

Numerical Recipes: Hand-In 4

Rens Kievit (s1948415)

May 26, 2023

Abstract

In this report we present the problems, solutions and scripts for the exercises from the fourth handout for the course Numerical Recipes.

Plotting styles in this report are set using the following code

```
1 def set_styles():
2     """For consistent plotting scheme"""
3     plt.style.use('default')
4     mpl.rcParams['axes.grid'] = True
5     plt.style.use('seaborn-darkgrid')
6     mpl.rcParams['font.family'] = 'serif'
7     mpl.rcParams['lines.linewidth'] = 1.5
8     mpl.rcParams['font.size'] = 14
```

Listing 1: Matplotlib Plotting Styles

1 Solar System Orbits

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from astropy.time import Time
4 from astropy.coordinates import solar_system_ephemeris
5 from astropy.coordinates import get_body_barycentric_posvel
6 from astropy import units as u
7 from astropy import constants as const
8 from plotting import set_styles
9
10 from algorithms import LeapFrog, RungeKutta4
11
12 # TODO update gravity
13 def grav_force(r, Rsun, Msun=const.M_sun.value,
14               G=const.G.to_value((u.AU**3)*(u.kg**(-1))*(u.d**(-2)
15               ))):
16     """Compute the N-D gravitational force acting on an object with
17     mass m1
18     from an object with mass m2 located at the origin"""
19     r_diff = r - Rsun
20     return - (G * Msun * r_diff) / (np.power(np.linalg.norm(r_diff)
21     , 3))
22
23     #return (G * Msun / (r_diff*r_diff)) * (r_diff/np.sum(np.sqrt(
24     r_diff*r_diff)))
25
26 def get_pos_vel(object_names,
27                 time="2021-12-07 10:00",
28                 plot=False):
29     """Get the positions in AU and velocities in AU/d of all solar
30     system objects at time."""
31     t = Time(time)
32     planets = {} # Dictionary in which we save our planet objects (
33     and the sun)
34
35     if plot:
36         fig, [ax1, ax2] = plt.subplots(1,2,figsize=(12,6),
37         tight_layout=True)
38
39     for obj in object_names:
40         # get the astropy data
41         with solar_system_ephemeris.set('jpl'):
42             obj_pos_vel = get_body_barycentric_posvel(obj, t)
43
44             x0 = obj_pos_vel[0].xyz.to_value(u.AU)
45             v0 = obj_pos_vel[1].xyz.to_value(u.AU/u.d)
46
47             if obj == 'sun':
48                 r_sun = x0 # Position of the sun, which we take to
49                 remain constant
50
51             # Feed r_sun into the grav_force function to get the
52             correct coordinates
53             acc_func = lambda x, v: grav_force(x, r_sun)
54             planets[obj] = LeapFrog(x0, v0, acc_func)
55
56     if plot:
57         if obj == 'sun':
58             m = '*'
59         else:
60             m = 'o'
```

```

52         ax1.scatter(x0[0], x0[1], marker=m)
53         ax2.scatter(x0[0], x0[2], label=obj.capitalize(),
marker=m)
54
55     if plot:
56         ax1.set_xlabel(r'$X$ [AU]')
57         ax1.set_ylabel(r'$Y$ [AU]')
58         ax2.set_xlabel(r'$X$ [AU]')
59         ax2.set_ylabel(r'$Z$ [AU]')
60         plt.suptitle('Initial Star and Planet Positions')
61
62         plt.legend()
63         plt.savefig('results/initial_positions.png', bbox_inches='
tight')
64
65     return planets
66
67 def make_orbits(planets, N, h, plot=False,
compare_to_rk4=False):
68     """
69     if plot:
70         fig1, axs1 = plt.subplots(1,2,figsize=(15,6))
71         fig2, ax2 = plt.subplots(1,1,figsize=(12,6))
72         fig3, ax3 = plt.subplots(1,1,figsize=(12,6))
73         fig4, axs4 = plt.subplots(1,2,figsize=(15,6))
74         i = 1
75
76     for name, obj in planets.items():
77         if name == 'sun':
78             # Just place it in all corresponding plots
79             axs1[0].scatter(obj.x0[0], obj.x0[1], color='yellow', s
=50, marker='*', label='Sun')
80             axs1[1].scatter(obj.x0[0], obj.x0[1], color='yellow', s
=50, marker='*')
81             axs4[0].scatter(obj.x0[0], obj.x0[1], color='yellow', s
=50, marker='*', label='Sun')
82             axs4[1].scatter(obj.x0[0], obj.x0[1], color='yellow', s
=50, marker='*')
83             continue # do not simulate the sun to itself
84
85         print(name.capitalize())
86         # Apply Leapfrog
87         obj.simulate_motion(N, h)
88
89         if plot:
90             axs1[0].scatter(obj.x0[0], obj.x0[1], c=f'C{i}', label=
name.capitalize())
91             axs1[1].scatter(obj.x0[0], obj.x0[1], c=f'C{i}')
92             axs1[0].plot(obj.x[:,0], obj.x[:,1], c=f'C{i}')
93             axs1[1].plot(obj.x[:,0], obj.x[:,1], c=f'C{i}')
94             ax2.plot(h*np.arange(obj.x.shape[0]), obj.x[:,2], c=f'C
{i}')
95
96
97
98     # Code for the bonus question
99     if compare_to_rk4:
100         RK4 = RungeKutta4(obj.x0, obj.v0, obj.acc_func)
101         RK4.simulate_motion(N, h)
102
103         # Plot the difference in x-positions between RK4 and LF
104         ax3.plot(h*np.arange(obj.x.shape[0]), (RK4.x[:,0] - obj
.x[:,0]), c=f'C{i}')

```

```

105         # Plot the orbits
106         axs4[0].scatter(RK4.x0[0], RK4.x0[1], c=f'C{i}', label=
107         name.capitalize())
108         axs4[1].scatter(RK4.x0[0], RK4.x0[1], c=f'C{i}')
109         axs4[0].plot(RK4.x[:,0], RK4.x[:,1], c=f'C{i}')
110         axs4[1].plot(RK4.x[:,0], RK4.x[:,1], c=f'C{i}')
111
112         i += 1
113
114     if plot:
115         # FIGURE 1
116         for ax1 in axs1:
117             ax1.set_xlabel(r'$X$ [AU]')
118             ax1.set_ylabel(r'$Y$ [AU]')
119             fig1.suptitle('Star and Planet Orbits over 200 Years with
LeapFrog')
120
121             # Right plot zoomed in on the rocky planets
122             axs1[1].set_xlim(-2, 2)
123             axs1[1].set_ylim(-2, 2)
124
125             # Place legend next to figures
126             fig1.subplots_adjust(right=0.85)
127             fig1.legend(loc='right')
128             fig1.savefig('results/orbits_1f.png', bbox_inches='tight')
129
130         # FIGURE 2
131         ax2.set_xlabel(r'Time in Days')
132         ax2.set_ylabel(r'$Z$ [AU]')
133         ax2.set_title('Planet z-positions over 200 Years with
LeapFrog')
134         fig2.savefig('results/zplane.png', bbox_inches='tight')
135
136         # FIGURE 3
137         ax3.set_xlabel(r'Time in Days')
138         ax3.set_ylabel(r'$X_{\mathrm{RK4}} - X_{\mathrm{LF}}$ [AU]')
139         #ax3.set_yscale('log')
140         ax3.set_title('Positional Differences')
141         fig3.savefig('results/rk4_1f_diff.png', bbox_inches='tight')
142
143         # FIGURE 4
144         for ax4 in axs4:
145             ax4.set_xlabel(r'$X$ [AU]')
146             ax4.set_ylabel(r'$Y$ [AU]')
147             fig4.suptitle('Star and Planet Orbits over 200 Years with
Runge-Kutta 4')
148
149             # Right plot zoomed in on the rocky planets
150             axs4[1].set_xlim(-2, 2)
151             axs4[1].set_ylim(-2, 2)
152
153             # Place legend next to figures
154             fig4.subplots_adjust(right=0.85)
155             fig4.legend(loc='right')
156             fig4.savefig('results/orbits_rk4.png', bbox_inches='tight')
157
158 def solar_system_sim():
159     set_styles()

```

```

160 object_names = ['sun', 'mercury', 'venus', 'earth', 'mars', '
    jupiter', 'saturn', 'uranus', 'neptune']
161
162 h = 0.5
163 N = int((365.25 * 200)/h) # 200yrs in total
164 print('N Steps = ', N)
165
166 plot = True
167 compare_to_rk4 = True
168 #####3
169
170 planets = get_pos_vel(object_names, plot=plot)
171 make_orbits(planets, N, h, plot=plot, compare_to_rk4=
    compare_to_rk4)
172
173 def main():
174     solar_system_sim()
175
176 if __name__ == '__main__':
177     main()

```

Listing 2: Code to collect the positions and velocities of solar system objects and apply the LeapFrog and Runge-Kutta 4 Algorithms.

In this section we will simulate the orbits of the planets in our Solar system. Here the forces between planets are negligible compared to the force between each planet and the sun. The orbit of each planet is therefore well defined by a single equation, which we can easily brute force with an ODE solver. The acceleration of the planet is given by the gravitational force, which is defined as

$$\vec{F}_G = \frac{-mMG}{||\vec{r}||^3} \vec{r}. \quad (1)$$

Where $M = 1.989 \times 10^{30}$ kg is the mass of the sun, m is the mass of the planet, $G = 6.6743 \times 10^{-11}$ m³ kg⁻¹ s⁻² is the gravitational constant. \vec{r} is the vector pointing from the planet to the sun. The acceleration is then simply given through $F = ma$, so the planet mass cancels out.

We get the initial positions and velocities of all 8 solar system planets and the sun using **astropy** following the procedure given on the handin at 10:00 on 2021-12-07. We show these initial positions in Figure 1.

1.1 Leapfrog

