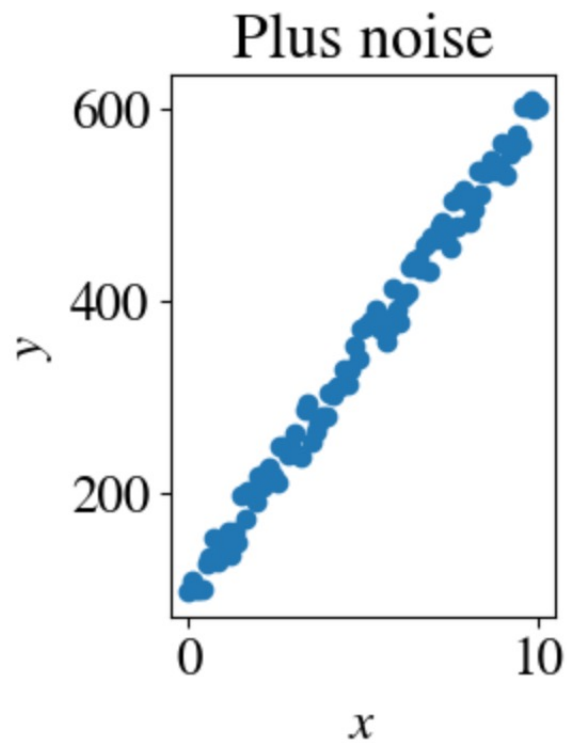
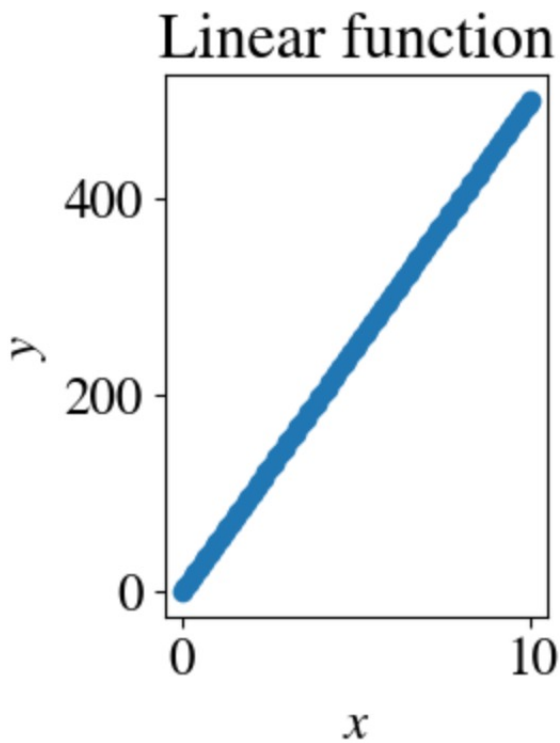


Lecture 2

Linear Numerical Regression

Problem Statement



- **1 Feature:** x (m values/samples)
- **1 Target:** y
- Supervised learning: m known labels y_{true}
- **Hypothesis:** $y_{\text{pred}} = \theta_0 + \theta_1 x$
- Example: $\theta_0 = 100, \theta_1 = 50$
- Mean-squared-error (MSE):
$$\frac{1}{2m} \sum_{i=1}^m (y_{i,\text{pred}} - y_{i,\text{true}})^2 = 106$$

Objective of Linear Regression:

Fit linear function to (generally noisy) data, i.e.:

Find slope (θ_1) & intercept (θ_0) with smallest MSE between y_{true} and $y_{\text{pred}} = \theta_0 + \theta_1 x$ for feature x .

How?

Brute force fit of slope

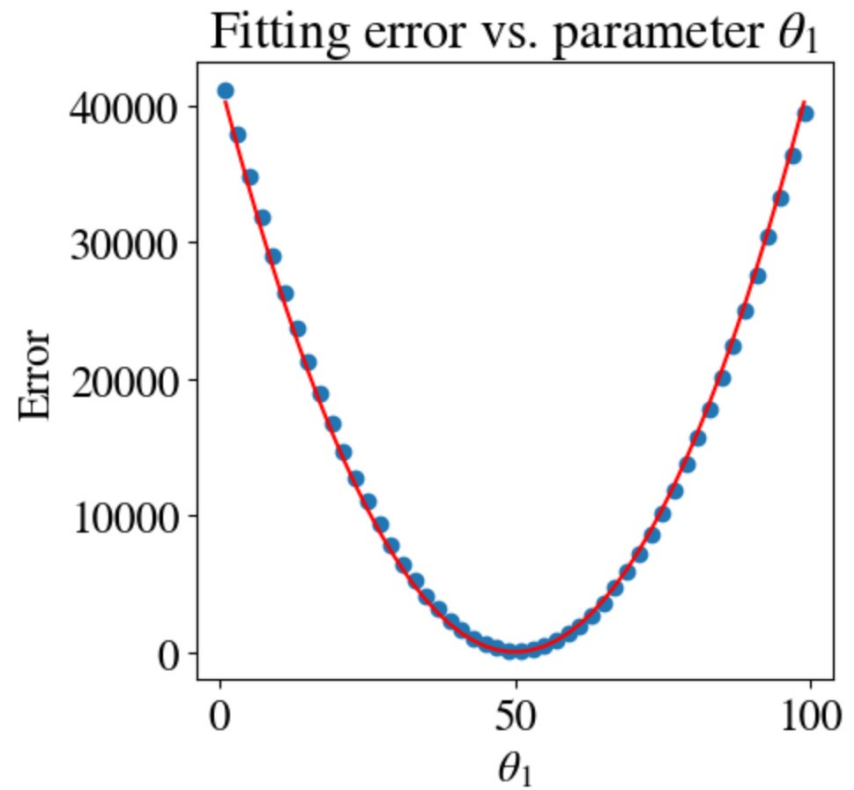
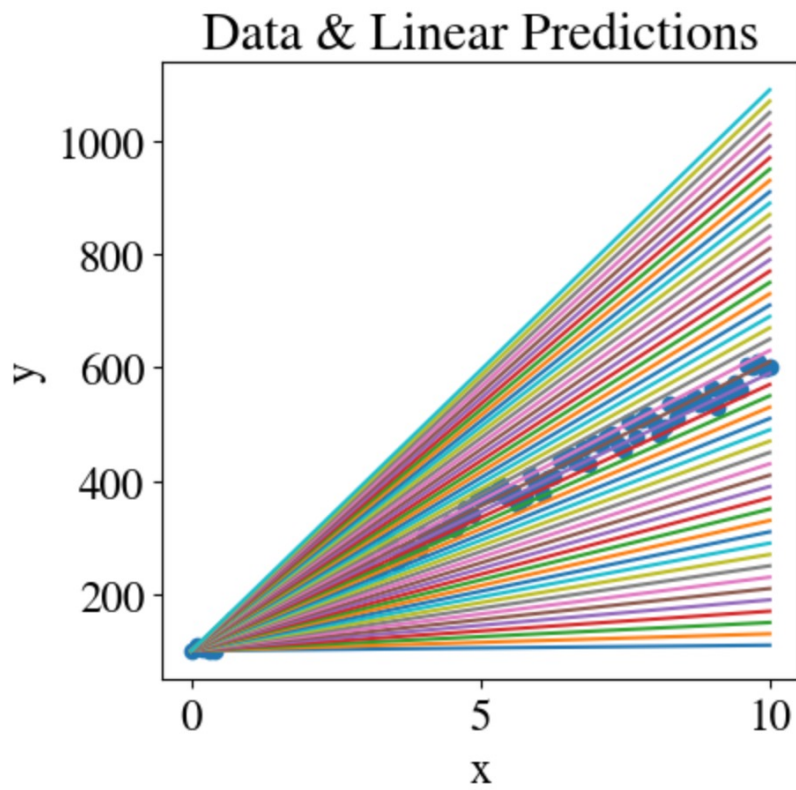
- Simplest example, consider known intercept $\theta_0 = 100$ and only fit θ_1 :



```
1  theta0 = 100
2  errors = np.zeros(50)
3  theta = np.zeros(50)
4
5  for i in range(50):
6      theta[i] = 1 + i*2
7      y_pred = theta0 + theta[i]*x
8      errors[i] = mean_squared_error(y_pred , y)/2
```

Result

- Best fit: $\theta_1 = 51.0$
- **CPU time: 650 ms**
- **MSE = 105.8**



Interpretation of error plot

- If 'true' values were exactly linear (slope of 50) without noise:

