



Department of Computational and Data Sciences



# Linear Models: Regression

DS 392 JAN2022

Deepak Subramani

Assistant Professor

Dept. of Computational and Data Science

Indian Institute of Science Bengaluru



① Frame the problem



② Data Munging : Xarray  
pandas  
statistics

$$\underbrace{m \times n}_{\text{Dim Prod}} \rightarrow \underbrace{\text{Red } x}_{\text{d < n}} \underbrace{m \times d}_{\text{d < n}}$$

③ ML Modeling

$$\textcircled{1} h_{\theta} \rightarrow$$

$$\textcircled{2} \theta :$$

① Train  
validate  
Test

$$x \times m \times n \rightarrow x \times m \times d$$

$$d \ll n$$

④ Fine tune ; Ensemble

$$\textcircled{1} \text{ Obj Fn}$$

$$x \times m \times d \xrightarrow{\text{ML}} h_{\theta}(x \times m \times d; \theta) \rightarrow y$$

Performance  
Metric

⑤ Present

⑥ Deploy

- a)
- b) Monitor



Department of Computational and Data Sciences

# Regression Models



- Predict the value of a continuous variable in a supervised setting
- Use other variables as predictors

$$\hat{y} = h_{\theta}(X; \theta)$$

- Variable with a hat,  $\hat{y}$  is usually the prediction from a model
- $h$  is called the hypothesis, or the ML Model parametrized by  $\theta$
- $\theta$  can be calculated, estimated or *learned* from data by solving the optimization problem

$$\Rightarrow \min \sum_{j=1}^m (\hat{y}^{(j)} - y^{(j)})^2$$



$$y = 4x^2 - 3x + 1$$

# Linear Regression



- In Linear Regression model  $h_\theta$  is a linear function in  $\theta$

$$\hat{y} = \theta_1 x_1 + \theta_2 x_2 = \underline{X\theta}$$

param:  $\theta_1, \theta_2$

$$\Rightarrow L = \left[ \underbrace{\theta_1 x_1^{(1)} + \theta_2 x_2^{(1)}}_{\hat{y}^{(1)}} - y^{(1)} \right]^2 + \left[ \theta_1 x_1^{(2)} + \theta_2 x_2^{(2)} - y^{(2)} \right]^2 + \left[ \theta_1 x_1^{(3)} + \theta_2 x_2^{(3)} - y^{(3)} \right]^2$$

$$\min_{\theta_1, \theta_2} L$$

$$+ \left[ \theta_1 x_1^{(4)} + \theta_2 x_2^{(4)} - y^{(4)} \right]^2$$

$$\theta_1^*, \theta_2^*$$

$$\frac{\partial L}{\partial \theta_1} \Big|_{\theta_1^*} = 0 \quad \& \quad \frac{\partial L}{\partial \theta_2} \Big|_{\theta_2^*} = 0 \quad \text{Jacobian}$$

Hessian  
Convex

$$\hat{w} = \theta_1 \cdot H + \theta_2 \cdot A$$

$$X = \begin{bmatrix} x_1 & x_2 \\ \vdots & \vdots \\ x_m & x_n \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}_{n \times 1}$$

$m \times n$   
 $3 \times 2$





$$L = [\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} - y^{(1)}]^2$$

$$\frac{\partial L}{\partial \theta_1} = 2 [\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} - y^{(1)}] \cdot x_1^{(1)} \\ + 2 [\theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} - y^{(2)}] x_1^{(2)} \\ + 2 [\theta_0 x_0^{(3)} + \theta_1 x_1^{(3)} - y^{(3)}] x_1^{(3)} = 0$$

$$\frac{\partial L}{\partial \theta_2} = 2 [\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} - y^{(1)}] \cdot x_2^{(1)} \\ = 0$$

Set of linear equations.