```

1 class LeapFrog():
2     """Class to store positions and velocities for any physical
    object
3     whose motion we want to simulate using leapfrog"""
4     def __init__(self, x0, v0, acc_func):
5         self.x0 = x0
6         self.v0 = v0
7         self.acc_func = acc_func
8         self.h = None # Only works for constant h at this moment
9         self.dim = x0.shape[0]
10
11     def setup_sim(self, N, h, continue_from_last=False):
12         """Simulate the motion under accelatrion func for N
    timesteps"""
13         if (self.h is not None) and (self.h != h):

```

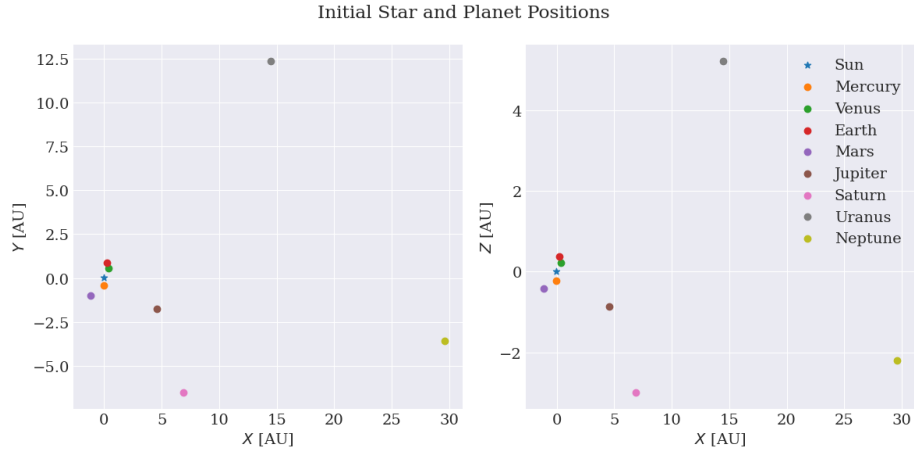


Figure 1: Initial positions of all 8 solar system planets and the sun at 10:00 on 2021-12-07 in the x-y plane (*left*) and the x-z plane (*right*)

```

14         raise ValueError(f'Provided h {h} does not match
existing h {self.h}')
15     self.h = h
16     if continue_from_last:
17         self.index = self.x.shape[0]
18         self.x = np.vstack([self.x, np.zeros([N, self.dim])])
19         self.v = np.vstack([self.v, np.zeros([2*N, self.dim])])
20         self.a = np.vstack([self.a, np.zeros([N, self.dim])])
21     else:
22         self.x = np.zeros([N, self.dim])
23         self.v = np.zeros([2*N, self.dim]) # Two times so we
can store v_i and v_i+0.5
24         self.a = np.zeros([N, self.dim])
25
26         self.x[0] = self.x0
27         self.v[0] = self.v0
28         self.v[1] = self._kickstart() # Get v at v_i+1/2 using
RK4
29
30         self.index = 1
31
32     def _kickstart(self):
33         """Given x0 and v0 and a function for acceleration, get the
velocity
34         at v_i+1/2. This is done standardized for RK4, but Euler
also works"""
35         h_rk4 = 0.5*self.h # take a half step
36         # Approximated positions after 0.25 and 0.5 of a step using
v0
37         x_half = self.x[0] + h_rk4 * self.v[0]
38         x_full = self.x[0] + h_rk4 * self.v[0]
39
40         # Find v_1/2 using RK4 and approximated positions
41         k1 = h_rk4 * self.acc.func(self.x[0], self.v[0])
42         k2 = h_rk4 * self.acc.func(x_half, self.v[0] + 0.5*k1)
43         k3 = h_rk4 * self.acc.func(x_half, self.v[0] + 0.5*k2)
44         k4 = h_rk4 * self.acc.func(x_full, self.v[0] + k3)

```

```

45         one_sixth = 1./6.
46         one_third = 1./3.
47         return self.v[0] + one_sixth*k1 + one_third*k2 + one_third*
k3 + one_sixth*k4
48
49     def simulate_motion(self, N, h, continue_from_last=False):
50         """Simulate motion under the given acceleration fuction"""
51         self._setup_sim(N, h, continue_from_last)
52
53         # Leapfrog algorithm
54         for i in range(self.index, self.index + N - 1):
55             # index in the velocity array. Offset because we leave
space for v_i
56             v_idx = (2*i) + 1
57             # currently can't use a = a(x, v) because we don't have
v_i
58             self.a[i] = self.acc_func(self.x[i-1], self.v[i-2])
59             self.v[v_idx] = self.v[v_idx-2] + self.h * self.a[i]
60             self.x[i] = self.x[i-1] + self.h * self.v[v_idx]

```

Listing 3: LeapFrog Class

We start by applying the leapfrog algorithm to simulate the orbits. This is the proper method of solving ordinary differential equations for physical systems, because it is reversible and therefore able to conserve energy. This is an important factor in making sure the solutions do not diverge to unphysical states. We have implemented the algorithm in a class, given above. The leapfrog algorithm works by kicking the object from positions x_i with velocity $v_{i+1/2}$. However, we are only given x_0 and v_0 . This means that we first have to use a regular ODE solver to find $v_{1/2}$. Here, this is done in the `_kickstart` function with a single iteration of Runge-Kutta 4.

We integrate the orbits with a step size $h = 0.5$ days for 200 years, this gives $N = 146100$ steps in total. We show the resulting orbits in Figure 2.

The orbits appear to behave very nicely. As mentioned previously, this algorithm conserves energy. This is visible in the fact that the orbits in the top plot of Figure 2 appear to be well defined lines that do not diverge away from or towards the Sun, which would be a violation of energy conservation. This same fact is visible in the bottom plot, the planets all appear to oscilate up and down in the z-plane, but the amplitude of the oscilations is constant. This also indicates the fact that energy conservation is not violated. Finally, this simulation shows that the orbit of Mercury is not entirely constant, but rather shifts around a little bit. This is not per se an inconsistent result, because its orbit remains closed.

1.2 Runge-Kutta 4

```

1 class RungeKutta4():
2     """Class to store positions and velocities for any physical
object
3     whose motion we want to simulate using leapfrog"""
4     def __init__(self, x0, v0, acc_func):
5         self.x0 = x0
6         self.v0 = v0
7         self.acc_func = acc_func
8         self.h = None # Only works for constant h at this moment
9         self.dim = x0.shape[0]

```

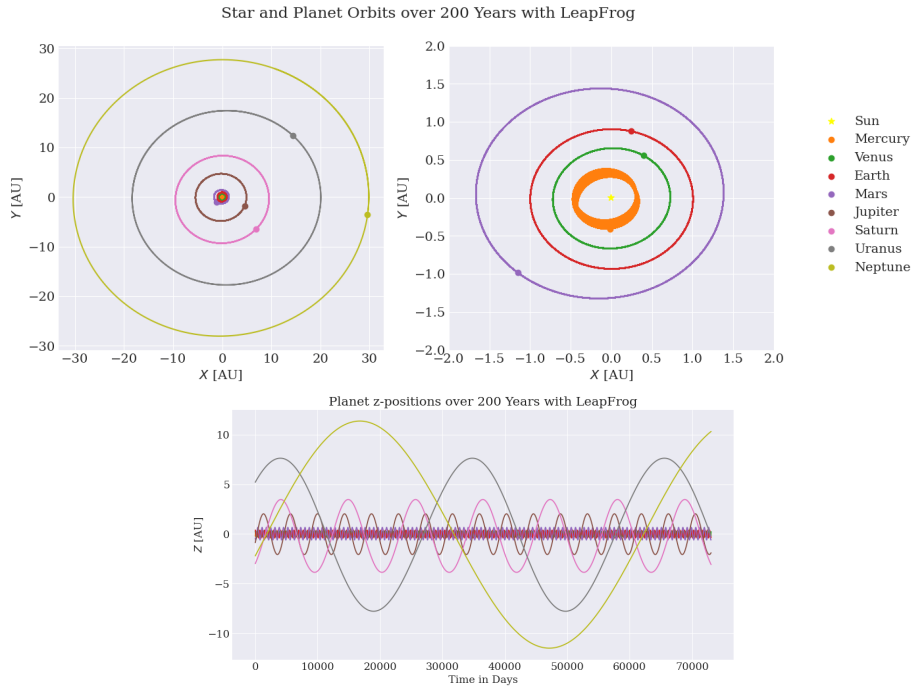


Figure 2: Integrated orbits of the 8 planets in the solar system over 200 years with a step size of 0.5 days using the LeapFrog algorithm and an initial kick with Runge-Kutta 4. *top*: Planetary orbits in the x-y plane, the plot in the right is zoomed in to better highlight the orbits of the rocky planets. *bottom*: Planetary positions on the z-axis as a function of time.


