
        """
        Optimization function for audio fitting.  Given a fundamental frequenc
y, freq0, the
        method performs a least squares fit for the audio sample using the mod
el:

        yhati[k] = c + \sum_{j=0}^{nterms-1} a[j]*cos(2*np.pi*k*freq0*(j+1)/s
r)
                                            + b[j]*sin(2*np.pi*k*freq0*(j+1)/s
r)

        The coefficients beta = [c,a[0],...,a[nterms-1],b[0],...,b[nterms-1]]
        are found by least squares.

        Returns:

        mse1:   The MSE of the best least square fit.
        mse1_grad:  The gradient of mse1 wrt to the parameter freq0
        """

        # TODO

        # Get time range and phase shift for the fundamental frequency
        nsamp = self.yi.shape[0]
```

```python
        t = np.arange(0,nsamp)/self.sr
        ph = 2*np.pi*t*freq0

        # Construct the regression matrix
        A = np.empty((nsamp,2*self.nterms+1))
        A[:,0]=1
        for i in range(self.nterms):
            A[:,2*i+1] = np.cos(ph*(i+1))
            A[:,2*i+2] = np.sin(ph*(i+1))

        # Construct the gradient of the data matrix
        Ag = np.zeros((nsamp,2*self.nterms+1))
        for i in range(self.nterms):
            Ag[:,2*i+1] = -(i+1)*2*np.pi*t*np.sin(ph*(i+1))
            Ag[:,2*i+2] = (i+1)*2*np.pi*t*np.cos(ph*(i+1))


        # Perform a LS fit and stores the fit in self.yhati
        betahat = np.linalg.lstsq(A,self.yi)[0]
        self.yhati = A.dot(betahat)

        # Compute the RSS per sample
        mse1 = np.mean((self.yi-self.yhati)**2)

        # Compute the gradient wrt to freq0
        mse1_grad = np.array([2*(self.yhati-self.yi).dot(Ag.dot(betahat))])/ns
amp

        return mse1, mse1_grad
```
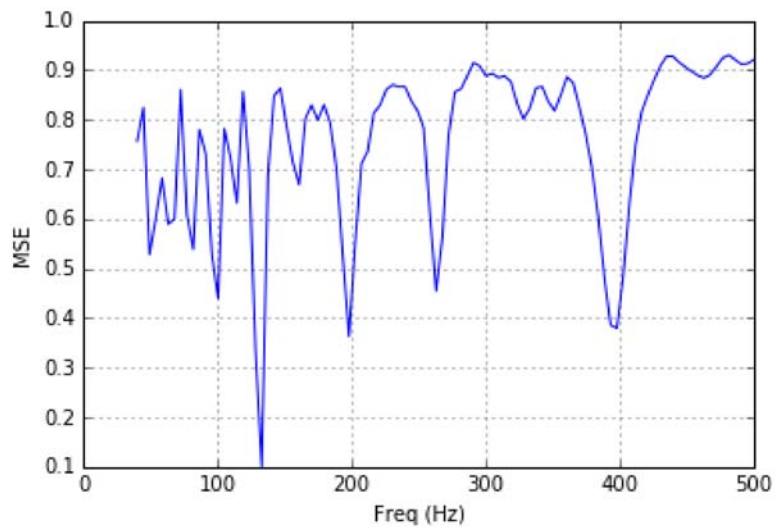
Instatiate an object, `audio_fn` from the class `AudioFitFn` with the samples `yi`. Then, using the `feval` method, compute and plot `mse1` for 100 values `freq0` in the range of 40 to 500 Hz. You should see a minimum around `freq0 = 130` Hz, but there are several other local minima.

In [34]:
```python
# TODO
audio_fn = AudioFitFn(yi)
freq0_test = np.linspace(40,500,100)
ntest = len(freq0_test)
mse1 = np.zeros(ntest)
for i, freq0 in enumerate(freq0_test):
    mse1[i], _ = audio_fn.feval(freq0)

plt.plot(freq0_test, mse1)
plt.grid()
plt.ylabel('MSE')
plt.xlabel('Freq (Hz)')
```

Out[34]: <matplotlib.text.Text at 0x1ca4b599d68>



Print the value of `freq0` that achieves the minimum `mse1`. Also, plot the estimated function `audio_fn.yhati` for that along with the original samples `yi`.
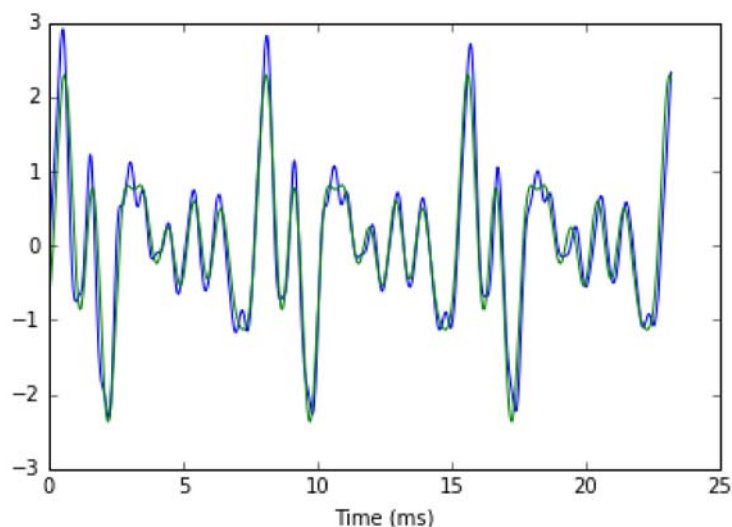
```
In [35]:  # TODO
          # Get the estimated frequency
          im = np.argmin(mse1)
          freq0_opt = freq0_test[im]
          print("Frequency estimate = %.2f Hz" % freq0_opt)