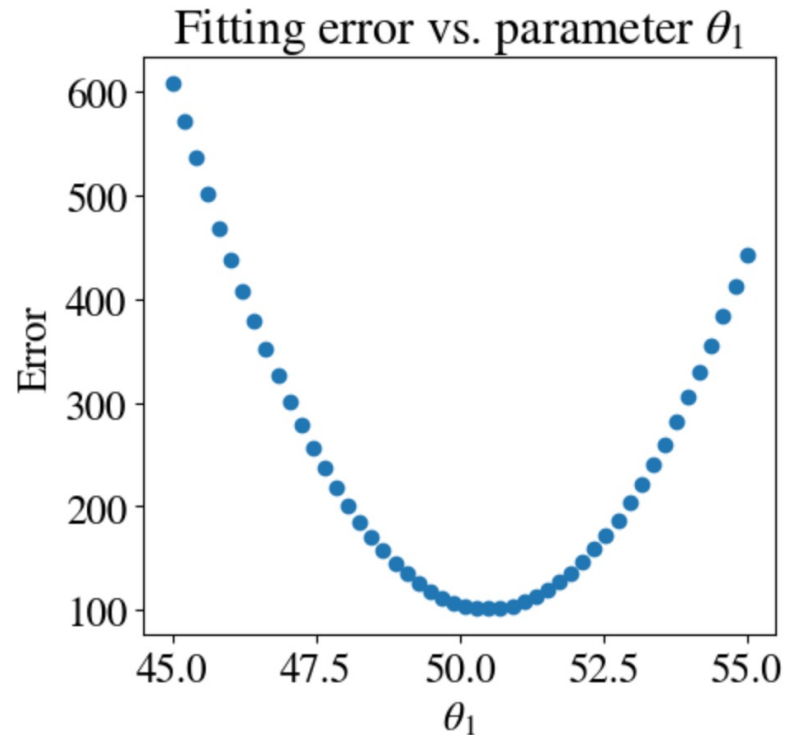
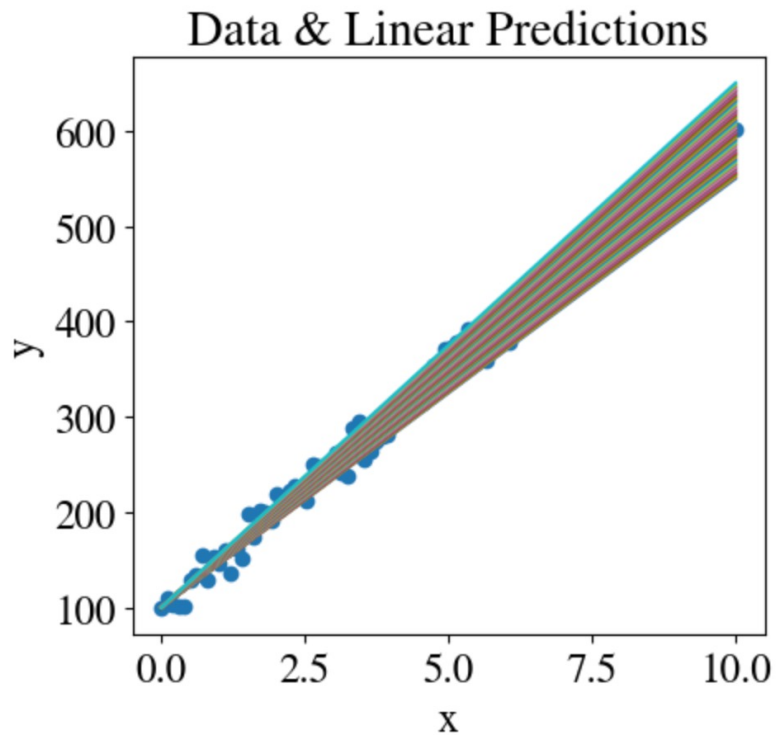
$$\begin{aligned}MSE &= \frac{1}{2m} \sum_{i=1}^m (y_{i,\text{pred}} - y_{i,\text{true}})^2 = \\&= \frac{1}{2m} \sum_{i=1}^m ((100 + \theta_1 x_i) - (100 + 50 x_i))^2 = \\&= (\theta_1 - 50)^2 \times \frac{1}{2m} \sum_{i=1}^m (x_i)^2 \approx 17(\theta_1 - 50)^2\end{aligned}$$

- How do we improve this?

Try narrower range

- $45 < \theta_1 < 55$, again in 50 steps

- Best fit: $\theta_1 = 50.5$
- **CPU time: 650 ms**
- **MSE = 101.6**



Brute force fit of slope and intercept

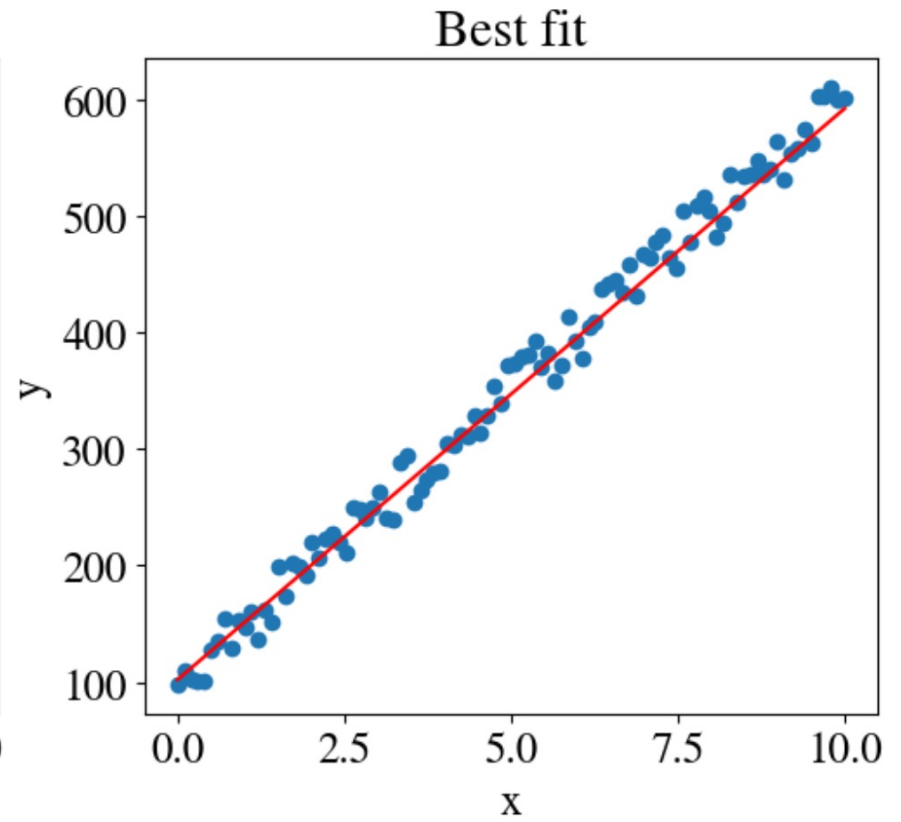
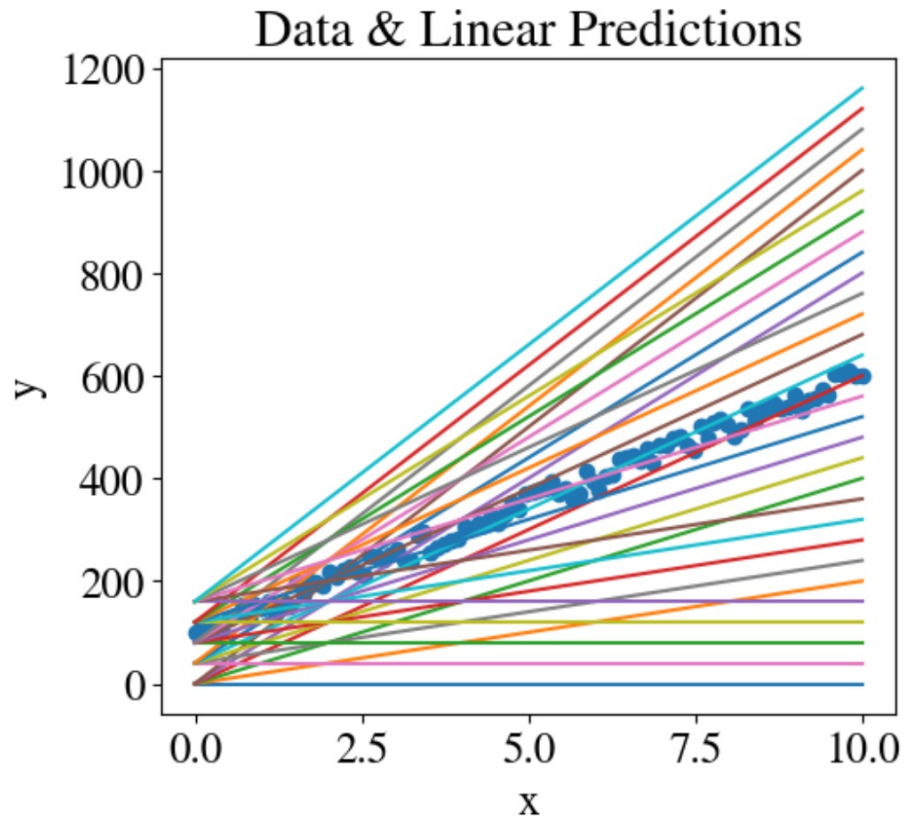
```
1  nrfits = 100
2
3  errors = np.zeros([nrfits,nrfits])
4  theta = np.zeros(nrfits)
5  theta0 = np.zeros(nrfits)
6
7  for j in range(nrfits):
8      theta0[j] = j * 200/nrfits
9      for i in range(nrfits):
10         theta[i] = i * 100/nrfits
11         y_pred = theta0[j] + theta[i] * x
12         errors[i,j] = mean_squared_error(y_pred , y)/2
```