```

10
11 def _setup_sim(self, N, h, continue_from_last=False):
12     """Simulate the motion under accelatrion funcn for N
13     timesteps"""
14     if (self.h is not None) and (self.h != h):
15         raise ValueError(f'Provided h {h} does not match
16         existing h {self.h}')
17     self.h = h
18     if continue_from_last:
19         self.index = self.x.shape[0]
20         self.x = np.vstack([self.x, np.zeros([N, self.dim])])
21         self.v = np.vstack([self.v, np.zeros([N, self.dim])])
22         self.a = np.vstack([self.a, np.zeros([N, self.dim])])
23     else:
24         self.x = np.zeros([N, self.dim])
25         self.v = np.zeros([N, self.dim])
26         self.a = np.zeros([N, self.dim])
27
28         self.x[0] = self.x0
29         self.v[0] = self.v0
30         self.index = 0
31
32 def simulate_motion(self, N, h, continue_from_last=False):
33     """Simulate motion under the given acceleration fucntion"""
34     self._setup_sim(N, h, continue_from_last)
35
36     one_third = 1./3.
37     one_sixth = 1./6.
38
39     # Leapfrog algorithm
40     for i in range(self.index, self.index + N - 1):
41         k1v = self.h * self.acc_func(self.x[i], self.v[i])
42         k1x = self.h * self.v[i]
43
44         k2v = self.h * self.acc_func(self.x[i]+0.5*k1x, self.v[
45 i]+0.5*k1v)
46         k2x = self.h * (self.v[i] + 0.5*k1v)
47
48         k3v = self.h * self.acc_func(self.x[i]+0.5*k2x, self.v[
49 i]+0.5*k2v)
50         k3x = self.h * (self.v[i] + 0.5*k2v)
51
52         k4v = self.h * self.acc_func(self.x[i]+k3x, self.v[i]+
53 k3v)
54         k4x = self.h * (self.v[i] + k3v)
55
56         self.v[i+1] = self.v[i] + one_sixth*k1v + one_third*k2v
57         + one_third*k3v + one_sixth*k4v
58         self.x[i+1] = self.x[i] + one_sixth*k1x + one_third*k2x
59         + one_third*k3x + one_sixth*k4x

```

Listing 4: Runge-Kutta 4 Class

As an experiment, we also integrate the orbits of the planets given the same initial conditions with a normal ODE solver that does not have the power to preserve energy. Similar to the LeapFrog method above, we also implement the RK-4 method using a class given above. However, this time we do not have to provide the velocities with an initial kick. We use the same step size and number of steps as before, and plot the resulting orbits in Figure 3

Visually, the orbits of all planets except Mercury appear exactly the same. The biggest difference lies in the fact that Mercury's orbit remains constant

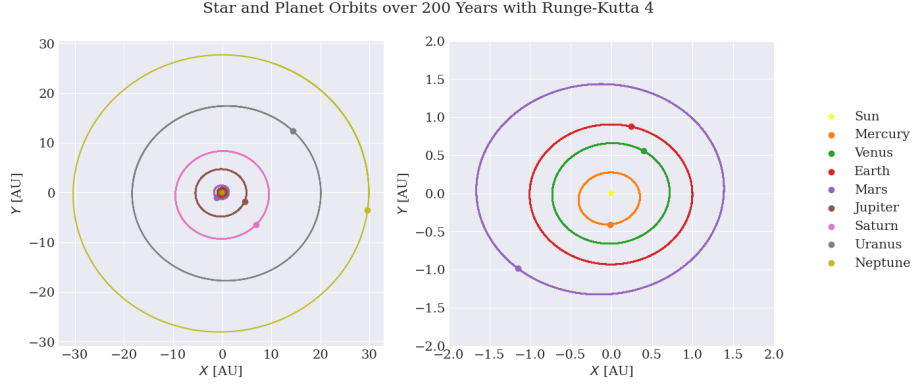


Figure 3: Integrated orbits in the x-y plane of the 8 planets in the solar system over 200 years with a step size of 0.5 days using the Runge-Kutta 4 method.

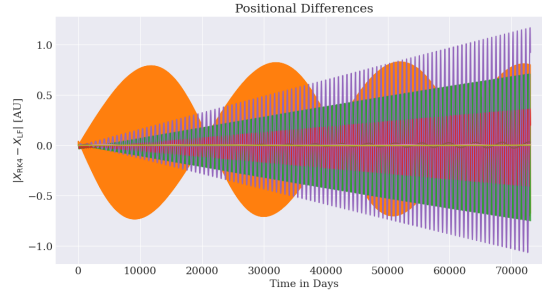


Figure 4: Difference in x-position between the orbit as integrated using the Runge-Kutta 4 method and the LeapFrog algorithm plotted as a function of time

over 200 years according to RK4, while it shifted around according to LeapFrog. However comparing two plots with each other visually is not the best method of investigating the performance of an algorithm. Therefore we present the difference in x-position between the two integrated orbits in Figure 4.

In that figure we see that the difference in x-position shows a somewhat periodic pattern that increases in amplitude over time. This indicates that the two solutions are slowly diverging from each other. The pattern we see corresponds to what we might have expected. Because RK-4 is a forward method, the velocity it predicts will always push the object slightly more away from the star than the true value. This small error accumulates over time leading to a larger orbit. This can physically be interpreted as increasing the amount of energy in the system, which is a violation of the laws of physics that LeapFrog was able to adhere to. Additionally, we can clearly see that there is a large discrepancy between the two orbits of Mercury that occurs from a very early point.

2 Computing Forces with FFT

```

1 #!/usr/bin/env python3
2 def get_densities():
3     """Code downloaded from Zorrry Belcheva's personal strw page.
4     Turned into a function for use in fourier.py"""
5     import numpy as np
6     np.random.seed(121)
7
8     n_mesh = 16
9     n_part = 1024
10    positions = np.random.uniform(low=0, high=n_mesh, size=(3,
11    n_part))
12
13    grid = np.arange(n_mesh) + 0.5
14    densities = np.zeros(shape=(n_mesh, n_mesh, n_mesh))
15    cellvol = 1.
16
17    for p in range(n_part):
18        cellind = np.zeros(shape=(3, 2))
19        dist = np.zeros(shape=(3, 2))
20
21        for i in range(3):
22            cellind[i] = np.where((abs(positions[i, p] - grid) < 1)
23            |
24            (abs(positions[i, p] - grid - 16)
25            < 1) |
26            (abs(positions[i, p] - grid + 16)
27            < 1))[0]
28            dist[i] = abs(positions[i, p] - grid[cellind[i].astype(
29            int)])
30
31            cellind = cellind.astype(int)
32
33            for (x, dx) in zip(cellind[0], dist[0]):
34                for (y, dy) in zip(cellind[1], dist[1]):
35                    for (z, dz) in zip(cellind[2], dist[2]):
36                        if dx > 15: dx = abs(dx - 16)
37                        if dy > 15: dy = abs(dy - 16)
38                        if dz > 15: dz = abs(dz - 16)
39
40                        densities[x, y, z] += (1 - dx)*(1 - dy)*(1 - dz
41                        ) / cellvol
42
43    return densities

```

Listing 5: Code for the Cloud-In-Cell method to distribute densities on the grid.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from plotting import set_styles
4 from cic import get_densities
5 from algorithms import fft_nd
6
7 def plot_at_zslices(data, savename, cb_label, suptitle,
8     z_slices=[4.5, 9.5, 11.5, 14.5]):
9
10    fig, axs = plt.subplots(2,2,figsize=(10,8),sharex=True,sharey=
11    True)
12
13    for ax, z in zip(axs.flatten(), z_slices):

```

```

13     # Integer z appear at the edges, but in our array we have
14     the centers
15     # so z=4.5 occurs at index 4
16     im = ax.imshow(data[:, :, int(z)], cmap='jet', origin='lower'
17 )
18     ax.set_title(rf'$Z$ = {z}')
19     ax.grid(False)
20     # Setup a colorbar for all
21     # https://stackoverflow.com/questions/13784201/how-to-have-one-
22     colorbar-for-all-subplots
23     fig.subplots_adjust(right=0.8)
24     cbar_ax = fig.add_axes([0.85, 0.11, 0.05, 0.77])
25     cb = fig.colorbar(im, cax=cbar_ax)
26
27     cb.set_label(cb_label)
28     for ax in axs[:, 0]:
29         ax.set_ylabel(r'$Y$')
30     for ax in axs[-1, :]:
31         ax.set_xlabel(r'$X$')
32     fig.suptitle(suptitle)
33     plt.savefig(f'results/{savename}.png', bbox_inches='tight')
34
35 def compute_forces(eps=1e-15):
36     # Make density grid and compute the density contrasts
37     mean_rho = 1024/(16**3)
38     rho = get_densities()
39     delta = (rho - mean_rho)/mean_rho # Density contrasts
40
41     # Plot slices at various z
42     plot_at_zslices(delta, savename='density_contrast_slices',
43 cb_label='Density Contrast', suptitle='Grid Density Contrast')
44
45     # Apply FFT to \delta to get k^2 \Phi~
46     phi_fft = fft_nd(delta)
47
48     # Compute all k = sqrt(k_x^2 + k_y^2 + k_z^2)
49     # The below only works for data with equal sized vertices in 3D
50     k_vals_sq = (np.arange(phi_fft.shape[0], dtype=np.float64))** 2
51     # Using this notation we force each array along a distinct axis
52     # resulting in a cube
53     k_cube_sq = k_vals_sq[None, None, :] + k_vals_sq[None, :, None]
54     + k_vals_sq[:, None, None]
55     k_cube_sq[k_cube_sq < eps] = eps # Fight divide by zero error
56     by setting zero to ~eps.m
57
58     potential = fft_nd(phi_fft/k_cube_sq, inverse=True)
59
60     # Plot the log of the FFT potential
61     plot_at_zslices(np.log10(np.abs(phi_fft)), savename='
62 fft_potential', cb_label=r'$\log_{10}(|\tilde{\Phi}|)$',
63 suptitle='FFT-Space Potential')
64     # Plot the potential
65     plot_at_zslices(np.real(potential), savename='potential_slices'
66 , cb_label='Potential', suptitle='Grid Potential')
67
68 def main():
69     set_styles()
70     compute_forces()
71
72 if __name__ == '__main__':
73     main()

