



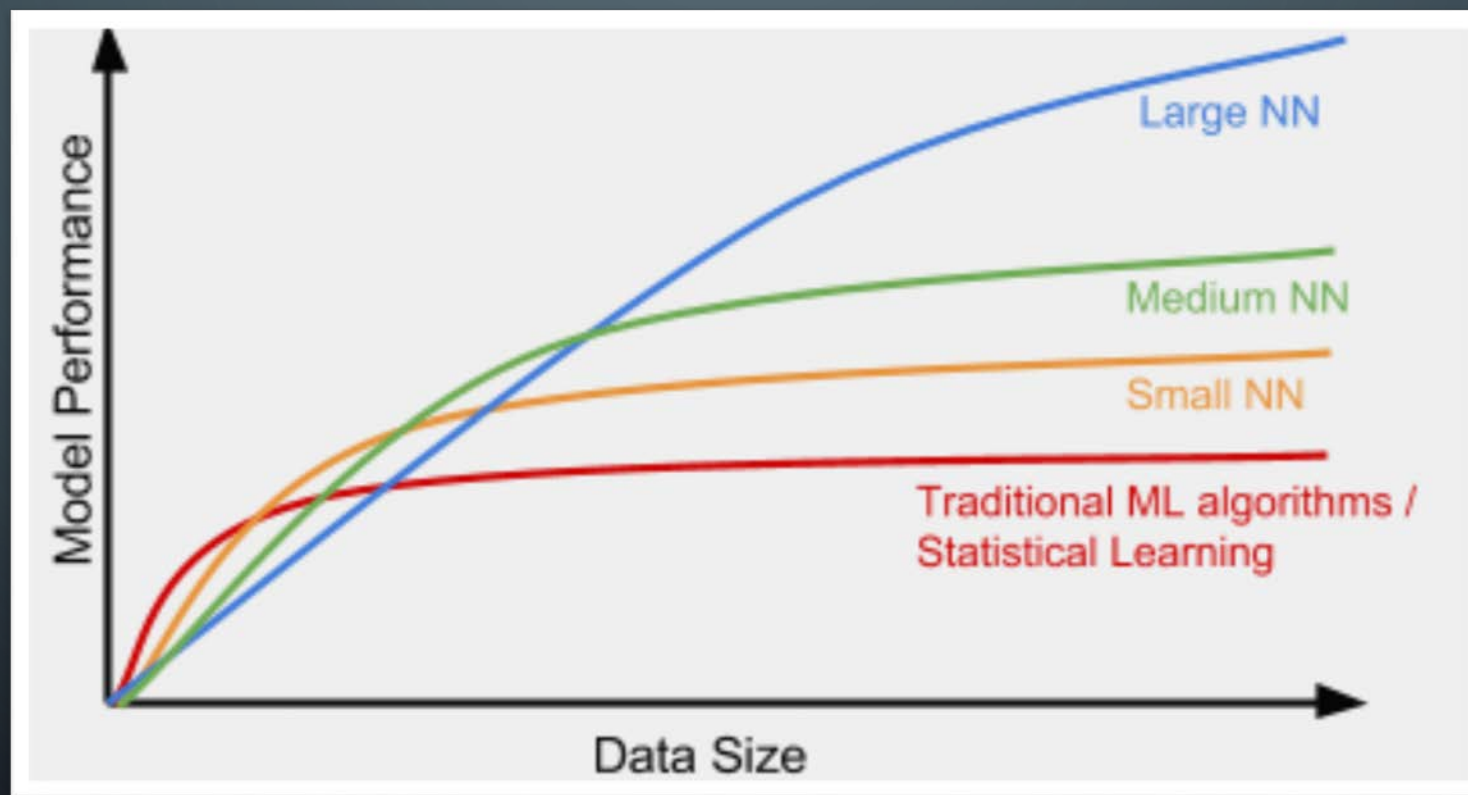
MACHINE LEARNING PROBLEMS

MOHAMMAD GHODDOSI

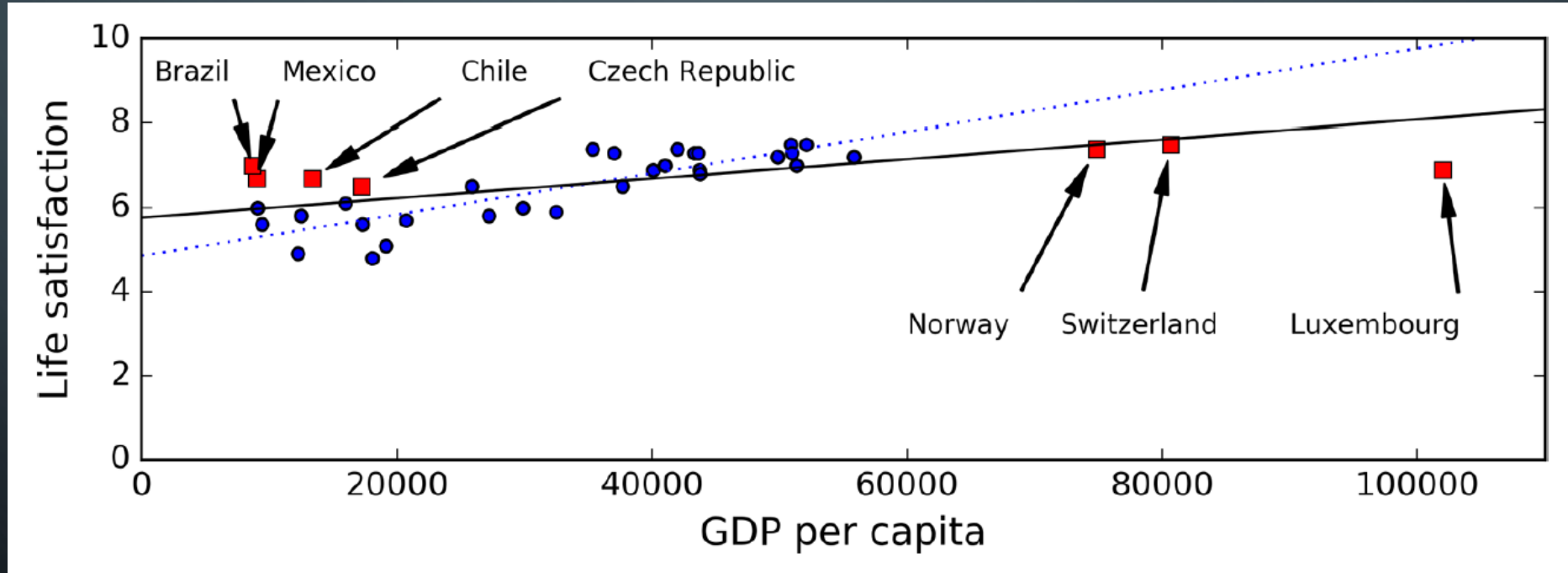
MACHINE LEARNING MAIN CHALLENGES

- Insufficient Quantity of Training Data
- Nonrepresentative data
- Poor-quality data
- Irrelevant feature
- Overfitting
- Underfitting