Better:

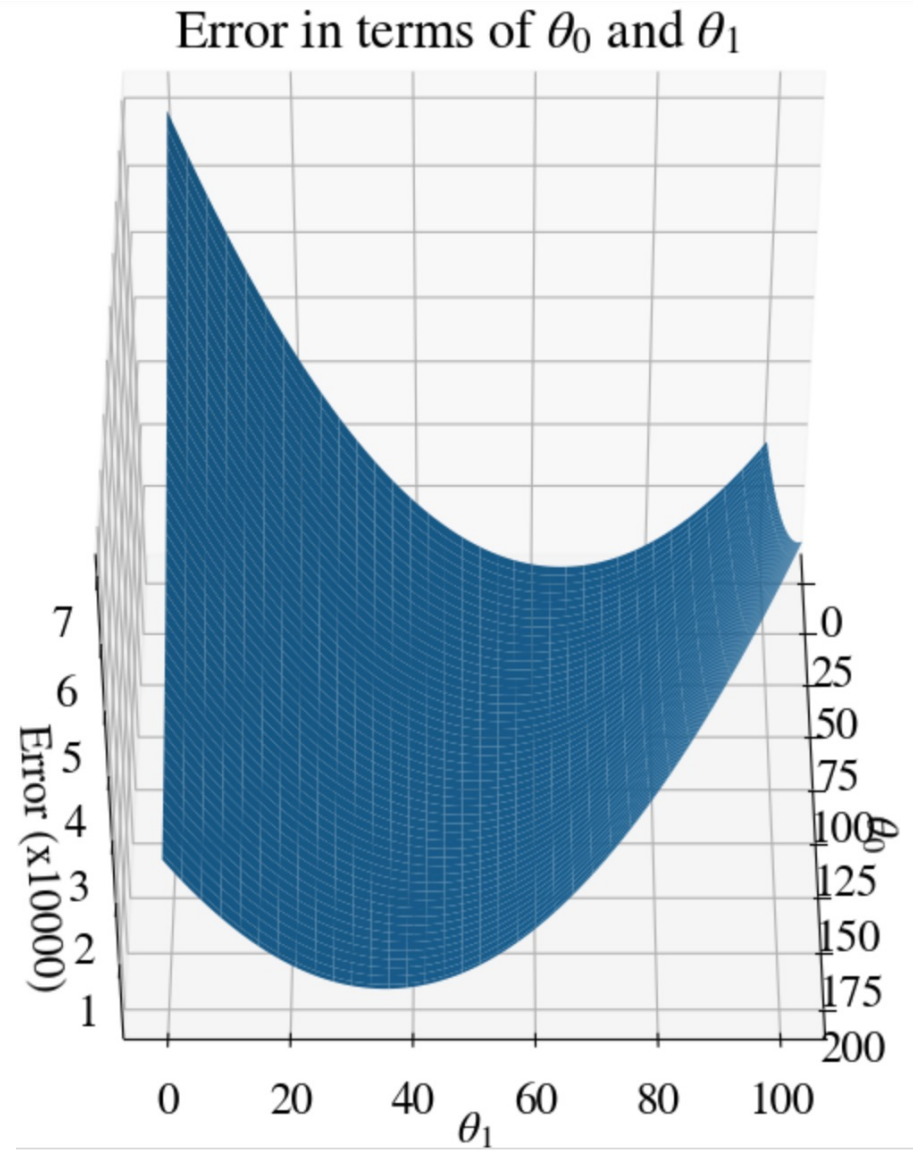
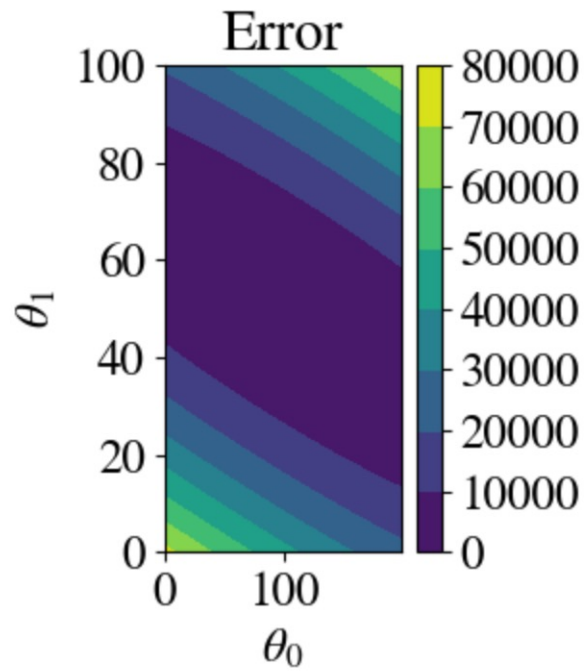
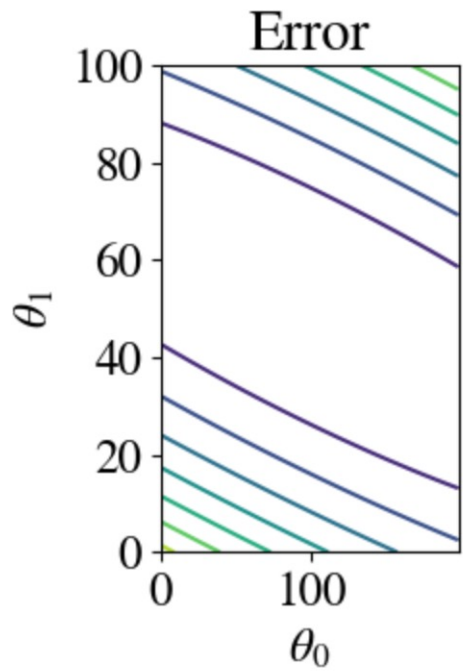
```
theta0 = np.linspace(0,200,nrfits)
theta = np.linspace(0,100,nrfits)
```

Result

- Best fit: $\theta_0 = 102$, and $\theta_1 = 49$
- 10,000 tries
- **CPU time: 5.1 s**
- **MSE = 102.2**