```

Listing 6: All code to apply the procedures outlined below to compute the

potential.

```

1 # FFT
2 def dft_recursive(x, inverse):
3     """Function to be called recursively by the FFT algorithm to
4     perform the DFT on
5     subsets of the array following the Danielson–Lanczos lemma. For
6     speed we make use
7     of trigonometric recurrence, therefore we never have to compute
8     a complex exponent."""
9     N = len(x)
10    if N > 2:
11        even = dft_recursive(x[::2], inverse)
12        odd = dft_recursive(x[1::2], inverse)
13        x = np.append(even, odd)
14
15    # If we want an iFFT, a -1 should appear in the exponent
16    if inverse:
17        inv_fac = -1.
18    else:
19        inv_fac = 1.
20
21    # Define the trig. recurrence variables
22    theta = 2.*np.pi/N
23    alpha = 2.*(np.sin(theta/2)**2)
24    beta = np.sin(theta)
25    cos_k = 1. # We start with k = 0; cos(0) = 1
26    sin_k = 0. # sin(0) = 1
27
28    for k in range(0, N//2):
29        k2 = k + N//2 # Index of the 'odd' number
30        t = x[k]
31
32        Wnk = cos_k + inv_fac*1j*sin_k #np.exp(inv_fac*2.j*np.pi*k/
33        N)
34        second_factor = Wnk * x[k2]
35
36        # one step of the fourier transform
37        x[k] = t + second_factor
38        x[k2] = t - second_factor
39
40        # Update trig.
41        cos_k_new = cos_k - alpha * cos_k - beta*sin_k
42        sin_k_new = sin_k - alpha * sin_k + beta*cos_k
43        cos_k, sin_k = cos_k_new, sin_k_new
44
45    return x
46
47 def fft(x, inverse=False):
48     """Apply the FFT algorithm to samples x using the recursive
49     Cooley–Tukey algorithm
50     If the length of x is not a power of 2, zeros are appended up
51     to the closest higher
52     power of 2. This function returns a complex array.
53     If inverse is set to True, a '-' sign is introduced in the
54     exponent of  $WN^k$ """
55     # Check the dimensionality of the incoming data
56     if len(x.shape) > 1:
57         return fft_nd(x, inverse)
58
59     # Check if N is a power of 2
60     N = len(x)

```

```

54     if (np.log2(N)%1) > 0: # Check if it is not an integer
55         diff = int(2**(np.ceil(np.log2(N))) - N) # amount of zeros
           to add to make N a power of 2
56         x = np.append(x, np.zeros(diff))
57         N = len(x)
58
59     # Cast x into a complex array so we can store
60     x = np.array(x, dtype=np.cdouble)
61     x_fft = dft_recursive(x, inverse)
62
63     if inverse:
64         x_fft /= N
65
66     return x_fft
67
68 def fft_nd(x, inverse=False):
69     """Apply the Fourier transform to multidimensional data. This
           can easily be done by performing
           the FFT algorithm along each axis separately consecutively"""
70     dim = len(x.shape)
71     func = lambda x: fft(x, inverse)
72     # Start with dim 0 and work up to the highest dimension
73     for i in range(dim):
74         x = np.apply_along_axis(func, i, x)
75     return x
76

```

Listing 7: Implementation of the recursive Cooley-Tukey algorithm with trigonometric recurrence.

In the previous section we could brute force the computation of the motion of the planets relatively easy because there are only eight objects, and eight relevant forces. However if interparticle interactions were important this task would have already been a lot more difficult as there would have been 56 forces to be taken into account at each timestep. Expanding beyond $N = 8$ rapidly worsens this problem. In this section we will investigate one of the methods to transform such a problem from $\mathcal{O}(N^2)$ to $\mathcal{O}(N \log(N))$, namely the Fast Fourier Transform (FFT). We will test this by computing the gravitational potential on a grid for a volume of $16 \times 16 \times 16$ units with 1024 particles, each with a mass of 1. The edges of the grid are connected such that $x = 16 \equiv 0$. The mass of each particle is assigned to a grid point using the Cloud-In-Cell method written by Zorrry Belcheva, given above. We start by filling out the grid, and present the density contrast $\delta = (\rho - \bar{\rho})/\bar{\rho}$ at 4 separate 2D slices at constant z Figure

To calculate the gravitational force, we first need the gravitational potential which we can compute using the Poisson equation

$$\nabla^2 \Phi = 4\pi G \rho = 4\pi G \bar{\rho} (1 + \delta) \quad (2)$$

We can simplify this because we only need to solve the spatial dependence of this equation, i.e. $\nabla^2 \Phi \propto \delta$. Using fourier transforms we can easily do this with by computing the fourier transform of the density contrast, using that to compute the transformed potential, and then applying an inverse fourier transformation to find the true potential. Mathematically, we can write this as

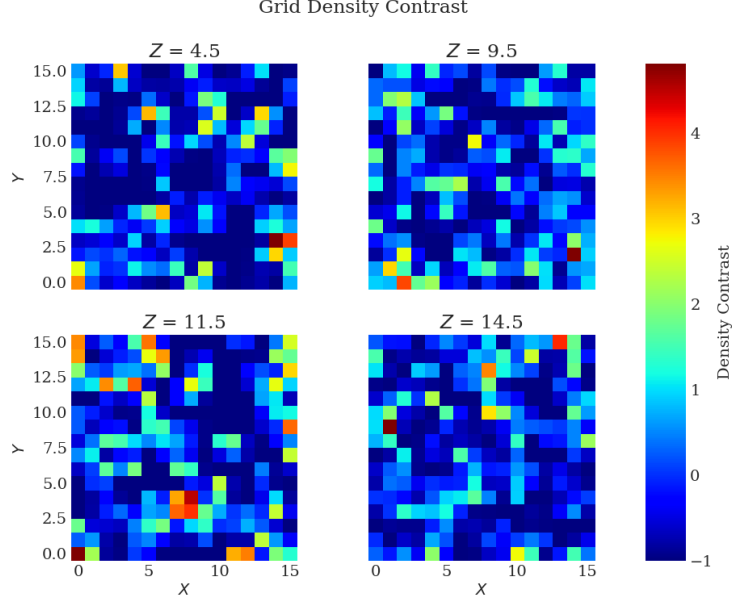


Figure 5: Density contrast δ in four separate 2D slices at constant z computed using the Cloud-In-Cell method on a 3D grid.

$$\begin{aligned}\delta &\xrightarrow{\text{FFT}} \tilde{\delta} \propto k^2 \tilde{\Phi} \\ \tilde{\Phi} &\propto \frac{\tilde{\delta}}{k^2} \\ \frac{\tilde{\delta}}{k^2} &\xrightarrow{\text{IFFT}} \Phi\end{aligned}$$

The only question that now remains is what k is. We can calculate this from the wavenumber vector as

$$k^2 = k_x^2 + k_y^2 + k_z^2. \quad (3)$$

Where k_x , k_y , and k_z are just equal to the grid coordinates and therefore take on integer values from 0 to 15.

The FFT algorithm in this work is implemented using the recursive Cooley-Tukey algorithm with trigonometric recurrence for additional speed up.

2.1 Results

We start with the first two steps from our potential calculation procedure, i.e. apply a fourier transform to the density contrast, and divide it by k^2 . In Figure 6 we present the log of the absolute value of $\tilde{\Phi}$. We take the absolute value because there is a complex component which we want to include in the figure.

We then apply the inverse Fourier Transform with the FFT algorithm to the data shown in Figure 6 to find the gravitational potential which we present in Figure 7.

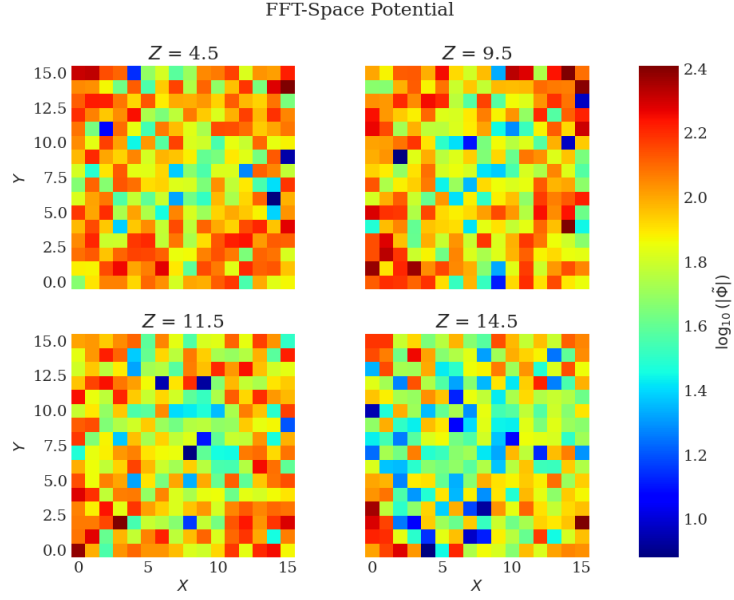


Figure 6: Logarithm of the absolute value of the fourier transformed gravitational potential, this is equal to the fourier transformed density contrast divided by k^2 .