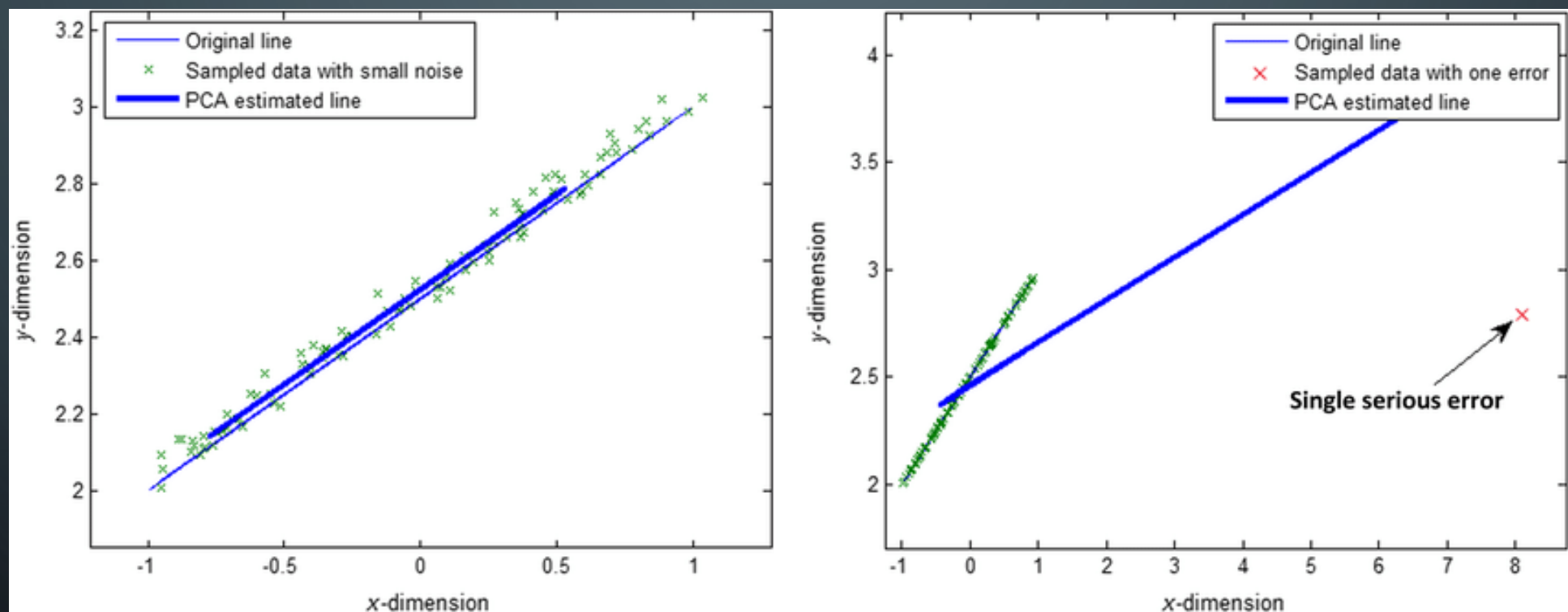
DATA SIZE



NONREPRESENTATIVE DATA



POOR-QUALITY DATA



HANDLING DATA QUALITY

- Noise
- Outlier
- Missing value

MISSING VALUES

- Get rid of the row
 - `df.dropna()`
- Get rid of the attribute
 - `df.drop(attribute)`
- Set missing values to some value
 - `sklearn.preprocessing.Imputer()`

SET MISSING VALUES TO SOME VALUE



- Out of range value (label as missing)
- Mean
- Median
- Machine learning models

IRRELEVANT FEATURE

- Garbage in, garbage out
- Feature selection
- Feature extraction
- Gathering new features