MSE in 2D

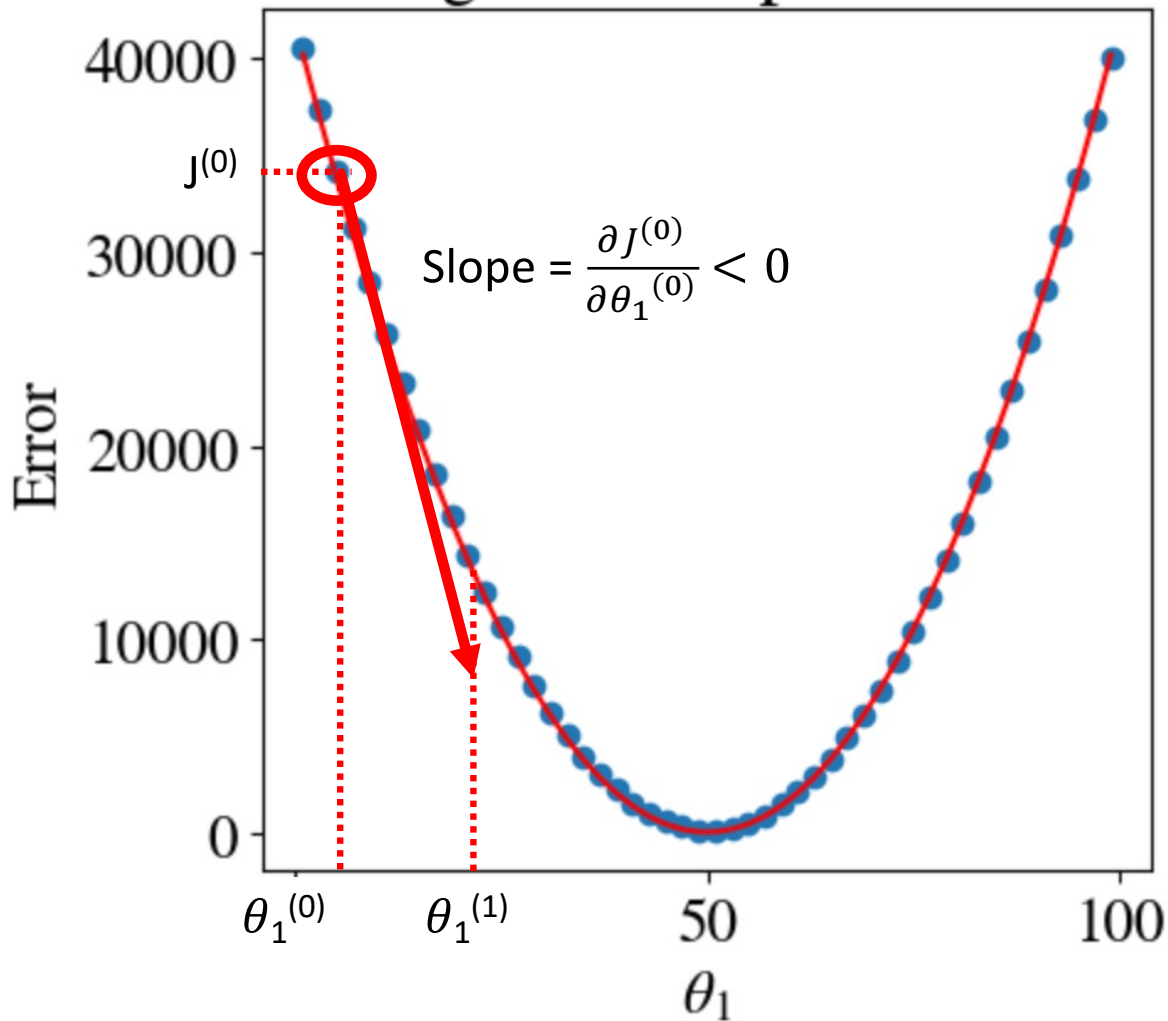


Gradient Descent Method

Navigating out of mountainous terrain



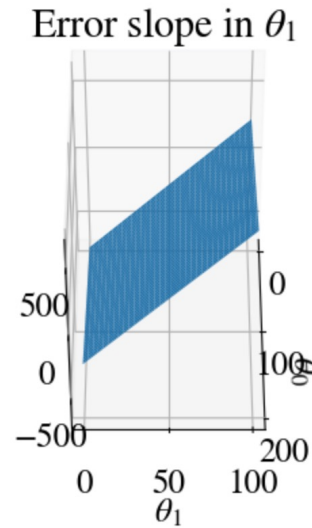
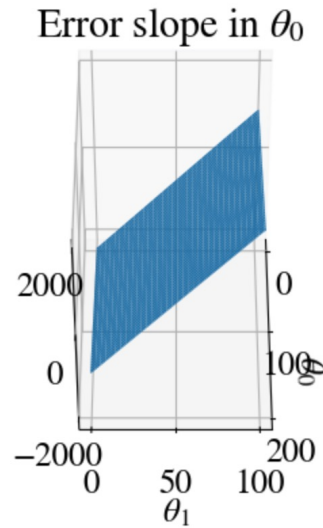
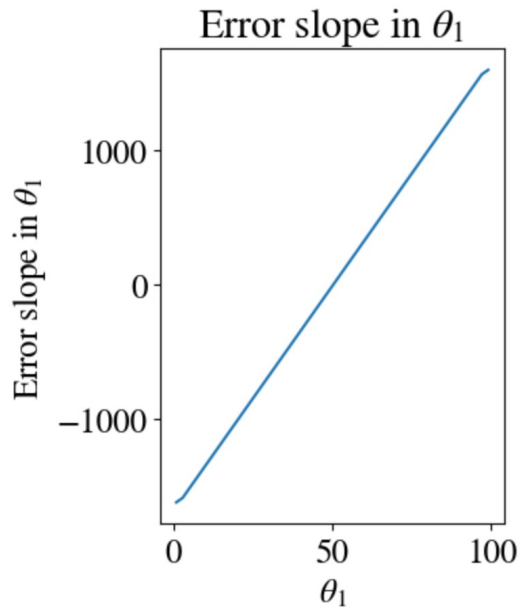
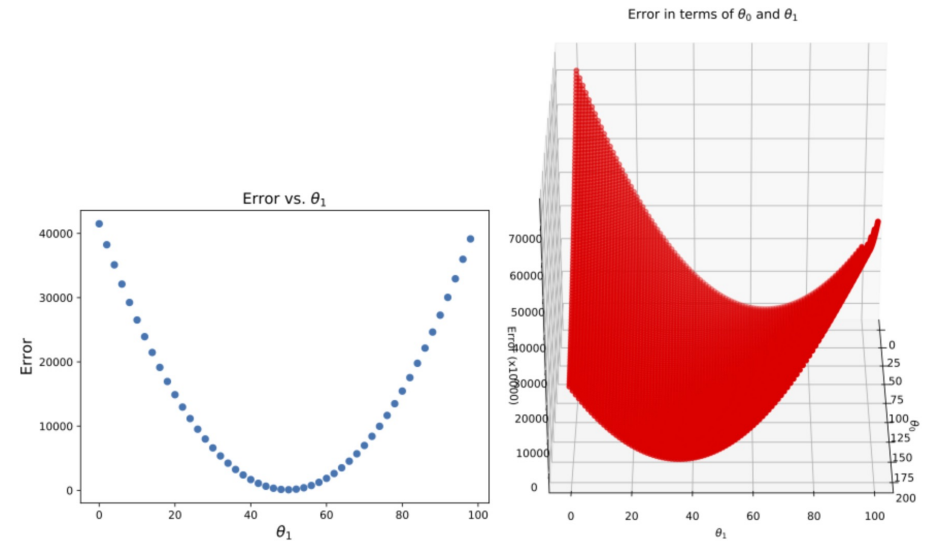
Fitting error vs. parameter θ_1



- Pick initial guess $\theta_1^{(0)}$,
- Calculate MSE, $J^{(0)}$,
- Calculate slope = **gradient**
 $\frac{\partial J^{(0)}}{\partial \theta_1^{(0)}}$,
- Pick new guess as:
$$\theta_1^{(1)} = \theta_1^{(0)} - \alpha \frac{\partial J^{(0)}}{\partial \theta_1^{(0)}}$$
- $J^{(1)} < J^{(0)}$ if α small enough
- α is the **learning rate**

Use derivative

- MSE is quadratic in θ_0 and θ_1 , so slope/gradient of errors is linear in θ_0 and θ_1



Explicit equations

$$J = \frac{1}{2m} \sum_{i=1}^m (h(x_i, \theta) - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2$$

- $\frac{\partial J}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i).$
- $\frac{\partial J}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)x_i.$

- $\theta_0^{\text{new}} = \theta_0 - \alpha \frac{\partial J}{\partial \theta_0} = \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i),$
- $\theta_1^{\text{new}} = \theta_1 - \alpha \frac{\partial J}{\partial \theta_1} = \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)x_i.$

Python code

```
1  %time
2  nrsteps = 10000
3  thetas = np.zeros([nrsteps,2])
4  errorsGD = np.zeros(nrsteps)
5
6  # Initial guess:
7  thetas[0,:] = [25,10]
8
9  # Step-size:
10 alphas = 0.01
11
12 for i in range(nrsteps-1):
13     hypothesis = thetas[i,0] + thetas[i,1]*x
14     errorsGD[i] = mean_squared_error(hypothesis,y) /2
15
16     thetas[i+1,0] = thetas[i,0] - (10*alphas/len(y)) * sum(hypothesis - y)
17     thetas[i+1,1] = thetas[i,1] - (alphas/len(y)) * sum((hypothesis - y)*x)
18
```

Still using many steps

Note: examples all start with 'bad' initial guess

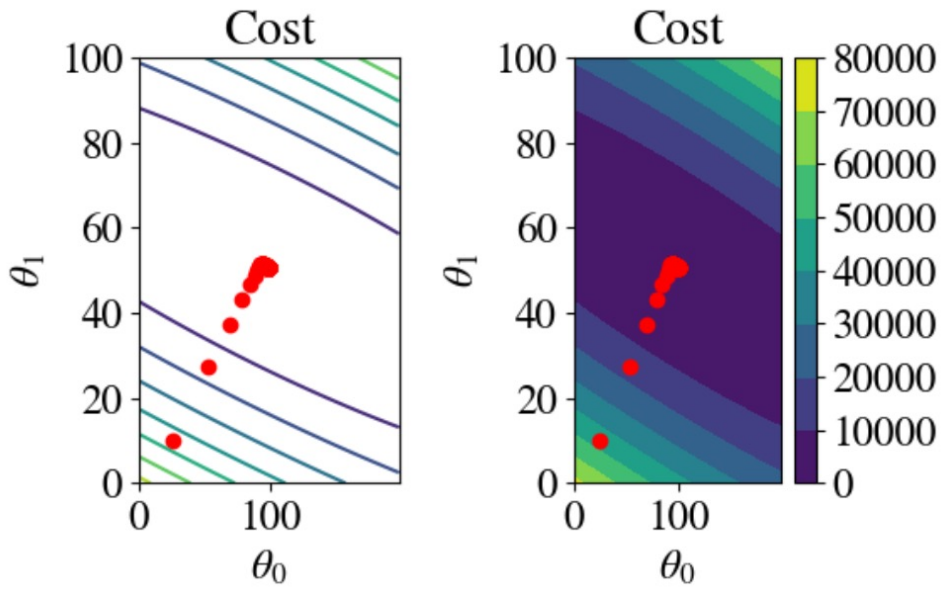
Learning rate (hyper-parameter)

Only 1 loop!