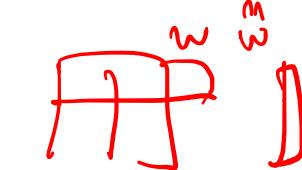
# Linear Regression

- In Linear Regression model  $h_\theta$  is a linear function in  $\theta$

$$\hat{y} = \underline{\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n} = X\theta$$

$$\Rightarrow \min \sum_{j=1}^m ([\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n]^{(j)} - y^{(j)})^2$$

$$\min(X\theta - y)^T (X\theta - y) \rightarrow S.S.E$$



$X: m \times n$

$5 \times 2$

$3 \times 3$

$$\hat{y} = X\theta$$

$$X\theta = \hat{y}$$

$$\boxed{\hat{y} = y}$$

$$\boxed{X\theta = y}$$

$$\Rightarrow \frac{\partial L}{\partial \theta} = 0$$

Normal Equation

$$X\theta = y$$

$$\boxed{\theta = X^{-1}y}$$

$$(X^T X)\theta = X^T y$$

$$\Rightarrow \theta = \underbrace{(X^T X)^{-1}}_{L.S.P} X^T y$$

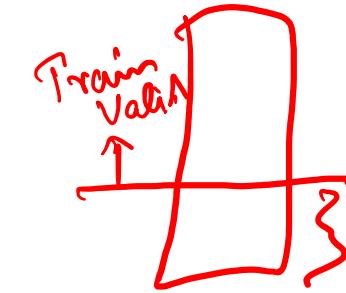
L.S.P

$$\underbrace{(X^T X)^{-1} X^T}_{\text{pseudo inverse}}$$

pseudo  
inverse



# Normal Equation



- The solution is the solution of the normal equation

$$\theta^* = \underbrace{(X^T X)^{-1}}_m \underbrace{X^T y}_n$$

- Solution:

- Direct Evaluation,
- Cholesky Decomposition,
- QR Factorization,
- SVD Method

} Matrix Factorization

$$X_{\text{test}} \theta^* = \hat{y}_{\text{test}}$$
$$y$$
$$L(y, \hat{y}_{\text{test}})$$

Code from scratch for  
Direct Evaluation

```
import numpy as np
X = 3 * np.random.rand(1000,1)
y = 2 + 4*X + np.random.randn(1000,1)
X_b = np.concatenate((np.ones((1000,1)),X),axis=1)
theta_star = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y)
theta_star
```



Department of Computational and Data Sciences

# Correlation is not Causation, But it helps Linear Models!



- While Data Cleaning, you will use correlation
- Predictors (columns of X) with high correlation to target (y) are good to make linear models
- What if X is correlated? - The issue is called multicollinearity
  - It is a problem if you want to identify contribution of one feature to the target
  - It is NOT a problem if you only want prediction accuracy



Department of Computational and Data Sciences

# Correlation is not Causation, But it helps Linear Models!



$$\boxed{X^T X \theta^* = X^T y}$$

- Let us see how correlation helps
- Remove mean from X and y [Standard Scaler] and see the normal eq.

$$E[(x - E[x])^2]$$

$$(X^T X) \theta^* = (X^T y); \frac{1}{m} (X^T X) \theta^* = \frac{1}{m} (X^T y);$$

$$\rightarrow \underbrace{C_{XX}}_{\sum n \times n} \theta^* = \underbrace{C_{XY}}_0$$

$$\sum n \times n \theta^* = n \times 1$$

$$\theta^* = C_{xy}$$

$$\{-1, 1\}$$

$$\theta_3 = D$$

$$\theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

p-value

Hypothesis Testing



# Scikit Learn

```
from sklearn.model_selection import train_test_split
```

```
from sklearn import linear_model
# Create linear regression object
lin_reg = linear_model.LinearRegression()
# Train the model using the training sets
lin_reg.fit(X_train, y_train)      SVD
# Make predictions using the testing set
y_pred = lin_reg.predict(X_test)
# The coefficients
print('Coefficients: \n', lin_reg.intercept_, lin_reg.coef_)
print('Prediction: \n', y_pred)
```

API → Sklearn  
Python → Scipy - nLS  
LAPACK → Fortran  
↓  
SVD

$(X^T X)^{-1} X^T y$



Department of Computational and Data Sciences



# lin\_reg.fit() – Behind the Scene

- Ordinary Least Squares (`scipy.linalg.lstsq`):
  - `scipy.linalg.lstsq` uses a LAPACK driver to solve the normal equations
  - Default is DGELSD – Double Precision, General Least Squares using SVD
  - DGELSY – Double Precision, General Least Squares using QR or LQ factorization
  - GE – General Matrices (Unsymmetric, no special property, convention in LAPACK)