          # Get the estimated function
          audio_fn.feval(freq0_opt)
          yhati = audio_fn.yhati

          plt.plot(t, yi)
          plt.plot(t, yhati)
          plt.xlabel('Time (ms)')
```

Frequency estimate = 132.93 Hz

Out[35]:  <matplotlib.text.Text at 0x1ca4bf4b5c0>



# Computing the Gradient

The above method found the estimate for `freq0` by performing a search over 100 different frequency values and selecting the frequency value with the lowest MSE. We now see if we can estimate the frequency with gradient descent minimization of the MSE. We first need to modify the `feval` method in the `AudioFitFn` class above to compute the gradient. Some elementary calculus (see the homework), shows that

```
dmse1(freq0)/dfreq0 = dmse(freq0,betahat)/dfreq0
```

So, we just need to evaluate the partial derivative of `mse = np.mean((yi-yhati)**2)` with respect to the parameter `freq0` holding the parameters `beta=betahat`. Modify the `feval` method above to compute the gradient and return the gradient in `mse1_grad`.

Then, test the gradient by taking two close values of `freq0`, say `freq0_0` and `freq0_1` and verifying that first-order approximation holds.

In [37]:
```python
# TODO
# Construct the audio fit object
i0 = 10
yi = yframe[:,i0]
audio_fn = AudioFitFn(yi,nterms=8)

# Two values of freq0
freq0_0 = 100
step=1e-6
freq0_1 = freq0_0 + step*np.random.randn(1)

# Compute the function and gradient at each point
mse1_0, mse1_grad0 = audio_fn.feval(freq0_0)
mse1_1, mse1_grad1 = audio_fn.feval(freq0_1)
dmse1_act = mse1_1-mse1_0
dmse1_est = mse1_grad0.dot(freq0_1-freq0_0)
print("Actual difference    %12.4e" % dmse1_act)
print("Predicted difference %12.4e" % dmse1_est)
```

```
Actual difference       -5.8694e-08
Predicted difference    -5.8694e-08
```

# Run the Optimizer

We cut and paste the optimizer from the gradient descent demo (./grad_descent.ipynb).

In [39]:
```python
def grad_opt_adapt(feval, winit, nit=1000, lr_init=1e-3):
    """
    Gradient descent optimization with adaptive step size

    feval:  A function that returns f, fgrad, the objective
            function and its gradient
    winit:  Initial estimate
    nit:    Number of iterations
    lr:     Initial learning rate

    Returns:
    w:   Final estimate for the optimal
    f0:  Function at the optimal
    """

    # Set initial point
    w0 = winit
    f0, fgrad0 = feval(w0)
    lr = lr_init

    # Create history dictionary for tracking progress per iteration.
    # This isn't necessary if you just want the final answer, but it
    # is useful for debugging
    hist = {'lr': [], 'w': [], 'f': []}

    for it in range(nit):
```

```
            # Take a gradient step
            w1 = w0 - lr*fgrad0

            # Evaluate the test point by computing the objective function, f1,
            # at the test point and the predicted decrease, df_est
            f1, fgrad1 = feval(w1)
            df_est = fgrad0.dot(w1-w0)

            # Check if test point passes the Armijo rule
            alpha = 0.5
            if (f1-f0 < alpha*df_est) and (f1 < f0):
                # If descent is sufficient, accept the point and increase the
                # learning rate
                lr = lr*2
                f0 = f1
                fgrad0 = fgrad1
                w0 = w1
            else:
                # Otherwise, decrease the learning rate
                lr = lr/2

            # Save history
            hist['f'].append(f0)
            hist['lr'].append(lr)
            hist['w'].append(w0)

        # Convert to numpy arrays
        for elem in ('f', 'lr', 'w'):
            hist[elem] = np.array(hist[elem])
        return w0, f0, hist
```

Now, run the optimizer with the feval function with a starting estimate for freq0 = 130 Hz. Use `lr_init=1e-3` and `f0_init=130`. Print the final frequency estimate. Also, print the midi number (https://newt.phys.unsw.edu.au/jw/notes.html) of the estimated frequency:

```
    midi_num = 12*log2(freq/440 Hz) + 69
```

If the note was exactly a musical note, `midi_num` should be an integer. But you will see that the frequency does not exactly lie on a note since the pitch in a viola bends around the note.