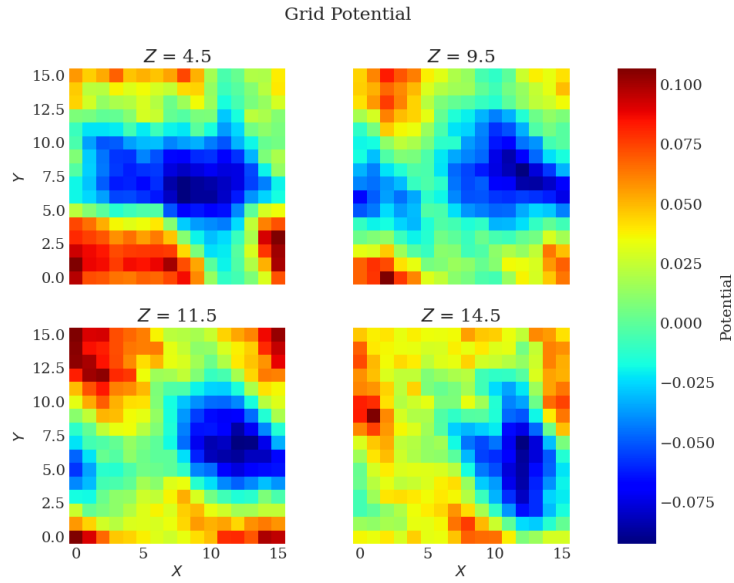


Figure 7: Gravitational potential of the density contrast grid shown in Figure 5, computed using fourier transformations as described in text.

We can see that the calculated potentials are quite smooth when compared to the more erratic distribution of mass in Figure 5. However this makes sense as we can see from the Poisson equation that the density contrast influences only the second order derivative of the potential. Comparing the density contrast and potential figures, we can see that the regions of high and low density correspond well to the regions of high and low potential respectively, which of course is exactly as expected for the gravitational potential. For example the low density region in the center of the $Z = 4.5$ slice corresponds to a potential well, and the high density 'square' in the $Z = 11.5$ slice at $Y = 3$ corresponds to a high potential arm reaching towards the center.

As a final note, a test run of the procedure outlined above with $k = 1$ resulted in a 'potential' that was exactly equal to the density contrast. While not physically relevant, this is a proof that the FFT algorithms properly worked because the inverse fourier transform of a fourier transformed dataset should be equal to the original data set.

3 Galaxy Classification

In this section we will attempt to implement a logistic regression algorithm to classify galaxies into two classes, ellipticals and spirals, using four parameters:

- κ_{co} , which indicates how much a galaxy is dominated by ordered rotation
- A color estimate (higher is redder, and therefore contains older stars)
- A measure of how extended the galaxy is
- Flux of an emission line tracing the star formation rate

We start by preprocessing all of the features by scaling them to have zero mean and unit standard deviation. This is implemented for each feature f_i separately as

$$f_{i,\text{scaled}} = \frac{f_i - \mu_i}{\sigma_i} \quad (4)$$

To get a sense of how each of our features is distributed, we plot their histograms in Figure 8, and present the first ten samples below.

```

1  -1.581862133006098681e+00 -6.612307828919693729e-03
   -3.475533387671284058e-02 4.149921715612928282e-03
2  1.574559067587914640e+00 -7.944354060046822097e-01
   2.031450234565139734e-01 1.605487869858183633e-02
3  1.531990475165863508e+00 8.573405014974248006e-01
   -2.079754178048100477e-01 4.166315770216809378e-02
4  1.297215471156623057e+00 1.226267345173265078e+00
   1.858865005815529270e-01 6.825469423304870997e-04
5  5.181259591131069930e-01 8.257714425578099871e-01
   -2.002589858418064583e-01 1.555766189074690373e-02
6  -1.647524013432828394e+00 -2.909640578232345343e-01
   -8.488083597139754743e-02 1.855013775209162245e-03
7  -1.312495960360358760e+00 -9.756536637063890627e-01
   6.454576379671385367e-02 6.772650820885613675e-03
8  9.418303994505743126e-02 -9.715157240517435788e-01
   -1.525183484703764025e-01 7.745013475795247196e-03
9  -7.308293662245395339e-01 -1.520148623547938893e+00
   -1.654403149798442940e-01 1.175440493844448313e-02
10 -1.198035514238125154e+00 -7.748995242729127542e-01
    -1.553888673156535449e-01 9.596892572192124507e-03

```

Listing 8: Scaled feature values for the first ten entries of the galaxy dataset

```

1  def hist(x, binmin, binmax, nbins, log=False, return_centers=False)
   :
2  :
3  :
4  if log:
5     bin_edges = np.logspace(np.log10(binmin), np.log10(binmax),
6     nbins + 1)
7  else:
8     bin_edges = np.linspace(binmin, binmax, nbins + 1)
9  if return_centers:
10     bin_centers = np.zeros(nbins)
11
12 histogram = np.zeros(nbins)
13 for i in range(nbins):
14     bin_mask = (x >= bin_edges[i]) * (x < bin_edges[i + 1])

```

```

13         if log:
14             histogram[i] = len(x[bin_mask])
15         else:
16             histogram[i] = len(x[bin_mask])
17         if return_centers:
18             bin_centers[i] = bin_edges[i] + 0.5 * (bin_edges[i+1] -
19             bin_edges[i])
20     if return_centers:
21         return histogram, bin_edges, bin_centers
22     return histogram, bin_edges

```

Listing 9: Code to make the histograms in Figure 8

```

1 def preprocess_data(fname, plot=False, nbins=None):
2     data = np.genfromtxt(fname)
3     features = data[:, :-1]
4     labels = data[:, -1]
5
6     # Rescale the features
7     for i in range(features.shape[1]):
8         mean = np.mean(features[:, i])
9         std = np.std(features[:, i])
10        features[:, i] = (features[:, i] - mean)/std
11
12    # Save to a text file
13    np.savetxt('results/scaled_features.txt', features)
14
15    if plot:
16        fig, axs = plt.subplots(2, 2, figsize=(8, 8), sharex=False,
17        sharey=False, tight_layout=True)
18        for i in range(features.shape[1]):
19            ax = axs.flatten()[i]
20            binmin = np.min(features[:, i])
21            binmax = np.max(features[:, i])
22            n, bin_edges, bin_centers = hist(features[:, i], binmin,
23            binmax, nbins, return_centers=True)
24
25            ax.step(bin_centers, n, where='mid', lw=3)
26            ax.set_title(f'Feature {i}')
27            ax.set_xlabel('Rescaled Values')
28
29            for ax in axs[:, 0]:
30                ax.set_ylabel('Counts')
31            plt.savefig('results/scaled_features_dist.png')
32
33    return features, labels

```

Listing 10: Code for the feature preprocessing and plotting.

It is immediately noticable that while most values for features 2 and 3 lie around 0, there are outliers at extremely high values. This is less true for feature 0 which shows a peak at 0 and has a relatively constant distribution in its wings. Feature 1 appears to be nicely distributed close to a Gaussian. Due to the few outliers in features 2 and 3 we will limit ourselves to the 2% and 98% percentiles for plotting in the rest of this work to better highlight the variations in the regions where the bulk of the datapoints reside. Ofcourse we do not discard these points for the training process.

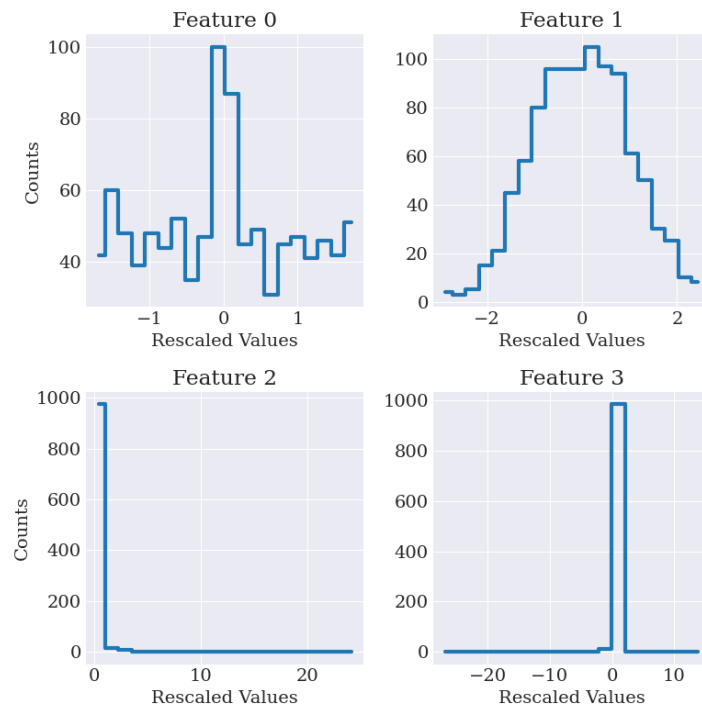


Figure 8: Distributions of the scaled features used for the classification algorithms in this section.