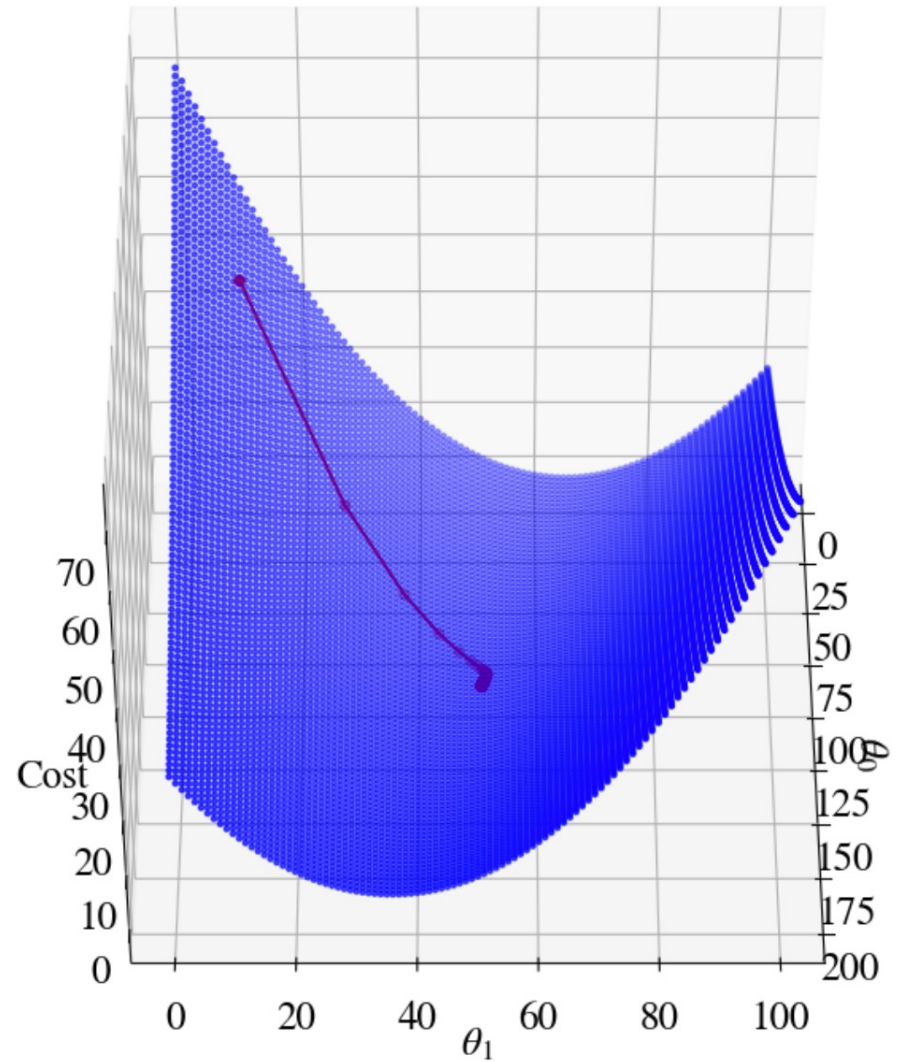
Cheating with 10x learning rate

How did we do?

- Best fit: $\theta_0 = 89.9$, and $\theta_1 = 50.7$
- 10,000 tries
- **CPU time: 3.6 s**
- **MSE = 101.4**

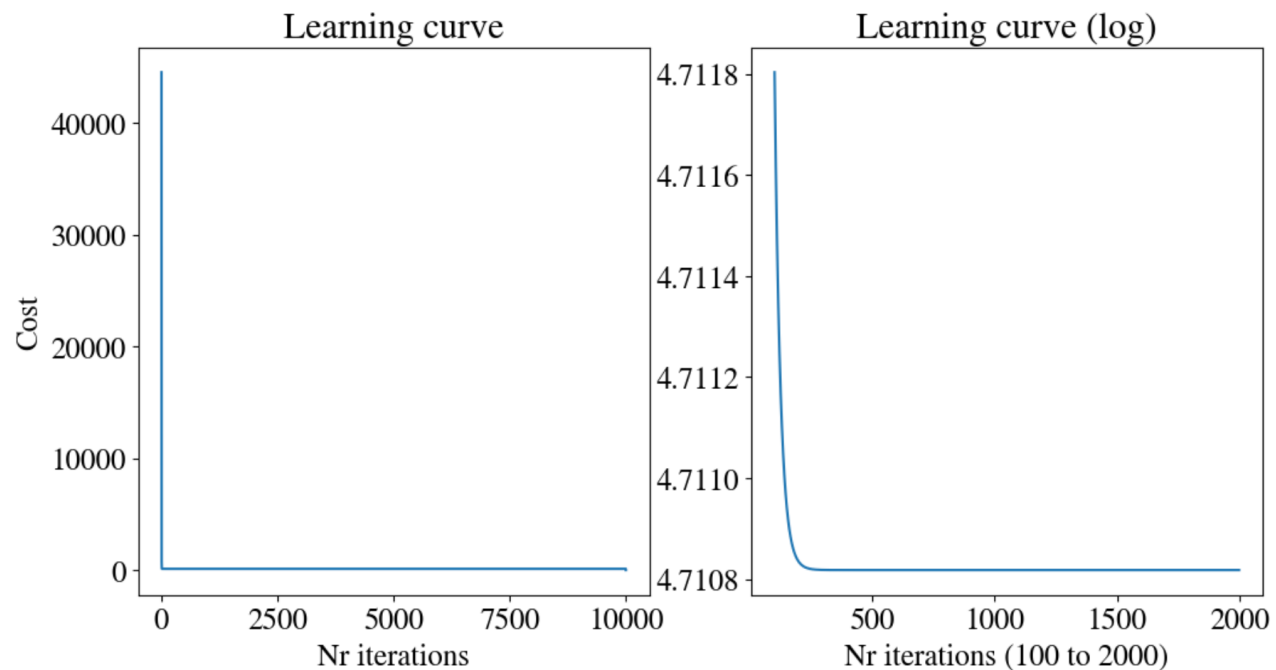


Demonstration of Gradient Descent Method



Learning curve

- Plot error as function of nr of iterations
- Clearly, best fit was already found after few iterations



Vectorization

- Add a column of all ones to the feature array \mathbf{x} ('bias')
- Define 2-element vector $\boldsymbol{\theta} = [\theta_0, \theta_1]$
- Then hypothesis (prediction) becomes simply:
- Vector $\mathbf{y} = h(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{x} \cdot \boldsymbol{\theta}$
- And 2D gradient is
- $\left(\frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}\right) = \nabla J = (h(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{y}) \cdot \mathbf{x} = (\mathbf{x} \cdot \boldsymbol{\theta} - \mathbf{y}) \cdot \mathbf{x}$

- $\frac{\partial J}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i).$
- $\frac{\partial J}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i) x_i.$

```
1 x_aug = np.zeros([m,2])
2 x_aug[:,0] = 1
3 x_aug[:,1] = x
4 x_aug
```

```
[ 1.          ,  3.93939394],
 [ 1.          ,  4.04040404],
 [ 1.          ,  4.14141414],
 [ 1.          ,  4.24242424],
 [ 1.          ,  4.34343434],
 [ 1.          ,  4.44444444],
 [ 1.          ,  4.54545455],
 [ 1.          ,  4.64646465],
 [ 1.          ,  4.74747475],
 [ 1.          ,  4.84848485],
 [ 1.          ,  4.94949495],
 [ 1.          ,  5.05050505],
 [ 1.          ,  5.15151515],
 [ 1.          ,  5.25252525],
 [ 1.          ,  5.35353535],
 [ 1.          ,  5.45454545],
 [ 1.          ,  5.55555556],
 [ 1.          ,  5.65656566],
 [ 1.          ,  5.75757576],
 [ 1.          ,  5.85858586],
 [ 1.          ,  5.95959596],
 [ 1.          ,  6.06060606],
 [ 1.          ,  6.16161616],
 [ 1.          ,  6.26262626],
 [ 1.          ,  6.36363636],
 [ 1.          ,  6.46464646],
 [ 1.          ,  6.56565657],
 [ 1.          ,  6.66666667],
 [ 1.          ,  6.76767677],
 [ 1.          ,  6.86868687],
 [ 1.          ,  6.96969697],
 [ 1.          ,  7.07070707],
 [ 1.          ,  7.17171717],
 [ 1.          ,  7.27272727],
 [ 1.          ,  7.37373737],
 [ 1.          ,  7.47474747],
 [ 1.          ,  7.57575758],
 [ 1.          ,  7.67676768],
```

Python code for non-vectorized GD

```
1  %time
2  nrsteps = 10000
3  thetas = np.zeros([nrsteps,2])
4  errorsGD = np.zeros(nrsteps)
5
6  # Initial guess:
7  thetas[0,:] = [25,10]
8
9  # Step-size:
10 alphas = 0.01
11
12 for i in range(nrsteps-1):
13     hypothesis = thetas[i,0] + thetas[i,1]*x
14     errorsGD[i] = mean_squared_error(hypothesis,y) /2
15
16     thetas[i+1,0] = thetas[i,0] - (10*alphas/len(y)) * sum(hypothesis - y)
17     thetas[i+1,1] = thetas[i,1] - (alphas/len(y)) * sum((hypothesis - y)*x)
18
```

