

Recommender Systems || Lecture 5

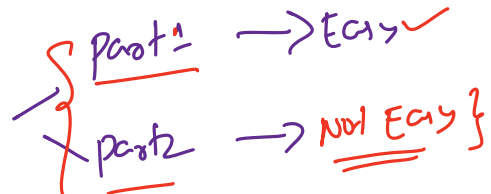
Summer 2022

Dr. Karthik Mohan

Univ. of Washington, Seattle

July 19, 2022

Logistics

- 1 Programming 2 due this Sunday A handwritten diagram in red ink. A blue line from the word 'Sunday' in the text below branches into two paths. The upper path is enclosed in a red curly bracket and points to the text 'part 1' (underlined in blue) followed by an arrow pointing to 'Easy' (with a checkmark). The lower path is also enclosed in a red curly bracket and points to the text 'part 2' (underlined in blue) followed by an arrow pointing to 'Not Easy' (underlined in red).

Logistics

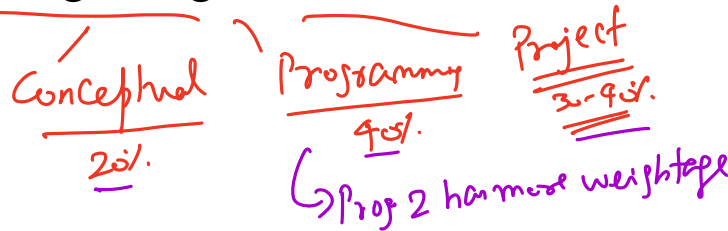
- 1 Programming 2 due this Sunday
- 2 Please spend good time on the programming assignment

Logistics

- 1 Programming 2 due this Sunday
- 2 Please spend good time on the programming assignment
- 3 Programming 2 part 1 focuses on SVD based recommendations.
Programming 2 part 2 is based on SGD (stochastic gradient descent) and matrix factorization.

Logistics

- 1 Programming 2 due this Sunday
- 2 Please spend good time on the programming assignment
- 3 Programming 2 part 1 focuses on SVD based recommendations. Programming 2 part 2 is based on SGD (stochastic gradient descent) and matrix factorization.
- 4 Any questions on grading? →



Logistics

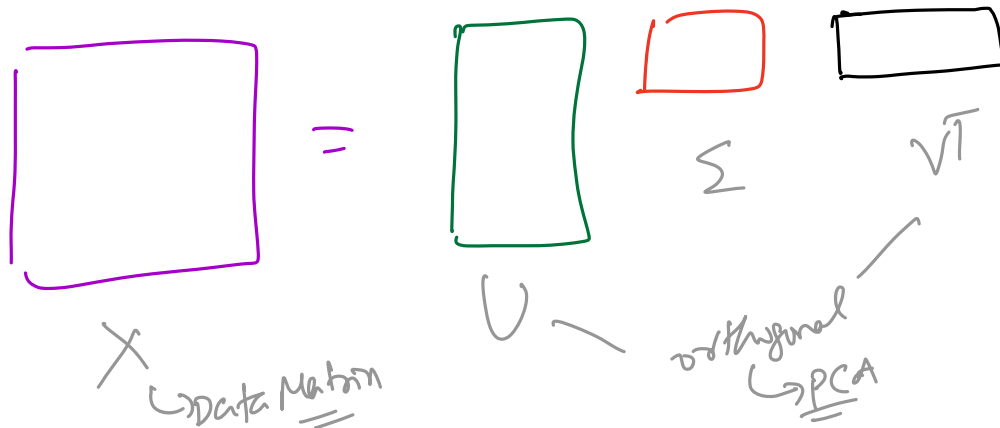
- 1 Programming 2 due this Sunday
- 2 Please spend good time on the programming assignment
- 3 Programming 2 part 1 focuses on SVD based recommendations.
Programming 2 part 2 is based on SGD (stochastic gradient descent) and matrix factorization.
- 4 Any questions on grading?
- 5 Any other questions?

Limitations of SVD

- ① PCA is based on SVD for dimensionality reduction. And we know PCA is not robust to outliers

Limitations of SVD

- 1 PCA is based on SVD for dimensionality reduction. And we know PCA is not robust to outliers
- 2 SVD is a nice baseline however it is very specific - Requires factors U, V to be orthogonal and is not flexible to regularization or other considerations.



Limitations of SVD

- ① PCA is based on SVD for dimensionality reduction. And we know PCA is not robust to outliers
- ② SVD is a nice baseline however it is very specific - Requires factors U , V to be orthogonal and is not flexible to regularization or other considerations.
- ③ Requires behavioral data - Can't handle hybrid formulations that include content as well

Limitations of SVD

- ① PCA is based on SVD for dimensionality reduction. And we know PCA is not robust to outliers
- ② SVD is a nice baseline however it is very specific - Requires factors U, V to be orthogonal and is not flexible to regularization or other considerations.
- ③ Requires behavioral data - Can't handle hybrid formulations that include content as well
- ④ **Solution and Extension?** Matrix Factorization!

→ Another Limitation :- can't work on an Incomplete matrix!

Matrix Factorization Methods

Matrix Factorization Problem

Let $X \in \mathcal{R}^{m \times n}$ be a data matrix, say for ratings of users vs movies. Rows represent users and columns represent movies, and entries represent the ratings. Then, mathematically, the matrix factorization problem is defined as:

more flexible on U, V than SVD!

$$\min_{U, V} \frac{1}{2} \|X - UV\|_F^2$$

CCCA, NMF, LCA, ...

where, $U \in \mathcal{R}^{m \times r}$ and r is the dimension corresponding to low-dimensional factors (U and V).

Matrix Completion

$1M$ customers & $100k$ products \approx 100 Billion possible entries!!