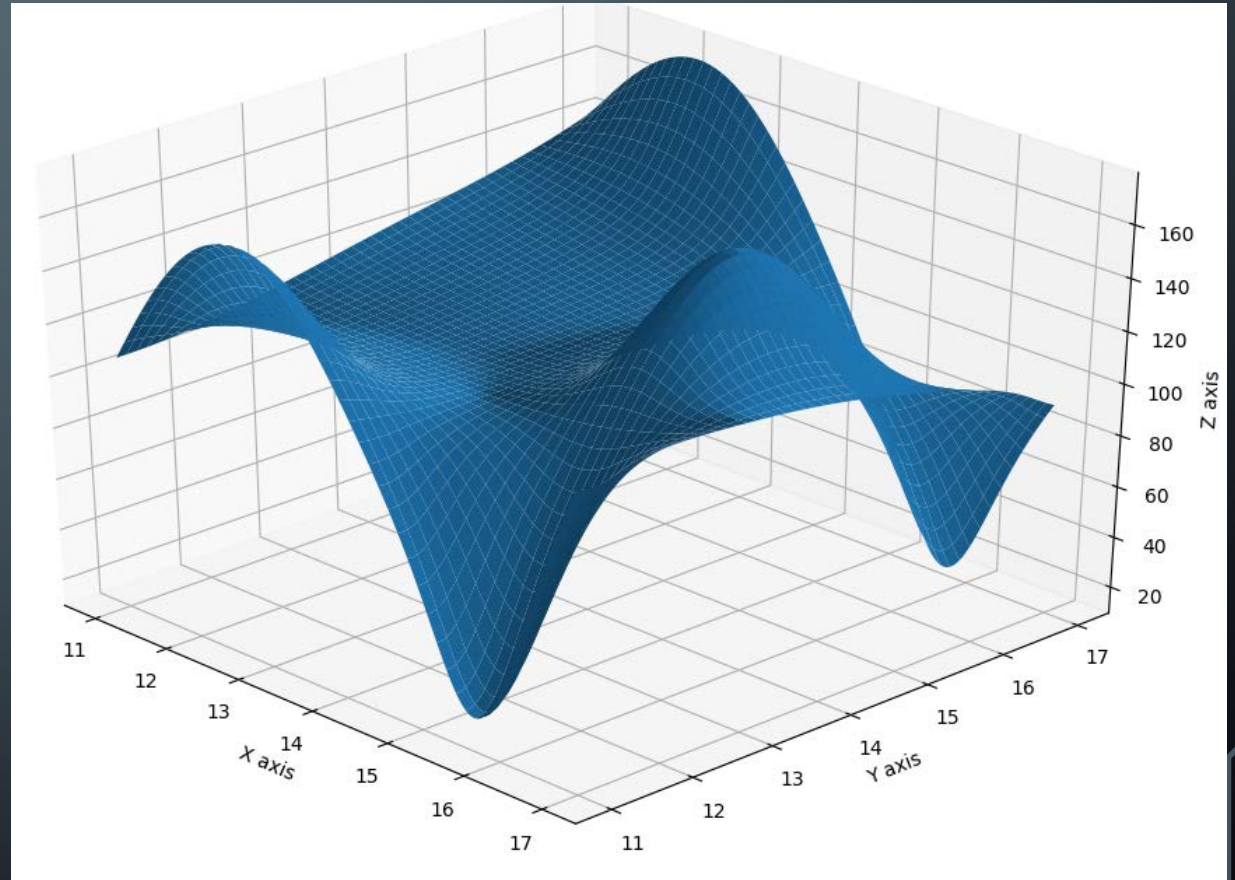
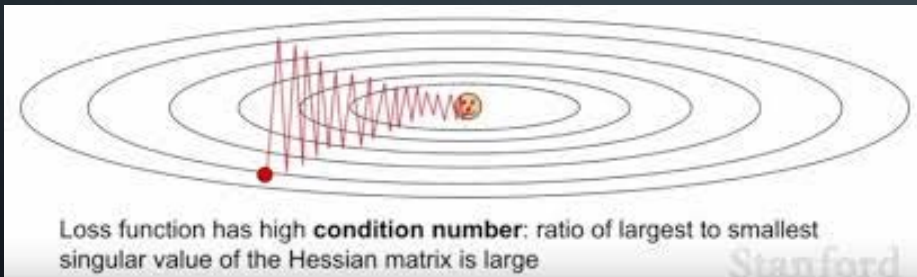


ONLINE / OFFLINE LEARNING

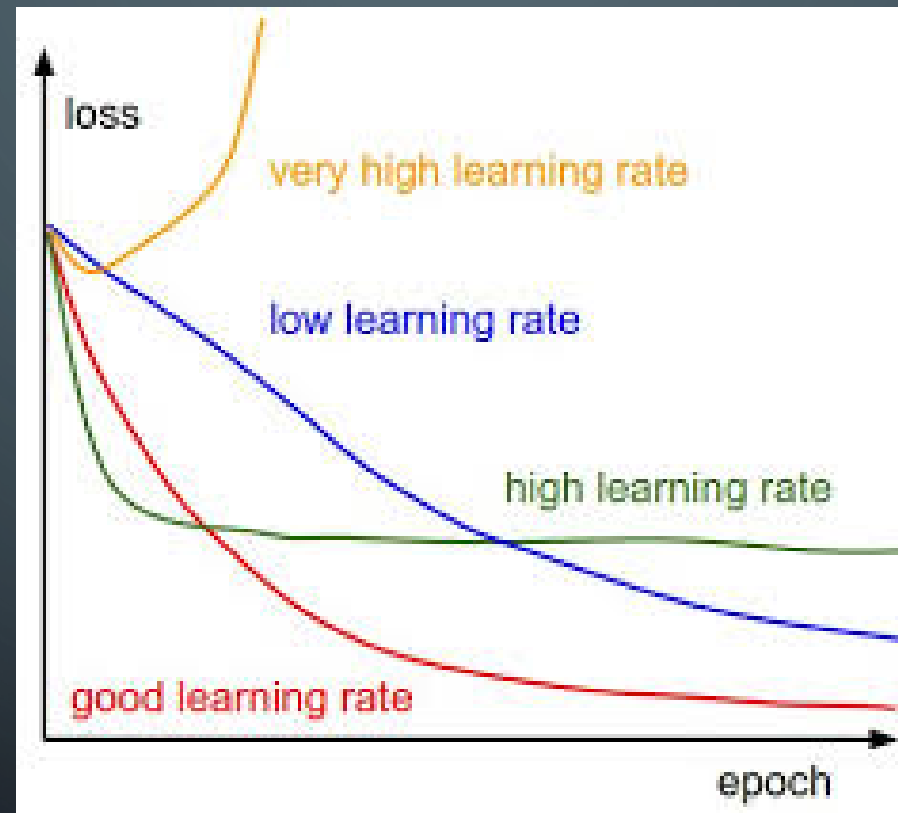
- We don't have all data
 - We need model to update during test time
 - Good for dynamic environments
- 
- 