Python code for vectorized gradient descent

```
1 %%time
2 nrsteps = 10000
3 m = len(y)
4 thetas = np.zeros([nrsteps,2])
5 errorsGD = np.zeros(nrsteps)
6
7 alpha_m = alphas/m
8
9 x_aug = np.zeros([m,2])
10 x_aug[:,0] = 1
11 x_aug[:,1] = x
12
13 # Initial guess:
14 thetas[0,:] = [25,10]
15
16 for i in range(nrsteps-1):
17     hypothesis = x_aug.dot(thetas[i,:])
18     errorsGD[i] = mean_squared_error(hypothesis,y) /2
19
20     thetas[i+1,:] = thetas[i,:] - alpha_m * (hypothesis - y).dot(x_aug)
```

One more optimization

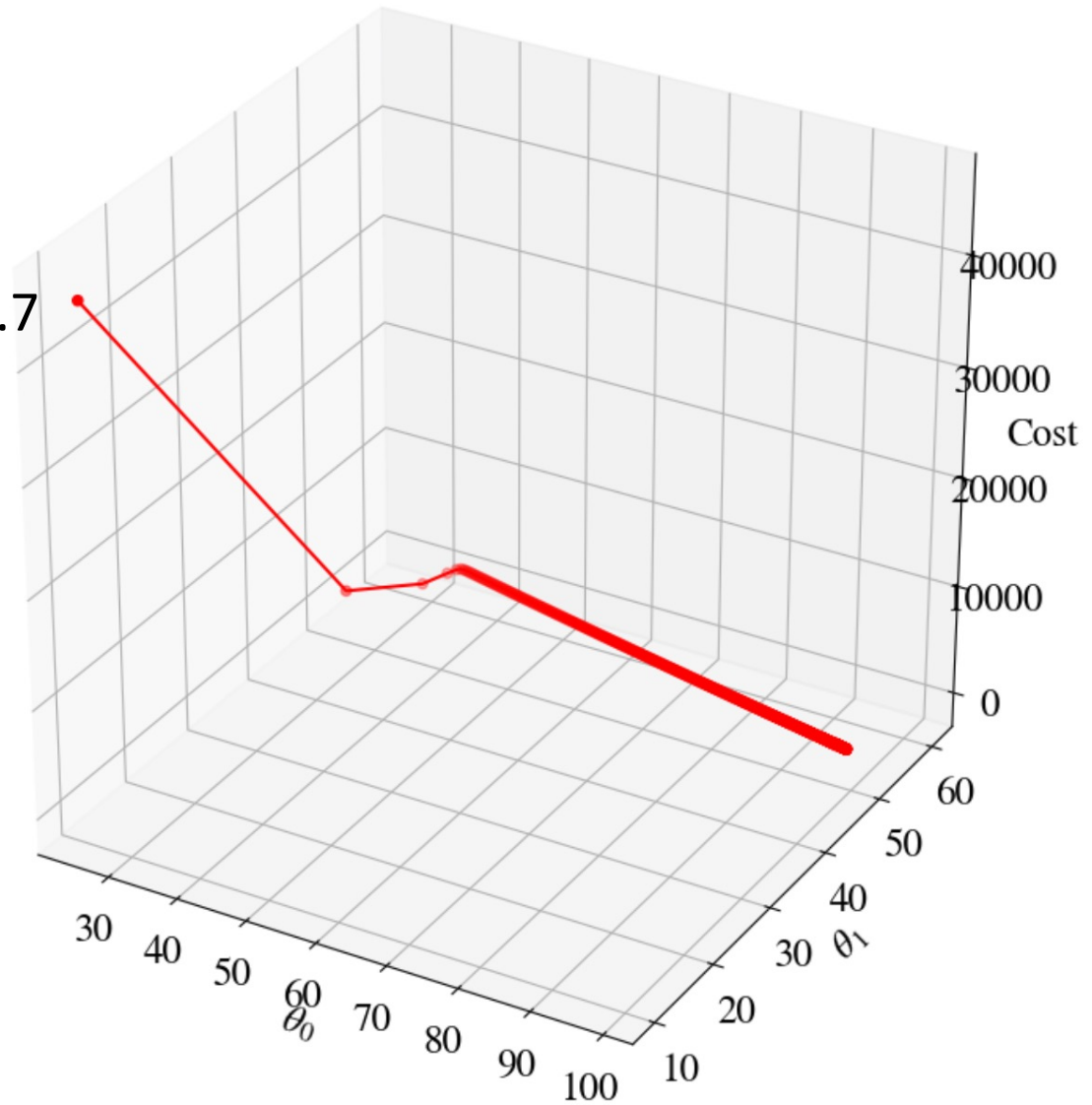
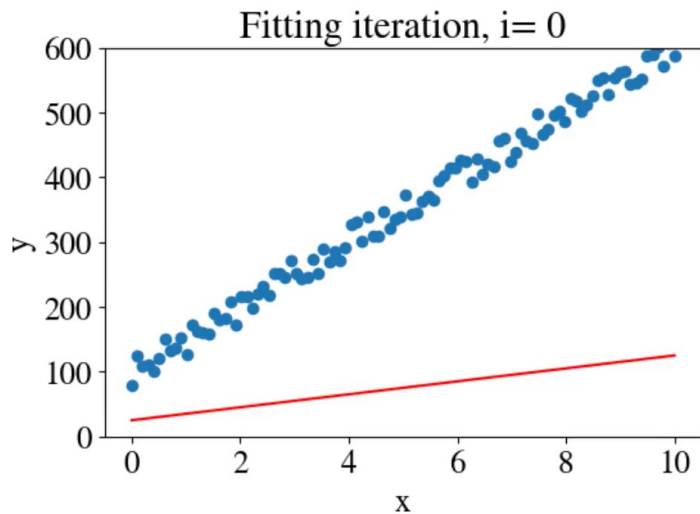
- Track changes, i.e., incremental reduction in MSE error, and stop when error is below certain tolerance ('while-loop'):

```
i=0
converged = False
tolerance = 1e-4
while not converged:
    hypothesis = x_aug.dot(thetas[i,:])
    errorsGD[i] = mean_squared_error(hypothesis,y) /2

    thetas[i+1,:] = thetas[i,:] - alpha_m * (hypothesis - y).dot(x_aug)
    if np.abs(errorsGD[i]-errorsGD[i-1]) < tolerance:
        converged = True
        print("Converged in", i, "iterations with thetas: ", thetas[i+1,:], "and error ", errorsGD[i])
    i+=1
```

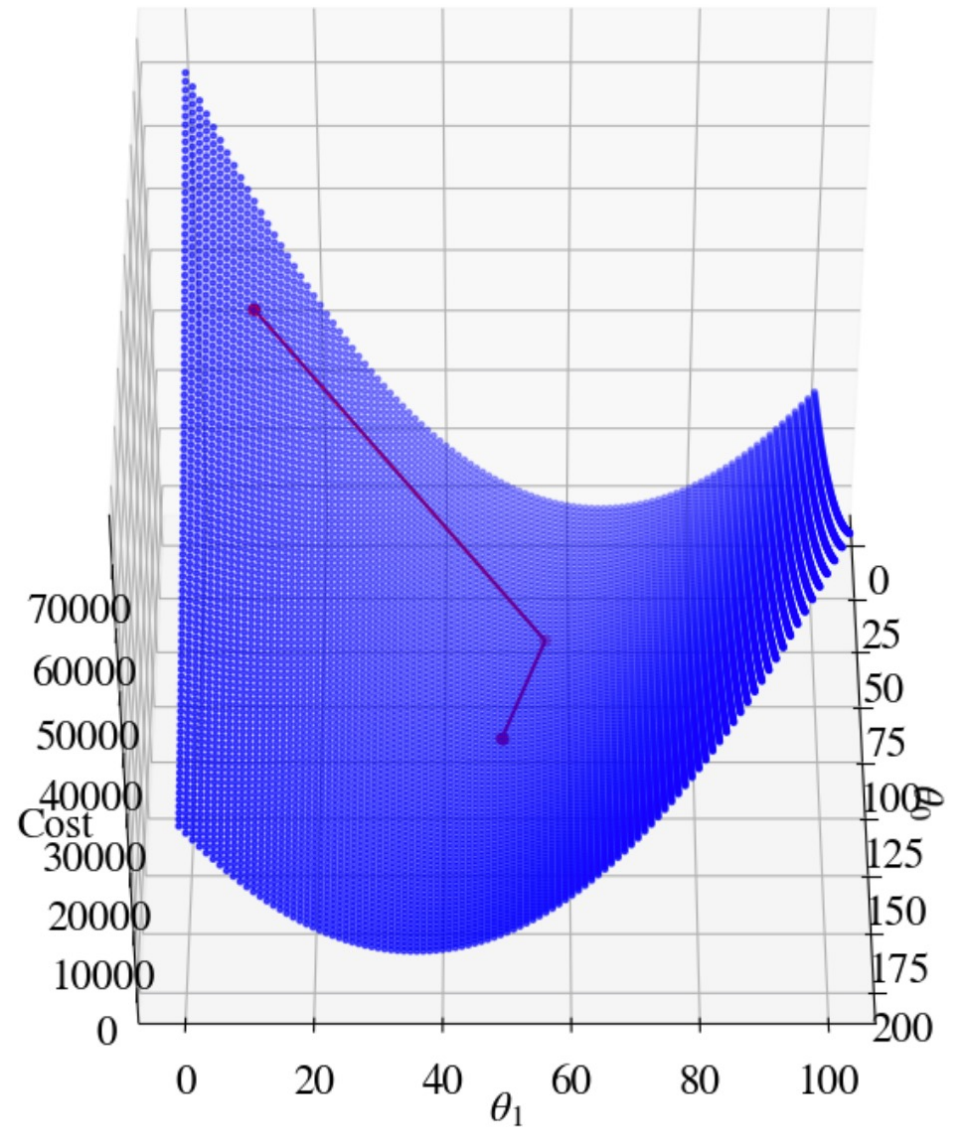
How did we do?