3.1 Two Feature Combinations

```

1 def classify(features, labels, lr, minim_type,
2             p=2):
3     # Start with 'simple' classification of two columns using
4     # constant step size
5     names = [r'$\kappa$-\mathrm{co}', 'Color', 'Extended', 'Emission Flux']
6     snames = ['kappa', 'color', 'extended', 'emission_flux'] # names for the savefiles
7     fig1, ax1 = plt.subplots(1,1, figsize=(9,4))
8     table_txt = "" # Fill out for the confusion matrix table
9
10    for i in range(features.shape[1]):
11        for j in range(i+1, features.shape[1]):
12            feats = features[:, [i,j]] # First two columns
13            # Add in the bias
14            feats = np.append(feats, np.ones(feats.shape[0])[:, np.newaxis], axis=1)
15            params, loss = logistic_regression(feats, labels, lr=lr, minim_type=minim_type)
16
17            # Make the confusion matrix and compute F1
18            logi = logistic_func(feats, params)
19            predictions = np.zeros(len(logi))
20            predictions[logi >= 0.5] = 1 # label = 1 if logi > 0.5
21            otherwise label = 0
22            conf_mat = make_confusion_matrix(labels, predictions)
23            f1 = compute_f1_score(conf_mat)
24
25            # Plot loss curve
26            ax1.step(np.arange(loss.shape[0]), loss, label=f'{names[i]} + \n{names[j]}')
27
28            # Plot data
29            fig2, ax2 = plt.subplots(1,1,)
30            ax2.scatter(feats[:,0], feats[:,1], s=3, c=labels, cmap='seismic')
31
32            # Make the 'cutoff line'
33            xx = np.linspace(np.min(features[:,0]), np.max(features[:,1]))
34            ax2.plot(xx, (-params[0]/params[1] * xx) + params[1], c='black', ls='—', lw=3)
35            ax2.set_xlabel(names[i])
36            ax2.set_ylabel(names[j])
37            ax2.set_title(f'Final Loss = {loss[-1]}')
38            ax2.set_xlim(np.percentile(feats[:,0], p), np.percentile(feats[:,0], 100-p))
39            ax2.set_ylim(np.percentile(feats[:,1], p), np.percentile(feats[:,1], 100-p))
40            fig2.savefig(f'results/2d_fit_{minim_type}-{snames[i]}-{snames[j]}.png', bbox_inches='tight')
41
42            # Add text for the table. Format is
43            # Feats. TN TP FN FP F1
44            conf_mat = np.array(conf_mat, dtype=int)
45            table_txt += f'{names[i]} + {names[j]} & {conf_mat[0][0]} & {conf_mat[1][1]} & {conf_mat[1][0]} & {conf_mat[0][1]} & {f1} \\\n'
46
47    ax1.set_xlabel('Iteration Number')

```

```

46     ax1.set_ylabel('Loss')
47     ax1.set_title('Loss over Time for 2D Fits')
48
49     # Shrink current axis by 25%
50     # https://stackoverflow.com/questions/4700614/how-to-put-the-
51     # legend-outside-the-plot
52     box = ax1.get_position()
53     ax1.set_position([box.x0, box.y0, box.width * 0.75, box.height
54     ])
55     ax1.legend(loc='center left', bbox_to_anchor=(1, 0.5))
56
57     fig1.savefig(f'results/2d_fit_losses_{minim_type}')
58
59     # write to file
60     with open(f'results/confmat_tab_{minim_type}.txt', 'w') as f:
61         table_txt = table_txt[:4]
62         f.write(table_txt)
63         f.close()
64
65 def galaxies():
66     set_styles()
67     fname = 'galaxy_data.txt'
68     plot = True
69     nbins = 20
70
71     lr = 0.1
72     minim_types = ['constant_step', 'line_minim']
73     #####
74
75     features, labels = preprocess_data(fname, plot=plot, nbins=
76     nbins)
77     for minim_type in minim_types:
78         classify(features, labels, lr, minim_type)

```

Listing 11: Code for fitting the classification algorithm and computing the F1 score for each two feature combination and both minimization types

```

1 # LOGISTIC REGRESSION
2 def logistic_func(X, theta):
3     """Estimate the labels of X given model paramters theta.
4     This only works for two object classification"""
5     z = np.dot(theta, X.T)
6     sigma = 1./(1.+np.exp(-z))
7     return sigma
8
9 def logistic_loss(X, Y, theta, hypothesis_func=logistic_func,
10                  return_gradient=False):
11     """Logistic loss functions for features X, labels Y and
12     parameters theta"""
13     h_theta = hypothesis_func(X, theta)
14     # Vectorized version of the logistic loss
15     loss = (-1./len(Y)) * np.sum((Y * np.log(h_theta) + (1. - Y) *
16     np.log(1. - h_theta)))
17     if return_gradient:
18         grad = np.sum((1./len(Y)) * X.T * (h_theta - Y), axis=1)
19         return loss, grad
20     return loss
21
22 def logistic_regression(X, Y, lr=0.1, eps=1e-6, max_iter=int(1e4),
23                        cost_func=logistic_loss,
24                        minim_type='constant_step'):

```

```

24 """Perform logistic regression on features X and labels Y
25 X should have shape (m, n); Y should have shape (m)"""
26 theta = np.ones(X.shape[1])
27 loss_ar = np.zeros(max_iter)
28 # Define a function where we only have to feed in theta,
29 # because X, Y are constant
30 loss_func = lambda theta, return_gradient=False: cost_func(X, Y
31 , theta, return_gradient=return_gradient)
32
33 for i in range(max_iter):
34     match minim_type:
35         # Use a constant learning rate to minimize
36         case 'constant_step':
37             loss, grad = loss_func(theta, return_gradient=True)
38             loss_ar[i] = loss
39             if np.abs(np.max(grad)) < eps:
40                 print('Gradient reached epsilon threshold')
41                 print(f'Final Loss = {loss}')
42                 return theta, loss_ar[:i+1]
43             theta -= lr * grad
44
45         # Step along -grad, but use line minimization to find
46         # the step size
47         case 'line_minim':
48             loss, grad = loss_func(theta, return_gradient=True)
49             loss_ar[i] = loss
50             if np.abs(np.max(grad)) < eps:
51                 print('Gradient reached epsilon threshold')
52                 print(f'Final Loss = {loss}')
53                 return theta, loss_ar[:i+1]
54
55             step_size = line_minimization(loss_func, theta,
56 grad)
57             theta += step_size*grad
58
59         # Use a downhill simplex to walk down the loss
60         # landscape
61         case 'simplex':
62             theta, _ = downhill_simplex(loss_func, theta,
63 eval_separate=True)
64             return theta, loss
65
66     print('Maximum number of iterations reached.')
67     return theta, loss_ar
68
69 # CONFUSION MAT. AND F1 SCORES
70 def make_confusion_matrix(labels, pred):
71     """
72     n_features = int(np.max(labels) + 1)
73     mat = np.zeros((n_features, n_features))
74     for i in range(n_features):
75         true = np.where(labels == i)[0]
76         for j in range(n_features):
77             if i == j:
78                 mat[i][i] = len(np.where(pred[true] == i)[0])
79             else:
80                 mat[i][j] = len(np.where(pred[true] == j)[0])
81
82     return mat
83
84 def compute_f1_score(mat):
85     """Compute the F1 score for a 2D confusion matrix. It is
86     defined as

```

```

79         F1 = 2 * (precision * recall)/(precision+recall)
80         with
81             precision = TP/(TP+FP)
82             recall = TP/(TP+FN)"""
83     precision = mat[1][1] / (mat[1][1] + mat[0][1])
84     recall = mat[1][1] / (mat[1][1] + mat[1][0])
85     return 2. * (precision * recall) / (precision + recall)

```

Listing 12: Logistic Regression, Confusion Matrix and F1 Score code

```

1  import numpy as np
2
3  # LINE MINIMIZATION WITH GOLDEN SECTION SEARCH
4  def parabola_min_analytic(a, b, c, fa, fb, fc):
5      """Analytically computes the x-value of the minimum of a
6      parabola
7      that crosses a, b and c
8      """
9      top = (b-a)**2 * (fb-fc) - (b-c)**2 * (fb-fa)
10     bot = (b-a) * (fb-fc) - (b-c) * (fb-fa)
11     return b - 0.5*(top/bot)
12
13 def make_bracket(func, bracket, w=(1+np.sqrt(5))/2, dist_thresh
14     =100, max_iter=10000):
15     """Given two points [a, b], attempts to return a bracket
16     triplet
17     [a, b, c] such that f(a) > f(b) and f(c) > f(b).
18     Note we only compute f(d) once for each point to save computing
19     time"""
20     a, b = bracket
21     fa, fb = func(a), func(b)
22     direction = 1 # Indicates if we're moving right or left
23     if fa < fb:
24         # Switch the two points
25         a, b = b, a
26         fa, fb = fb, fa
27         direction = 1 # move to the left
28
29     c = b + direction * (b - a) * w
30     fc = func(c)
31
32     for i in range(max_iter):
33         if fc > fb:
34             return np.array([a, b, c]), i+1
35         d = parabola_min_analytic(a, b, c, fa, fb, fc)
36         fd = func(d)
37         if np.isnan(fd):
38             #print(f'New point d:{d} gives fd:{fd}. Breaking
39             function')
40             return np.array([a,b,c]), i+1
41         # We might have a bracket if b < d < c
42         if (d>b) and (d<c):
43             if fd > fb:
44                 return np.array([a, b, d]), i+1
45             elif fd < fc:
46                 return np.array([b, d, c]), i+1
47             # Else we don't want this d
48             #print('no parabola, in between b and c')
49             d = c + direction * (c - b) * w
50         elif (d-b) > 100*(c-b): # d too far away, don't trust it
51             #print('no parabola, too far away')
52             d = c + direction * (c - b) * w