# Computational Cost

- The Cholesky-based algorithm is particularly useful when  $m \gg n$  and it is practical to store  $X^T X$  but not  $X$  itself.
- It can also be less expensive than the alternatives when  $m \gg n$  and  $X$  is sparse.
- However, this approach must be modified when  $X$  is rank-deficient or ill conditioned to allow pivoting of the diagonal elements of  $X^T X$ .
- The QR approach avoids squaring of the condition number and hence may be more numerically robust, but often comes at an extra cost.
- While potentially the most expensive, the SVD approach is the most robust and reliable of all.



Department of Computational and Data Sciences

# Computational Cost



- Solving the normal equation by QR or SVD method is  $O(mn^2)$  – square in the number of features, and linear in the number of data points
- If we increase features by a factor 2, then computational cost increases by a factor 4.
  - Both the Normal Equation and the SVD approach get very slow when the number of features grows large (e.g., 100,000).
- As they are linear in the number of data points, they can handle large training sets efficiently, provided data can fit in memory.
- We need a better optimization algorithm when train data is large.



# Need for Iterative Algorithms

$m \times n$

- Large problem size – lots of feature vectors and lots of data
- We used squared loss for the accuracy, what if we have another accuracy metric?
- What about nonlinear accuracy metrics? Nonlinear function  $h()$
- Traditionally, nonlinear function (e.g., neural network) were trained with squared loss using a class of algorithms known as nonlinear least squares
- Nonlinear least squares algorithm Gauss-Newton and Levenberg Marquardt algorithms were very popular to train shallow neural network models
- Now we use gradient descent algorithms which is what is called as a “first-order method”.
- Let us see these iterative algorithms now.

$\times \theta$   
 $\times \sin(\theta)$



Department of Computational and Data Sciences

# Gradient Descent - Intuition



- Suppose you are lost in the mountains in a dense fog, and you can only feel the slope of the ground below your feet.
- A good strategy to get to the bottom of the valley quickly is to go downhill in the direction of the steepest slope.
- This is exactly what Gradient Descent does: it measures the local gradient of the error function with regard to the parameter vector  $\theta$ , and it goes in the direction of descending gradient.
- Once the gradient is zero, you have reached a minimum!



 $\theta$ 

$$\frac{\partial L}{\partial \theta} = 0$$

$$\Rightarrow \theta_{(k+1)} = \underbrace{\theta_{(k)}}_{\uparrow} - \underbrace{\eta}_{\circlearrowleft} \nabla L \Big|_{\theta(k)}$$

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial \theta_1}, \\ \frac{\partial L}{\partial \theta_2}, \\ \vdots \end{bmatrix}$$

 $\theta$ : param $\eta$ : hyper param

tolerance

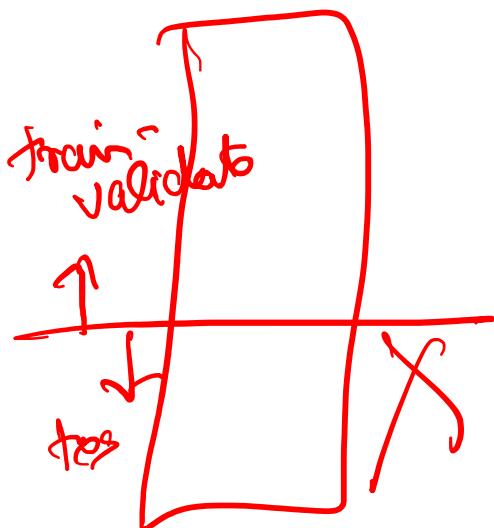
max\_iter

$\{\eta$ : Learning Rate  
 $\text{tol}$   
 $\text{max\_iter}$

$$\theta^* = \theta_{99}$$

$$\underline{DL} \approx 0 \quad \underline{\text{tol}}$$

$$\frac{\theta_{(k+1)} - \theta_{(k)}}{\eta}$$





# Gradient Descent - Variants

- Batch Gradient Descent

- Use the full data and calculate the Jacobian of the objective function

$$\nabla f(\theta) = \sum_{j=1}^m r_j(\theta) \nabla r_j(\theta) = J(\theta)^T r(\theta) = X^T (X\theta - y)$$