- Best fit: $\theta_0 = 98.7$, and $\theta_1 = 50.7$
- 10,000 tries => 1,099 tries
- **CPU time: 5.1 s => 248 ms**
- **MSE = 102.2 => 101.4**



Even faster!

- Use *second derivative* to determine *adaptive* learning rate α
- (Truncated) Newton method
- Best fit: $\theta_0 = 98.9$, and $\theta_1 = 50.7$
- 10,000 tries => **2 steps** (!!!)
- **CPU time: 5.1 s => 5 ms** (1/1000×)
- **MSE = 102.2 => 101.4**



Stochastic Gradient Descent

Stochastic Gradient Descent

- Many ML problems involve (many) thousands of features (columns of input data) and millions of data points (rows of input data)
- Gradient Descent can require too much RAM
- Key idea behind *stochastic* GD: Only fit to 1 point at a time!
 1. Initial guess of slope and intercept $\theta = [\theta_0^{(0)}, \theta_1^{(0)}]$
 2. Randomly pick a point $(x_i, y_{i,true})$ and predict target $y_{i,pred} = \theta_0^{(0)} + \theta_1^{(0)} x_i$
 3. Compute MSE error $J^{(0)}$ and its derivatives, only w.r.t. that single point
 4. Update fitting parameters to $\theta = [\theta_0^{(1)}, \theta_1^{(1)}]$
 5. Return to Step 2 and iterate/loop through all points x_i
 6. If accuracy is not satisfactory after looping through all points, set $[\theta_0^{(0)}, \theta_1^{(0)}]$ to last values and start again from Step 1

Python code I/II

- Randomly shuffle $(x_i, y_{i,true})$ pairs and set aside memory for θ and errors

```
7  p = np.random.permutation(len(xtmp))
8  xshuf = xtmp[p]
9  yshuf = ytmp[p]
10
11  m = len(yshuf)
12  nrouter = 10
13  nrsteps = m
14  # Allocate memory for variables:
15  savethetas = np.zeros([(nrsteps-1)*nrouter,2])
16  saveerrors = np.zeros((nrsteps-1)*nrouter)
17  thetas = np.zeros([nrsteps,2])
```

Python code II/II

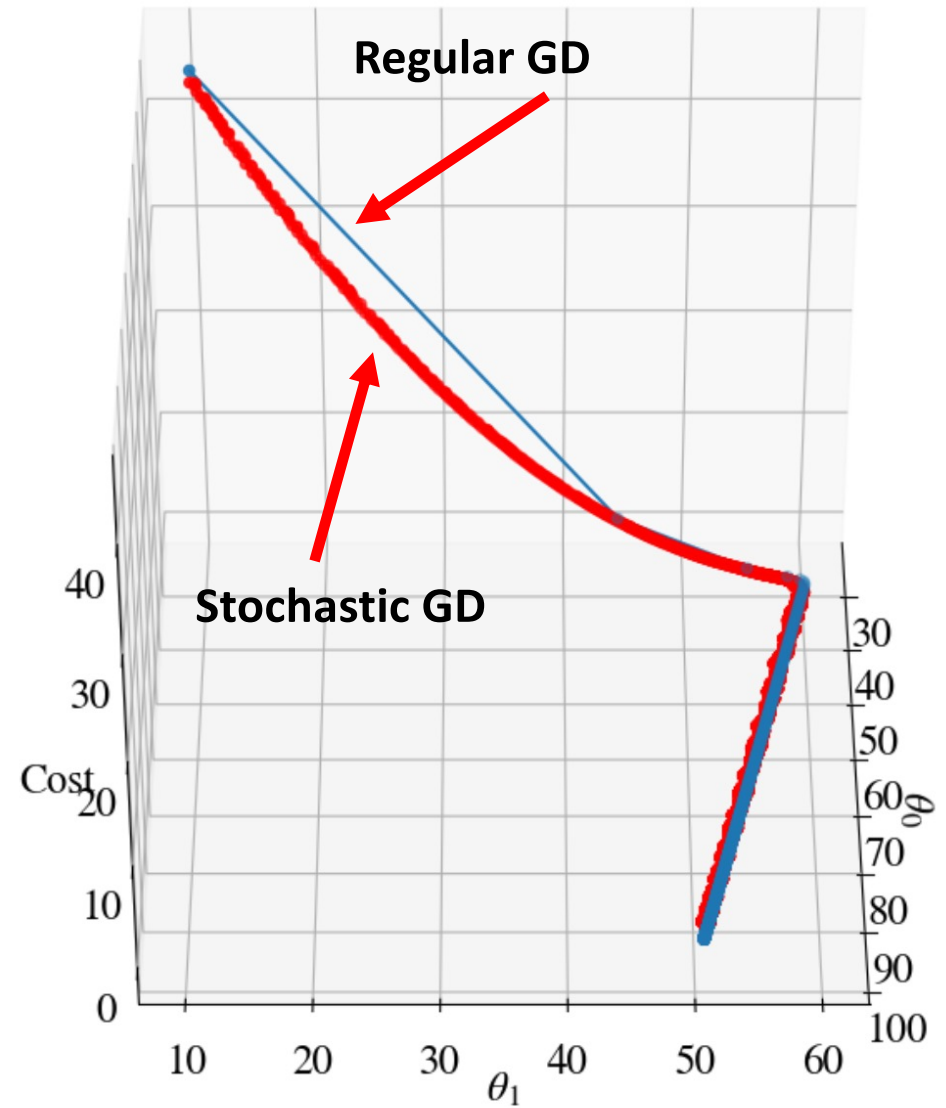
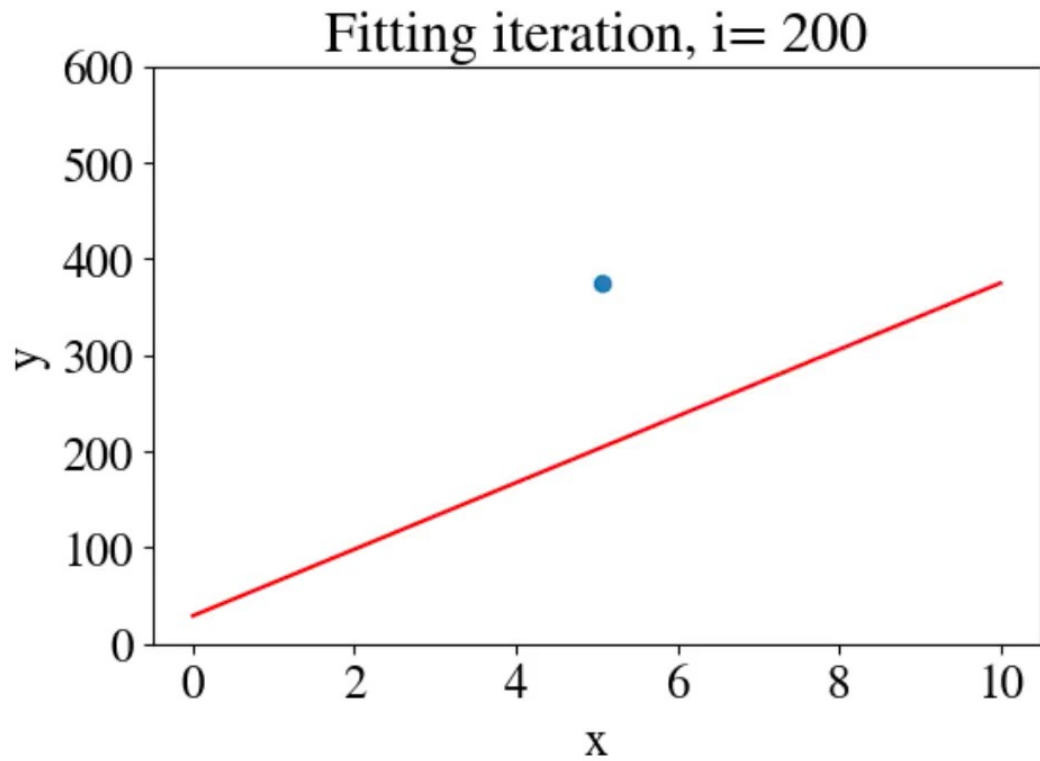
- Non-vectorized implementation:

```
1 # Initial guess:
2 thetas[0,:] = [25,10]
3 alpha = 1
4
5 cnt = 0
6 for j in range(nrouters):
7     for i in range(m-1):
8         thetas[i+1,0] = thetas[i,0] - (alpha/m) * (thetas[i,0] + thetas[i,1]*xshuf[i] - yshuf[i])
9         thetas[i+1,1] = thetas[i,1] - (alpha/m) * ((thetas[i,0] + thetas[i,1]*xshuf[i] - yshuf[i])*xshuf[i])
10
11         # Compute error:
12         saveerrors[cnt] = mean_squared_error(thetas[i+1,0]+thetas[i+1,1]*xshuf, yshuf) /2
13         cnt = cnt+1
14     # Set initial guess for the next run through all input values:
15     thetas[0,:] = thetas[m-1,:]
```

Only 1 point!