PROBLEMS IN OPTIMIZERS

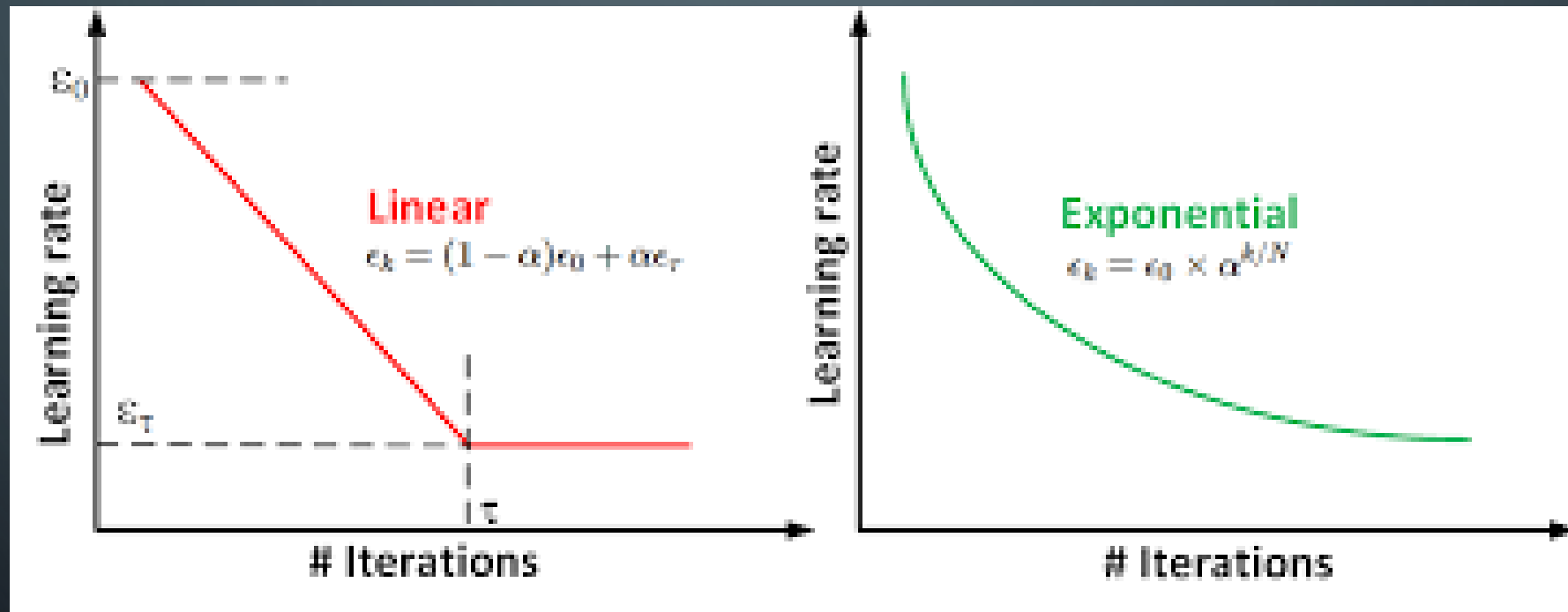
- Plateau
- Saddle point
- Local minimum
- Zig-zag moves