- Iteration of the optimization algorithm

$$\theta_{(k+1)} = \theta_{(k)} - \eta_{(k)} \nabla f(\theta_k)$$

$\eta_{(k)}$  → Learning schedule

## ~~Stochastic~~ Stochastic Gradient Descent

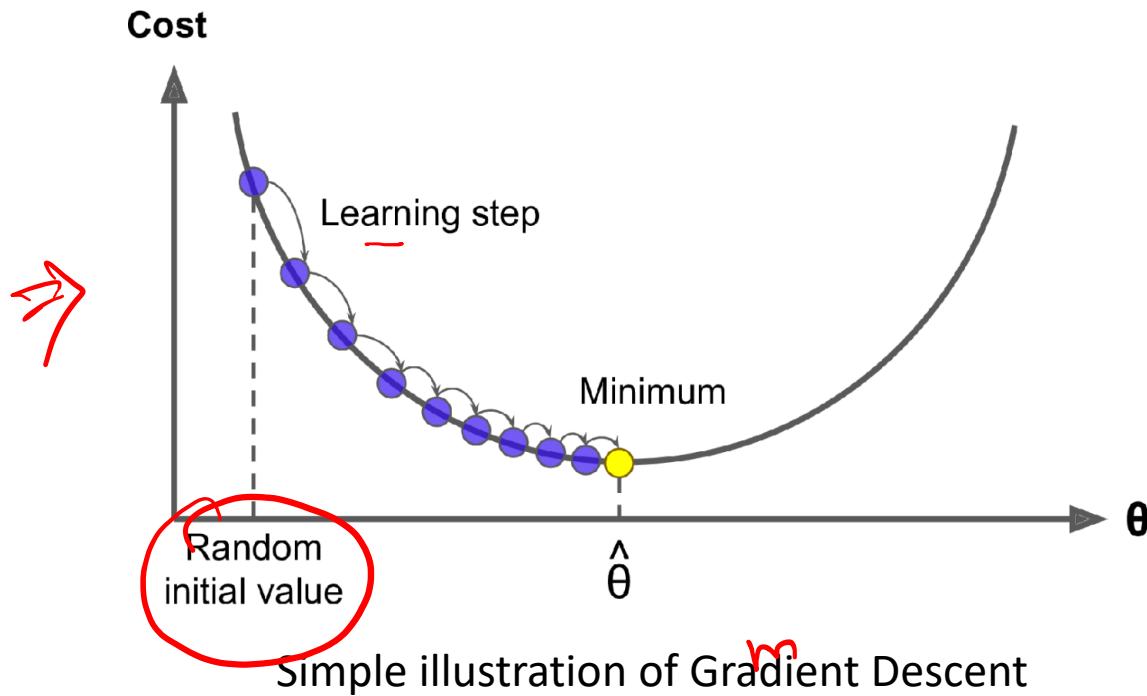
- Use only one data point at a time.
  - Usually m iterations is called an epoch