```



```

48         elif d < b:
49             pass# print('d smaller than b')
50
51         # we shifted but didn't find a bracket. Go again
52         a, b, c = b, c, d
53         fa, fb, fc = fb, fc, fd
54
55     print('WARNING: Max. iterations exceeded. No bracket was found.
56           Returning last values')
57     return np.array([a, b, c]), i+1
58
59 def golden_section_search(func, bracket, target_acc=1e-5, max_iter=
60     int(1e5)):
61     """Once we have a start 3-point bracket surrounding a minima,
62     this function iteratively
63     tightens the bracket to search of the enclosed minima using
64     golden section search."""
65     w = 2. - (1.+np.sqrt(5))/2 # 2 - golden ratio
66     a, b, c = bracket
67     fa, fb, fc = func(a), func(b), func(c)
68     for i in range(max_iter):
69         # Set new point in the largest interval
70         # We do this separately because the bracket propagation can
71         # just not be generalized sadly
72         if np.abs(c-b) > np.abs(b-a): # we tighten towards the
73             right
74             d = b + (c-b)*w
75             fd = func(d)
76             if fd < fb: # min is in between b and c
77                 a, b, c = b, d, c
78                 fa, fb, fc = fb, fd, fc
79             else: # min is in between a and d
80                 a, b, c = a, b, d
81                 fa, fb, fc = fa, fb, fd
82             else: # we tighten towards the left
83                 d = b + (a-b)*w
84                 fd = func(d)
85                 if fd < fb: # min is in between a and b
86                     a, b, c = a, d, b
87                     fa, fb, fc = fa, fd, fb
88                 else: # min is in between d and c
89                     a, b, c = d, b, c
90                     fa, fb, fc = fd, fb, fc
91
92     if np.abs(c-a) < target_acc:
93         return [b,d][np.argmin([fb, fd])], i+1 # return the x
94         point corresponding to the lowest f(x)
95
96 def line_minimization(func, x_vec, step_direction, method=
97     golden_section_search, minimum_acc=1e-5):
98     """
99     # Make a function f(x+lmda*n)
100     minim_func = lambda lmda: func(x_vec + lmda * step_direction)
101     bracket_edge_guess = [0, 1] # keeps the steps relatively small
102     to combat divergence
103     bracket, _ = make_bracket(minim_func, bracket_edge_guess) #
104     make a 3-point bracket surrounding a minimum
105
106     # Use a 1-D minimization method to find the 'best' lmda
107     minimum, _ = method(minim_func, bracket, target_acc=minimum_acc
108     )
109     return minimum

```

```

99
100
101 # CODE FOR THE N-DIMENSIONAL DOWNHILL SIMPLEX
102 def compute_centroid(A):
103     """Compute the centroid of N points in N dimensional space. x
    should
104     be an NxN array of N vectors with N dimensions (in that order).
    The function
105     then returns one ndarray of length N with the centroid
    coordinates."""
106     return (1./A.shape[1]) * np.sum(A, axis=0)
107
108
109 def downhill_simplex(func, start, shift_func=lambda x: x+1,
    max_iter=int(1e5), target_acc=1e-5,
    eval_separate=False):
110     """Finds the minimum of a function using the downhill simplex
    method
111     INPUT:
112         func: A function taking only one variable as input with
    dimension N
113         start: N-Dimensional numpy array where the function starts
    searching for a minimum
114         shift_func: A function taking only one float as input
    dictating how to mutate the
115             initial simplex vertices
116
117     OUTPUT:
118
119     """
120     dim = start.shape[0] # = N
121     # Store N+1 vertice vectors in this matrix. This ordering fails
    if we feed it directly to func,
122     # but it allows us to choose a vertex as vertices[i]. The
    function problem we solve by just
123     # transposing this this matrix.
124     vertices = np.zeros((dim+1, dim))
125     func_vals = np.zeros(dim+1) # Store the f(X) values in here
126
127     # Create the simplex, add slight variation to each vector
    except the first using 'shift_func'
128     vertices[0] = start
129     for i in range(dim):
130         vertices[i+1] = vertices[0]
131         vertices[i+1][i] = shift_func(vertices[i+1][i])
132
133     if eval_separate: # slightly slower, but dodges difficutlies
    with matrix/vector stuff
134         for i in range(dim+1):
135             func_vals[i] = func(vertices[i])
136     else:
137         func_vals = func(vertices.T)
138     # Start algorithm
139     for i in range(1, max_iter+1):
140         # Sort everything by function value
141         sort_idx = merge_sort(key=func_vals)
142         vertices = vertices[sort_idx]
143         func_vals = func_vals[sort_idx]
144
145     # Check if we have reached our accuracy level by comparing
    the best and worst function evals.
146     accuracy =(np.abs(func_vals[-1] - func_vals[0])/np.abs

```

```

148 (0.5*(func_vals[-1] + func_vals[0]))
149
150     if accuracy < target_acc:
151         print(accuracy, target_acc)
152         return vertices[0], i # corresponds to func_vals[0], so
153         the best point
154
155     # Compute the centroid of all but the last (worst) point
156     centroid = compute_centroid(vertices[:-1])
157     # Try out a new points
158     x_try = 2.* centroid - vertices[-1]
159     f_try = func(x_try)
160
161     if f_try < func_vals[-1]:
162         # There is improvement in this step
163         if f_try < func_vals[0]:
164             # We are the best point. Try expanding
165             x_exp = 2*x_try - centroid
166             f_exp = func(x_exp)
167             if f_exp < f_try:
168                 # expanded point is even better. Replace x_N
169                 vertices[-1] = x_exp
170                 func_vals[-1] = f_exp
171             else:
172                 # x_try was good, x_exp is not better
173                 vertices[-1] = x_try
174                 func_vals[-1] = f_try
175         else:
176             # Better than x_N, not better than x_0. Just accept
177             the point
178             vertices[-1] = x_try
179             func_vals[-1] = f_try
180
181     else:
182         # This point is worse than what we had. First try
183         contracting, new x_try
184         x_try = 0.5*(centroid+vertices[-1])
185         f_try = func(x_try)
186         if f_try < func_vals[-1]:
187             # contracting improved x
188             vertices[-1] = x_try
189             func_vals[-1] = f_try
190         else:
191             # Nothing worked, just contract all points towards
192             the best points
193             vertices = 0.5*(vertices[0] + vertices) # x0 doesnt
194             shift here: 0.5(x0 + x0) = x0
195             # need to evaluate all but x0 because they shifted
196             if eval_separate:
197                 for i in range(dim):
198                     func_vals[i+1] = func(vertices[i+1])
199             else:
200                 func_vals[1:] = func(vertices[1:].T)
201
202     print('Maximum Number of Evaluations Reached')
203     return vertices[0], i
204
205 # SORTING
206 def sort_subarrays(a1, a2, k1, k2):
207     """Takes two subarrays 1,2 with indices a1, a2 and values k1,
208     k2 and combines these

```

```

203 into array with indices a such that the k are sorted in
204 ascending order"""
205 N1 = len(a1)
206 N2 = len(a2)
207
208 # We built up a new instance of our sorting array. This is not
209 # memory efficient
210 # TODO: Study the algorithm below to try and find a better
211 # method
212 a_sorted = np.zeros(N1+N2)
213
214 if N1 == 0:
215     return a2
216 if N2 == 0:
217     return a1
218
219 # Walk through the left- and right- sub arrays separately
220 idx1 = 0
221 idx2 = 0
222 while True:
223     if k1[idx1] > k2[idx2]:
224         # Then place the second element to the left of the
225         # first element and take
226         # one step to the right in the right sub-array
227         a_sorted[idx1+idx2] = a2[idx2]
228         idx2 += 1
229
230         # Then, if there are elements remaining in the right
231         # array, keep
232         # placing them to the left as long as they're smaller
233         if idx2 < N2: # need this if statement to save us from
234             indexing errors in the while
235                 while (k1[idx1] > k2[idx2]):
236                     a_sorted[idx1+idx2] = a2[idx2]
237                     idx2 += 1
238
239                 if idx2 >= (N2):
240                     # No more elements left in the right array,
241                     # we can fill out with the left array
242                     for j in range(idx1, N1):
243                         a_sorted[j+idx2] = a1[j]
244                     return a_sorted
245
246                 # Now the element from the left array is smaller
247                 # than the first remaining element
248                 # from the right array, we can safely place it
249                 a_sorted[idx1+idx2] = a1[idx1]
250                 idx1 += 1
251             else:
252                 # No more elements left in the right sub-array
253                 for j in range(idx1, N1):
254                     a_sorted[j+idx2] = a1[j]
255                 return a_sorted
256
257             else:
258                 a_sorted[idx1+idx2] = a1[idx1]
259                 idx1 += 1
260
261             # Check if we have reached the end of the left sub-array
262             # If we have, fill out the rest of the array with the right
263             # sub-array
264             if idx1 == N1:

```

```

256         for j in range(idx2, N2):
257             a_sorted[idx1+j] = a2[j]
258         return a_sorted
259
260
261
262 def merge_sort(a=None, key=None):
263     """Sorts the array or list using merge sort. This function
264     iteratively
265     builds up the array from single elements which are sorted by
266     sort_subarrays
267     Note, in principle one should only provide either 'a' or 'key':
268
269     - If 'a' is provided, that array is sorted in ascending order,
270     and returned
271     by this function
272     - If 'key' is provided, this function returns the indices
273     corresponding to the
274     order in which the array would be sorted
275     - If both 'a' and 'key' are provided, this function assumes '
276     key' has already
277     been previously shuffled, and just swaps indices of the
278     preexisting 'a'
279
280     RETURNS: 'a': numpy array
281     """
282     if key is not None:
283         key = np.array(key)
284
285     if a is None and key is not None:
286         a = np.arange(len(key))
287
288     a = np.array(a)
289     subsize = 1
290     N = len(a)
291     is_sorted = False
292     # Build up the array sorting arrays of increasing subsize
293     while not is_sorted:
294         subsize *= 2
295         if subsize > N:
296             is_sorted = True # After this iteration, the array is
297             sorted
298
299         for i in range(int(np.ceil(N/subsize))):
300             # We need the min(.., N) to ensure that we do not
301             exceed the length of the
302             # array with our indexing
303             subarray1 = a[i*subsize: i*subsize+int(0.5*subsize)] #
304             First half of the interval
305             subarray2 = a[i*subsize+int(0.5*subsize): np.min(((i+1)
306             *subsize, N))]
307             if key is not None:
308                 key1, key2 = key[subarray1], key[subarray2]
309                 sorted_sub = sort_subarrays(subarray1, subarray2,
310                 key1, key2)
311             else:
312                 # we feed in 'subarrayx' twice because if we only
313                 sort a, a is its own key
314                 sorted_sub = sort_subarrays(subarray1, subarray2,
315                 subarray1, subarray2)
316
317             a[i*subsize:subsize*(i+1)] = sorted_sub