GRADIENT DESCENT – LR



GRADIENT DECENT – DECAY LR

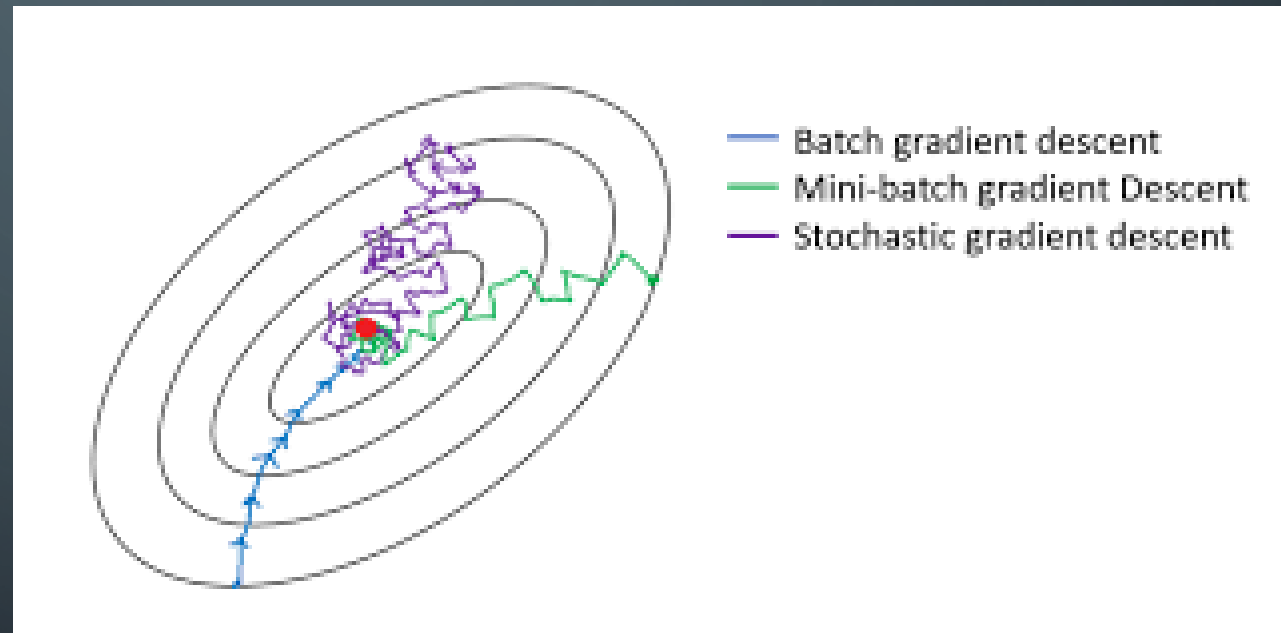


STOCHASTIC GRADIENT DESCENT (SGD)

- Stochastic moves
- avoiding local minimum
- avoiding saddle points
- avoiding plateau
- Computational complexity

MINI-BATCH GRADIENT DESCENT

- Mini-batch
- Not too stochastic
- Fast
- Scalable
- Batch size
- epoch



MOMENTUM

- Like momentum in physics
- Remember update at each step
- Determine next update using
 - Gradient
 - Pervious updates
- avoiding zig-zag moves

MOMENTUM FORMULA

- Normal GD:

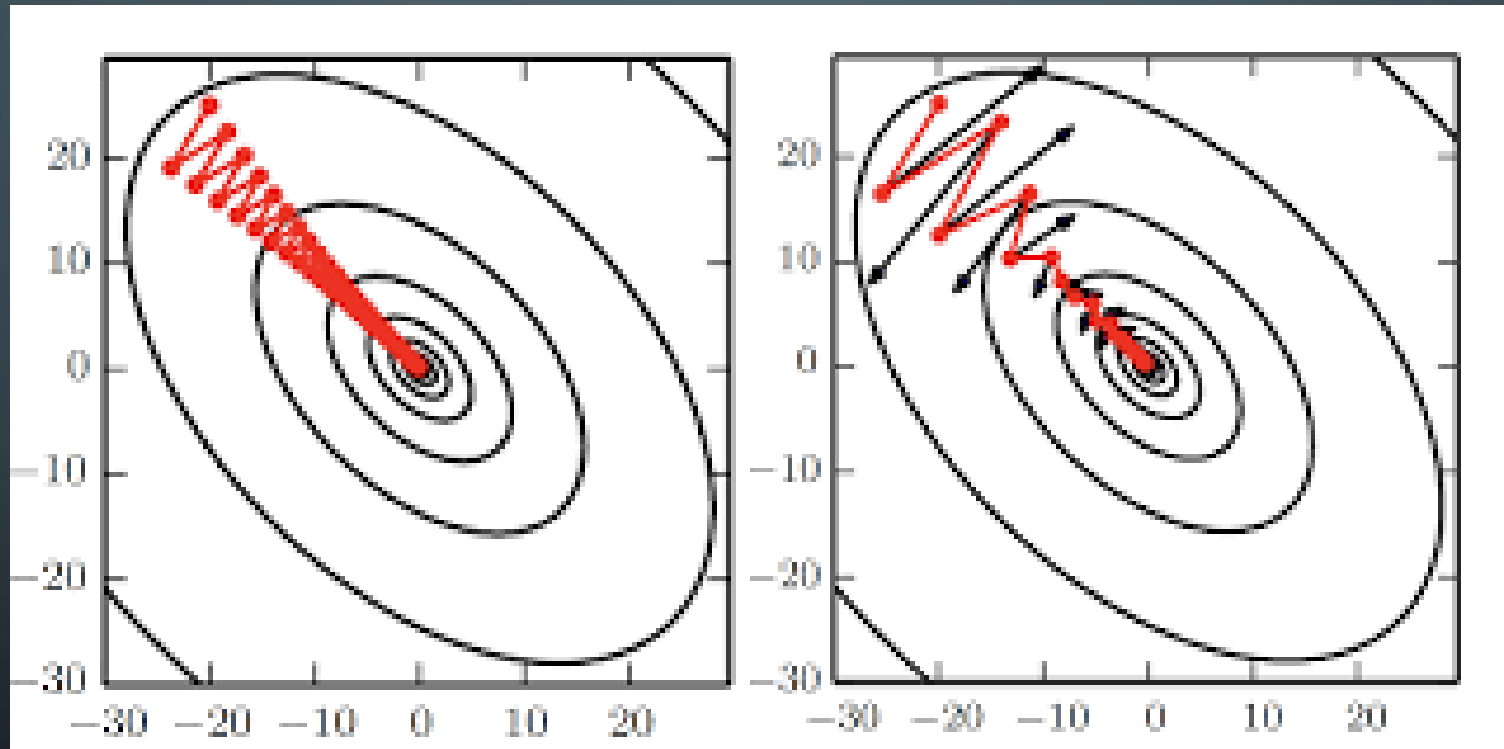
$$\Delta W = -\eta \nabla J$$

- Momentum GD:

$$v = \alpha v - \eta \nabla J$$

$$\Delta W = v$$

MOMENTUM EFFECT



ADAPTIVE LEARNING RATE

- Increase or decreases learning rate during training
- different learning rates for different parameters
- AdaDelta
- AdaGrad
- RMSprop
- Adam

ADAGRAD

- scaling learning rates inversely proportional to the square root of the sum of all the historical squared values of the gradient
- For w_i where $\frac{\partial J}{\partial w_i}$ is small, $\Delta\eta_i$ is small
- For w_j where $\frac{\partial J}{\partial w_j}$ is large, $\Delta\eta_j$ is large
- Not good in some nonconvex functions

$$r = r + (\nabla J \odot \nabla J)$$

$$\Delta W = -\frac{\eta}{\delta + \sqrt{r}} \nabla J$$

RMSPROP

- Better than AdaGrad in nonconvex functions
- Like AdaGrad but with leakage

$$r = \rho r + (1 - \rho)(\nabla J \odot \nabla J)$$

$$\Delta W = -\frac{\eta}{\delta + \sqrt{r}} \nabla J$$

ADAM

- Using both momentum and RMSprop

$$r = \rho r + (1 - \rho)(\nabla J \odot \nabla J)$$

$$v = \alpha v - \frac{\eta}{\delta + \sqrt{r}} \nabla J$$

$$\Delta W = v$$

VISUALIZATION

- <https://emiliendupont.github.io/2018/01/24/optimization-visualization/>

EVOLUTIONARY COMPUTING

- Another optimization methods
- Based on generation
- Natural selection
- Survival of the fittest