$$\nabla f = \sum_{j=1}^m \rightarrow \text{Full}$$

$$\nabla f = \sum_{j=j_{rand}}^m$$



# Gradient Descent



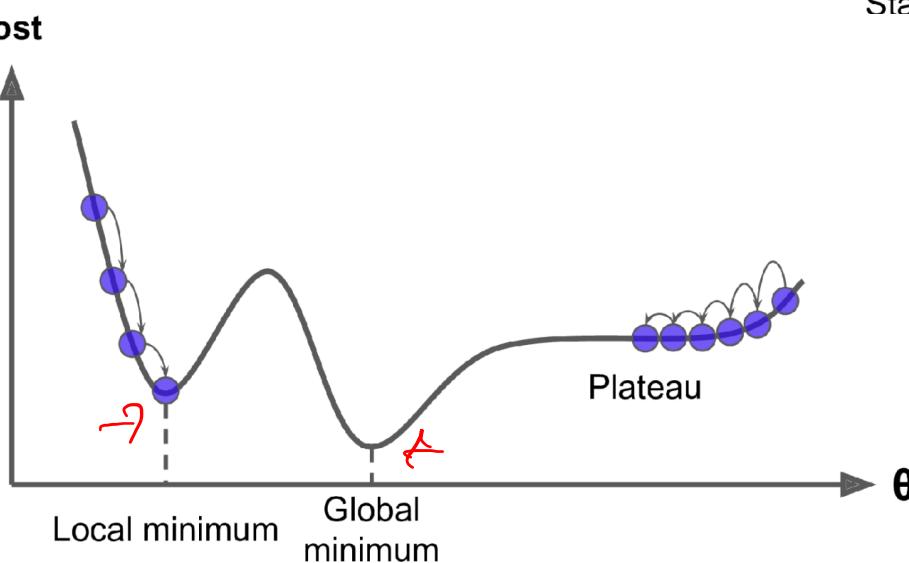
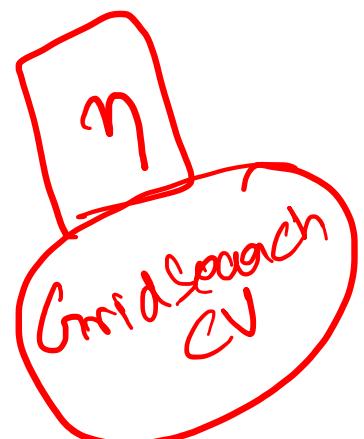
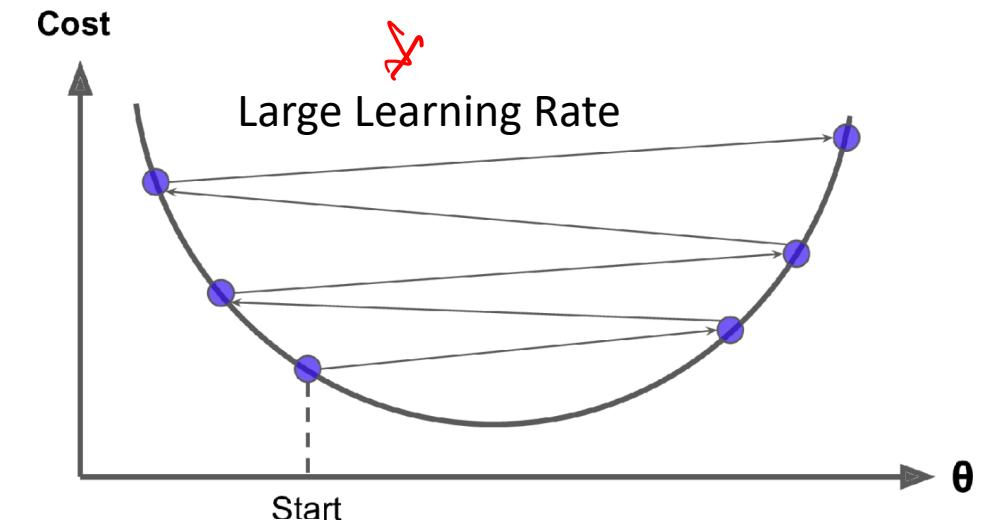
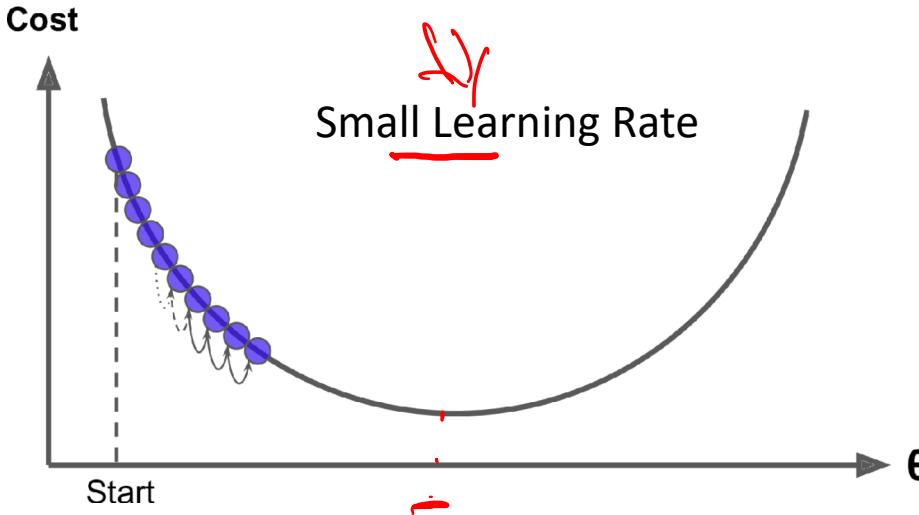
## Optimization vs Learning:

- In Optimization, we have to compute the step length accurately
- In learning, we don't care as much about being exactly optimal at every iteration, so we choose a learning rate  $\eta_k$  as a hyperparameter, instead of an exact step length
- Define Hyperparameter – They are parameters of the training/model that are not trainable with data directly, but are chosen by the designer/developer/user using some other information source (usually, validation)



Department of Computational and Data Sciences

# Gradient Descent



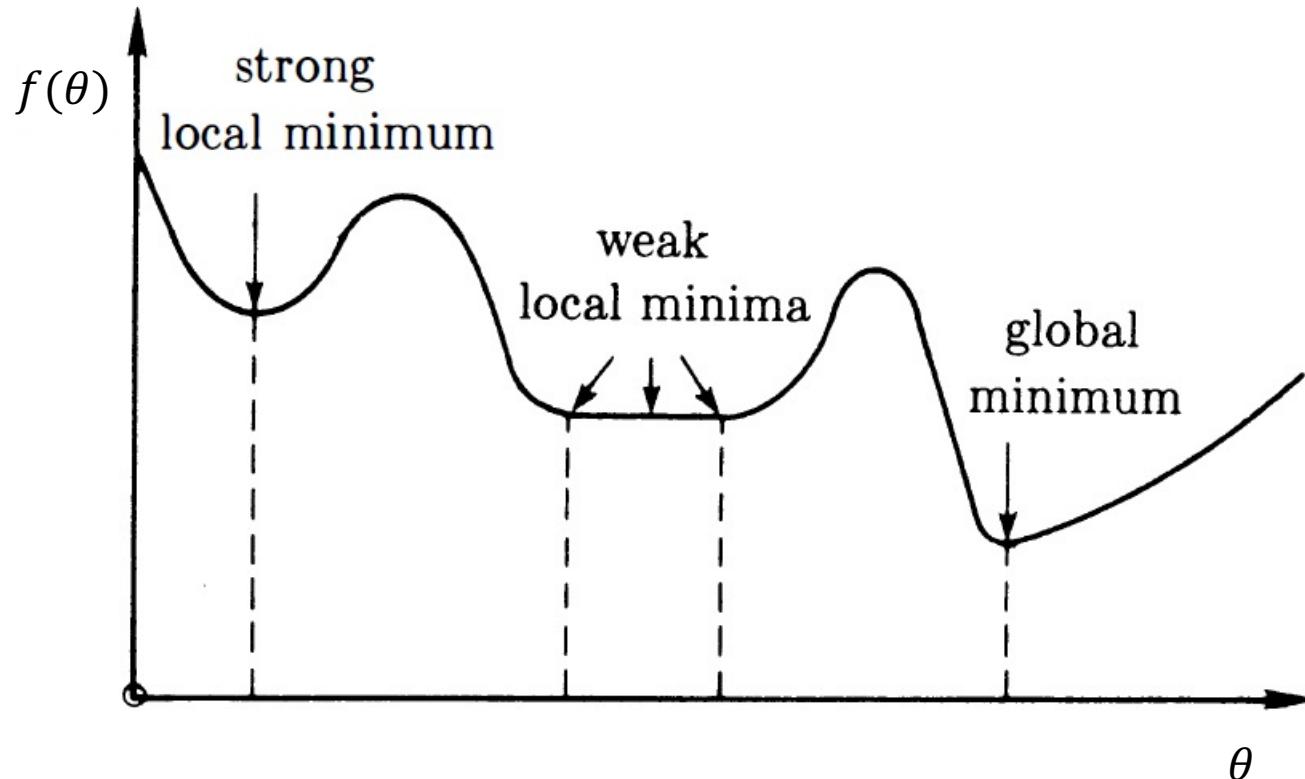
- Gradient Descent may miss global minima
- Very slow in a Plateau
- Problem Scaling Issues
- To ensure fast convergence, use scaled features