```

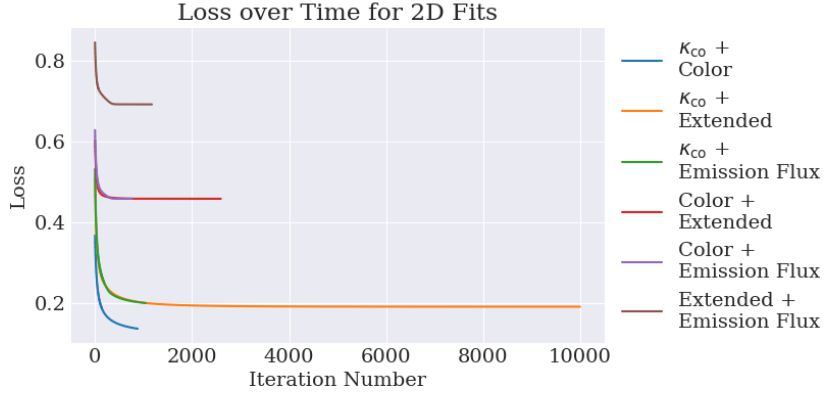


Figure 9: Loss value as a function of iteration number for each of the six possible combinations of two out of four features using logistic regression and a constant step size as minimization method.

```

305
306     return a

```

Listing 13: Code for 'smarter' function minimization

We begin by investigating how well the logistic regression algorithm is able to predict the galaxy class if it is only provided with two of the four available features. We do this by running the algorithm by each of the six possible combinations of two features and investigating the resulting loss values. In each case we also include a bias feature consisting of only ones. As an extra test we also look at the differences between using a constant step size $\eta = 0.1$ and using line minimization using golden section search to discover the minimum. We show the resulting loss curves, the evolution of the loss value as a function of iteration, in Figure 9 for the learning rate method and in Figure 10 for the line minimization method.

We can see that under both minimization methods, all loss values converge to approximately the same values. However line minimization converges ~ 100 – 1000 times faster than using a learning rate of 0.1 . From these plots we can gather that the combination of κ_{CO} and color estimate is the best combination of parameters to estimate galaxy class. Meanwhile, the extended measure and the flux of the emission line tracing star formation rate appear to be a very bad combination of predictors. To get a better grasp of the distribution of the features, and the prediction capability of our models, we plot all six of the two dimensional combinations of distributions of the features coloured by their galaxy class together with the decision boundary as learned by the logistic regression model in Figure 11. We only plot the results of the line minimization method here.

In the last panel, corresponding to the combination with the highest loss, we can see that the two classes are very much intertwined making it very difficult for the algorithm to separate the two. The figure also shows us why the combination of color and κ_{CO} work well together to predict the galaxy class. The two classes appear quite clearly as two distinct clusters with only a little bit of overlap in the

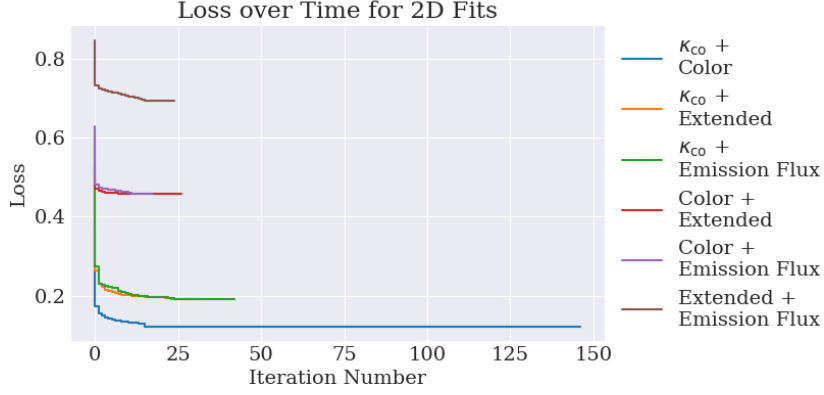


Figure 10: Loss value as a function of iteration number for each of the six possible combinations of two out of four features using logistic regression and line minimization as minimization method.

center. The model attempted to draw its decision boundary between these two clusters, although visually it looks like the true best decision boundary might have been placed slight more to the bottom left.

To investigate these results quantitatively we also compute the confusion matrix and the F1 score for each of these fits. The F1 score is defined as

$$F1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (5)$$

With

$$\text{precision} = \frac{TP}{TP + FP} \quad (6)$$

$$\text{recall} = \frac{TP}{TP + FN} \quad (7)$$

Both the precision and the recall are values between 0 and 1. Therefore the F1 score is also always a value between 0 and 1, where higher numbers correspond to better predictions. We calculate these values for all combinations of features, and present the results in Table 1.

Features	TN	TP	FN	FP	F1
κ_{co} + Color	476	472	28	24	0.9477911646586344
κ_{co} + Extended	439	440	60	61	0.8791208791208791
κ_{co} + Emission Flux	440	438	62	60	0.8777555110220441
Color + Extended	390	392	108	110	0.7824351297405189
Color + Emission Flux	389	393	107	111	0.7828685258964143
Extended + Emission Flux	370	91	409	130	0.2524271844660194

Table 1: Confusion matrix values laid out horizontally and the F1 score computed using these for each combination of two features.

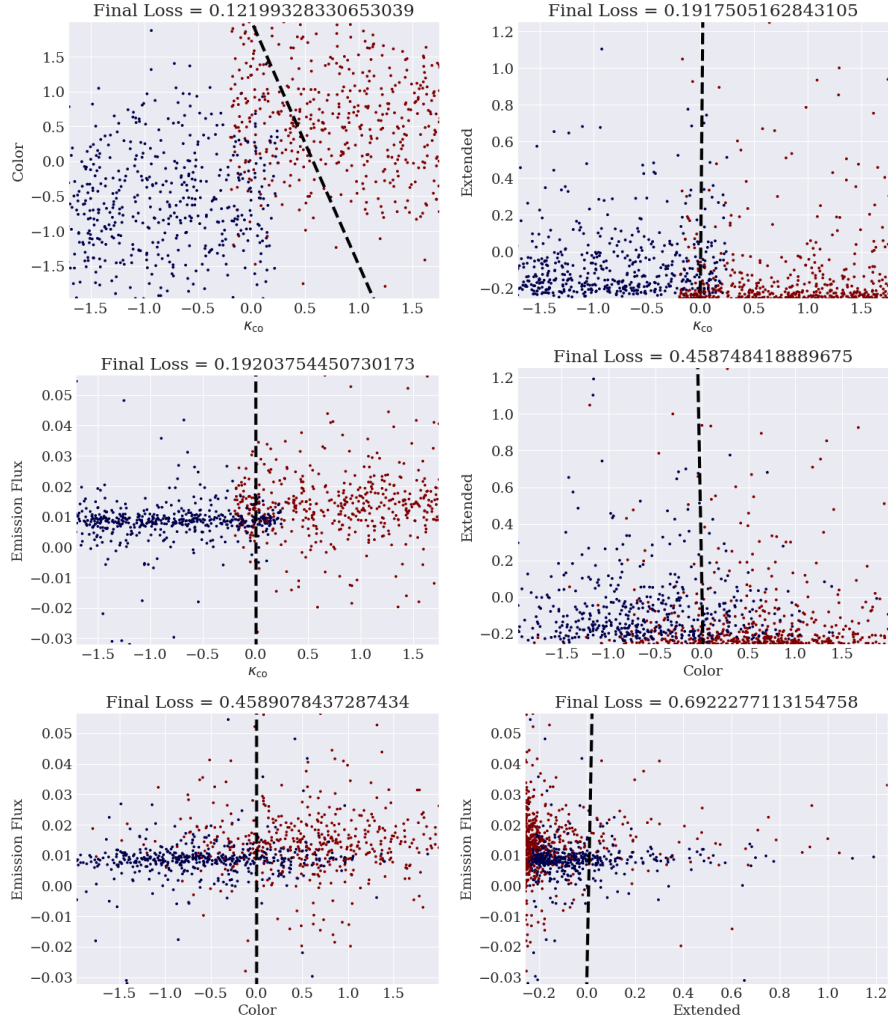


Figure 11: Distributions of each combination of features with the outer 2% cut out due to outliers. Black dashed lines indicate the decision boundary of a classification algorithm trained with logistic regression and line minimization using golden section search. The titles indicate the converged loss value.

Correspondingly to what we can see both in the loss values and visually, the first combination performs best and the last model performs worst. The rest of the models all hover around $F1 \sim 0.8$. Interestingly, while most of the predictors appear to be equally bad at recreating the label for elliptical and spiral galaxies, the final model is a lot worse at detecting when a galaxy in its training set is a spiral (positive class, corresponding to the red dots in Figure 11). This can also be seen in the last panel of that figure because the decision boundary is placed to the right of the majority of galaxies, and there are a lot of blue (negative class, ellipticals) to the left of it still.