- Best fit: $\theta_0 = 94.4$, and $\theta_1 = 50.9$
- 100,000 tries ☹️
- **CPU time: 23 s** ☹️
- **MSE = 110.4** ☹️



Mini-Batch Gradient Descent

Mini-Batch Gradient Descent

- Best of both worlds:
- Instead of fitting to all points at once (regular gradient descent), or to a single point at a time (stochastic gradient descent), fit to small mini-batches in each iteration
- Example: for 1,000 data-points, loop through 100 iterations that each fit to 10 points
- Repeat if error still too large (using θ from previous pass as new initial guess)

Python code

- Vectorized

```
1 batch = 10
2 m = len(xshuf)
3 alpha_m = 0.01/batch
4
5 nrsteps = int(m/batch)
6 x_i = np.zeros([batch,2])
7 x_i[:,0] = 1
8 thetas = np.zeros([nrsteps+1,2])
9 thetas[0,:] = [25,10]
10
11 mb_iter=0
12 for j in range(1):
13     for i in range(nrsteps):
14         x_i[:,1] = xshuf[i*batch:(i+1)*batch]
15         y_i = yshuf[i*batch:(i+1)*batch]
16         hypothesis = x_i.dot(thetas[i,:])
17         thetas[i+1,:] = thetas[i,:] - alpha_m * (hypothesis - y_i).dot(x_i)
18         mb_iter += 1
19
20 thetas[0,:] = thetas[i+1,:]
```

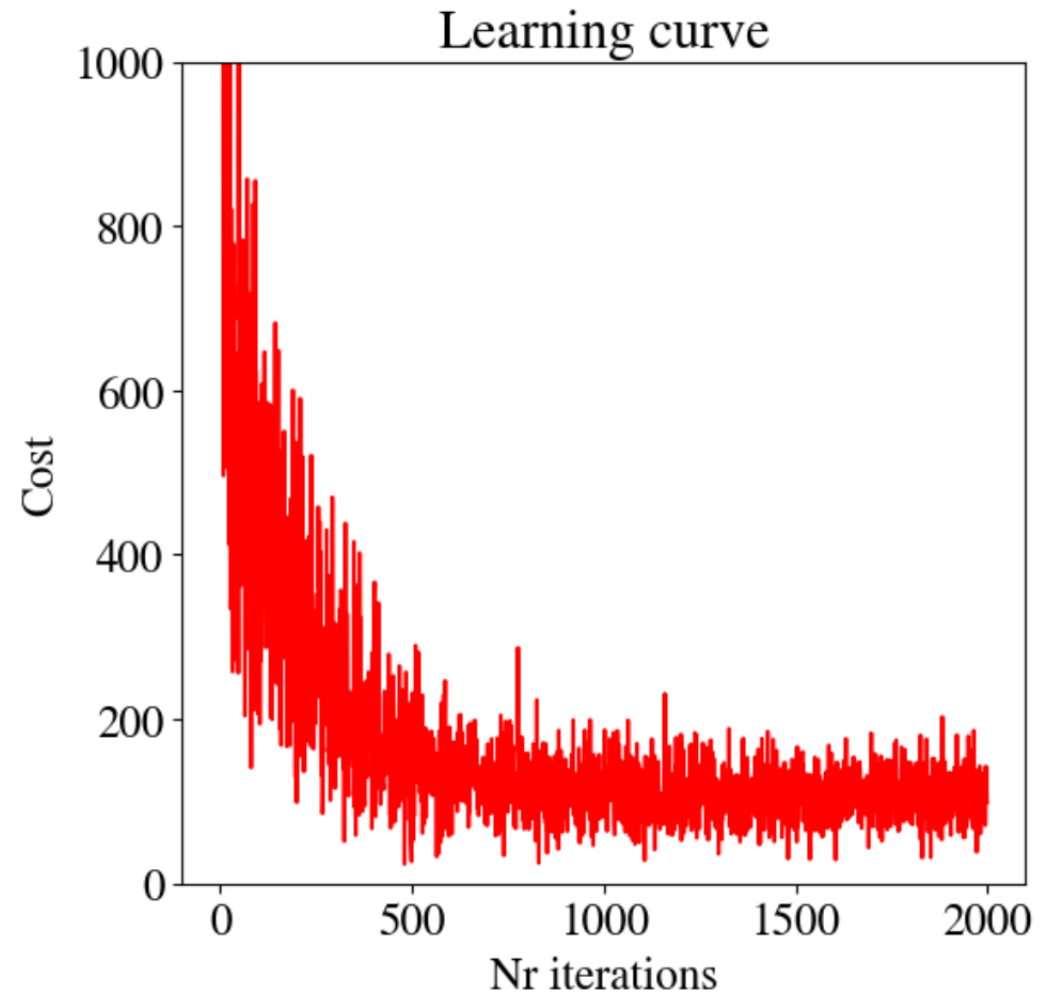
Batches of 10 points

Single pass

Vectorized

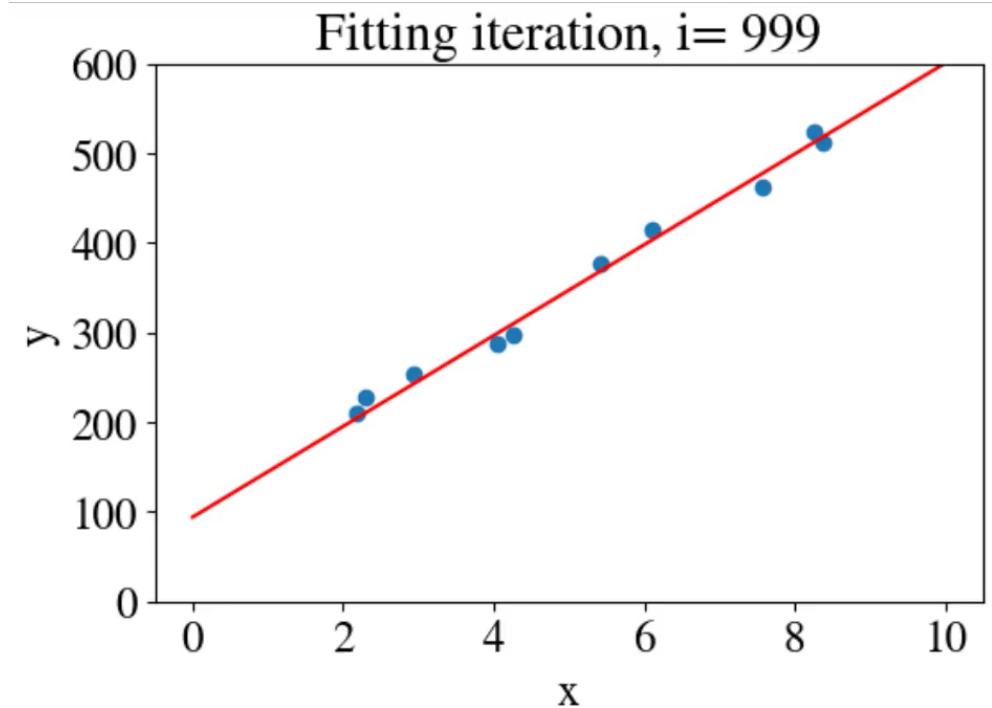
Learning curve

- Noisier than before
- Error for each batch of 10 points can be higher/lower, but still converge to global minimum



How did we do?

- Two passes:
- Best fit: $\theta_0 = 99.4$, and $\theta_1 = 49.8$
- 2,000 tries 😊
- **CPU time: 15.2 ms** 😊
- **MSE = 107.5** 😊
- ✓ Memory efficient
- ✓ CPU efficient (10× faster than vectorized regular GD)
- ✓ Accurate



Summary

Summary

- Linear regression fits linear **hypothesis** to labeled **features + targets**
- Double nested loop of trials for slope and intercept is inefficient
- **Gradient Descent methods** (GD) takes advantage of error derivatives: slopes of the error landscape
- Standard GD uses a user-defined **learning rate** that can be too small (slow convergence) or too large (overshooting)
- **Accelerated GD** uses second derivative as *adaptive* learning rate
- **Stochastic GD** fits 1 point at a time for less memory
- **(Mini-)Batch GD** fits, say, 10 points per batch: more efficient when memory allows