Department of Computational and Data Sciences



# Example of Non-Convex Function



**Figure 3b.** Examples of minima in the univariate case.



Department of Computational and Data Sciences

# Scikit Learn – SGDRegressor



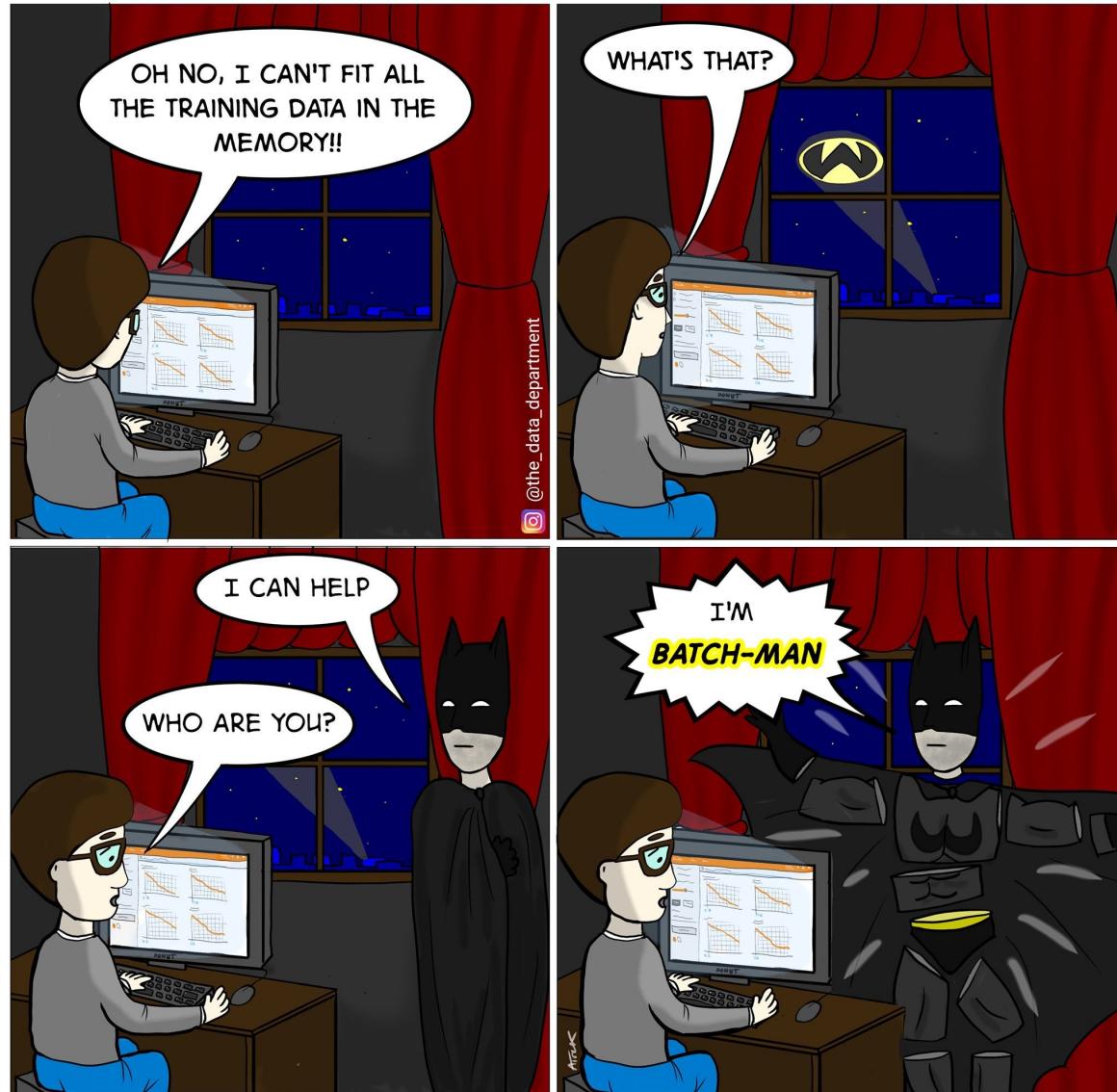
```
from sklearn.linear_model import SGDRegressor  
sgd_reg = SGDRegressor(max_iter=1000, tol=1e-3, penalty=None,  
eta0=0.1)  
sgd_reg.fit(X_train,y_train.ravel())
```

```
print(sgd_reg.intercept_, sgd_reg.coef_)
```

Note: .ravel() method of numpy array is used to make y\_train a 1-d array which is expected by SGDRegressor  
penalty is used to specify regularization which we will discuss shortly



Department of Computational and Data Sciences



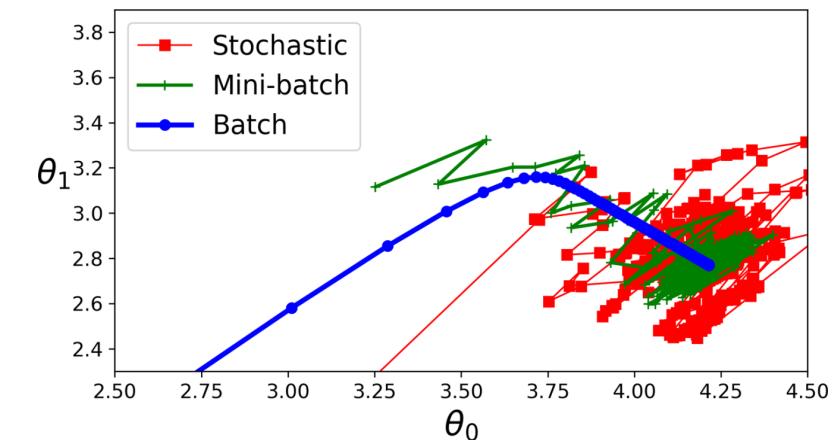


Department of Computational and Data Sciences

# Mini-batch Gradient Descent



- Mini-batch GD is in between Batch GD (Full data set at once) and Stochastic GD (one instance at once)
- Mini-batch computes the gradient on small random sets of instances called mini-batches
- The main advantage of Mini-batch GD over Stochastic GD is that you can get a performance boost from hardware optimization of matrix operations, especially when using GPUs.
- Min-batch SGD is less erratic than SGD





# Comparison of algorithms for Linear Regression

Algorithm	Large $m$	Out-of-core support	Large $n$	Hyperparams	Scaling required	Scikit-Learn
Normal Equation	Fast	No	Slow	0	No	N/A
SVD	Fast	No	Slow	0	No	LinearRegression
Batch GD	Slow	No	Fast	2	Yes	SGDRegressor
Stochastic GD	Fast	Yes	Fast	$\geq 2$	Yes	SGDRegressor
Mini-batch GD	Fast	Yes	Fast	$\geq 2$	Yes	SGDRegressor

Geron Chapter 4, Page 128, Table 4-1

Out of core algorithms are those that can process data that are too large to fit into a computer's main memory at once.



# Audience Poll

1. If there are  $m$  train data points, how many data points does batch GD use?
  - i.  $m$ ,
  - ii. 1,
  - iii.  $m/2$ ,
  - iv. one pre-determined fraction of  $m$
2. Which of the following is TRUE
  - i. Batch GD is a stochastic algorithm - FALSE
  - ii. When SGD sees all points once, it is called an epoch - TRUE
  - iii. Mini Batch GD is a deterministic algorithm - FALSE
  - iv. Batch GD uses one data point at a time - FALSE
3. Which of the following are the hyperparameters of Batch GD
  - i. Learning Rate, ✓
  - ii. Number of Iterations, ✓
  - iii. Size of Batch, ✗
  - iv. Sampling Strategy ✗

$$\left[ \begin{array}{c} \downarrow \\ \vdots \end{array} \right]^2 + \left[ \begin{array}{c} \downarrow \\ \vdots \end{array} \right]^2 + \left[ \begin{array}{c} \downarrow \\ \vdots \end{array} \right]^2$$