$U \downarrow V \downarrow$

$1M \times 1000$ $1000 \times 100k$ \Rightarrow 0.01% (In reality)

\Rightarrow 0 (Number of parameters learned) = ? $O(1B)$

Matrix Factorization Methods

Matrix Factorization Problem

Let $X \in \mathcal{R}^{m \times n}$ be a data matrix, say for ratings of users vs movies. Rows represent users and columns represent movies, and entries represent the ratings. Then, mathematically, the matrix factorization problem is defined as:

$$\min_{U, V} \frac{1}{2} \|X - UV\|_F^2$$

where, $U \in \mathcal{R}^{m \times r}$ and r is the dimension corresponding to low-dimensional factors (U and V).

Quick question

Based on the set up above, What's the dimension of the matrix V ?

Matrix Factorization vs Matrix Completion

Missing ratings

When we have missing ratings in X , instead of filling it in, we can actually only use the ratings we know to learn the factors U and V . This formulation is called the **matrix completion** problem.

Matrix Factorization vs Matrix Completion

Missing ratings

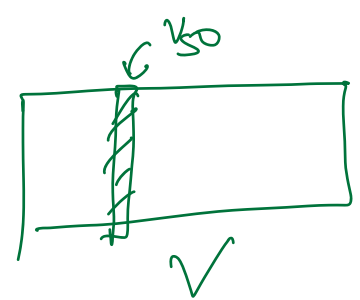
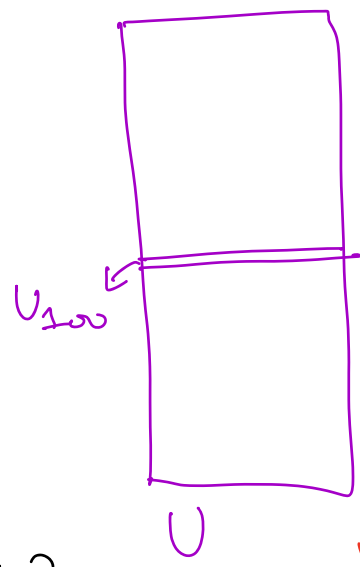
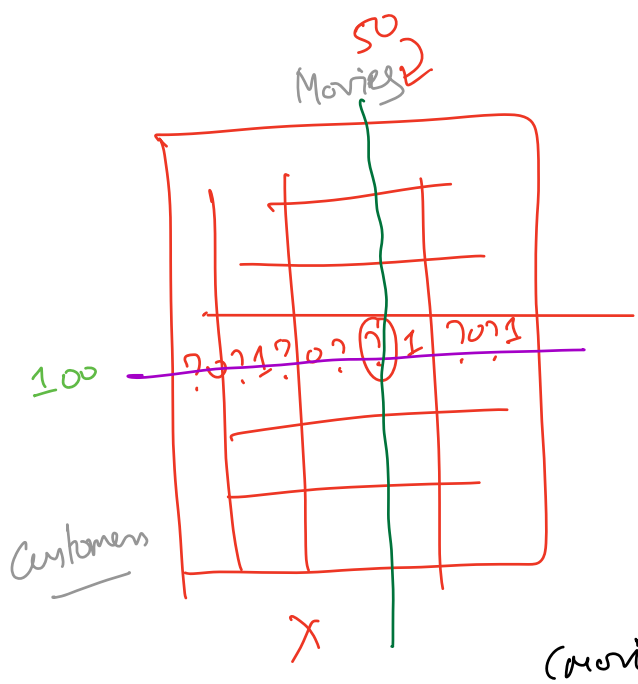
When we have missing ratings in X , instead of filling it in, we can actually only use the ratings we know to learn the factors U and V . This formulation is called the **matrix completion** problem.

Matrix Completion 1 Customer \rightarrow rate / purchase
SO product
(work)
SO $< 0.2\%$
work

Matrix Completion

$$\min_{U, V} \frac{1}{2} \sum_{(i,j) \in R} (X_{ij} - U_{i,\cdot}^T V_{\cdot,j})^2, \quad \text{Optimization}$$

where R is the set of tuples (i, j) for which a rating is known apriori. How does this compare with matrix factorization with imputation?



$U_{100} \cdot V_{50} = ?$
 High \Rightarrow Customer 100 likes Movie 50
 Low \Rightarrow Customer 100 may not like Movie 50!

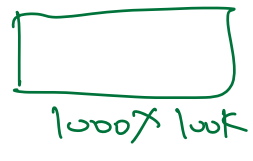
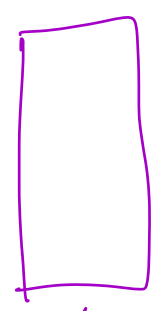
(Movies)
 (Products)

1MM customers, 100k products

\Downarrow
 0.1% of this matrix! } 1 customer purchases 100 products

\Downarrow
 100MM training data points

Test data points $\approx O(100B)$
 unknowns!!



$1MM \times 1000$
 \Downarrow
 # parameters = $O(1B)$

$O(100MM) \rightarrow$ Train data points
 $O(100B) \rightarrow$ unknowns!!
 $O(1B) \rightarrow$ Parameters \ll

underfitting / overfitting

Matrix Factorization Extensions

overfitting
↓
use

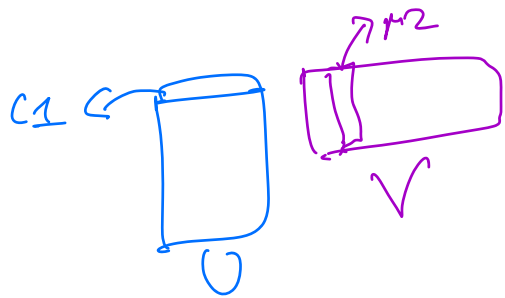
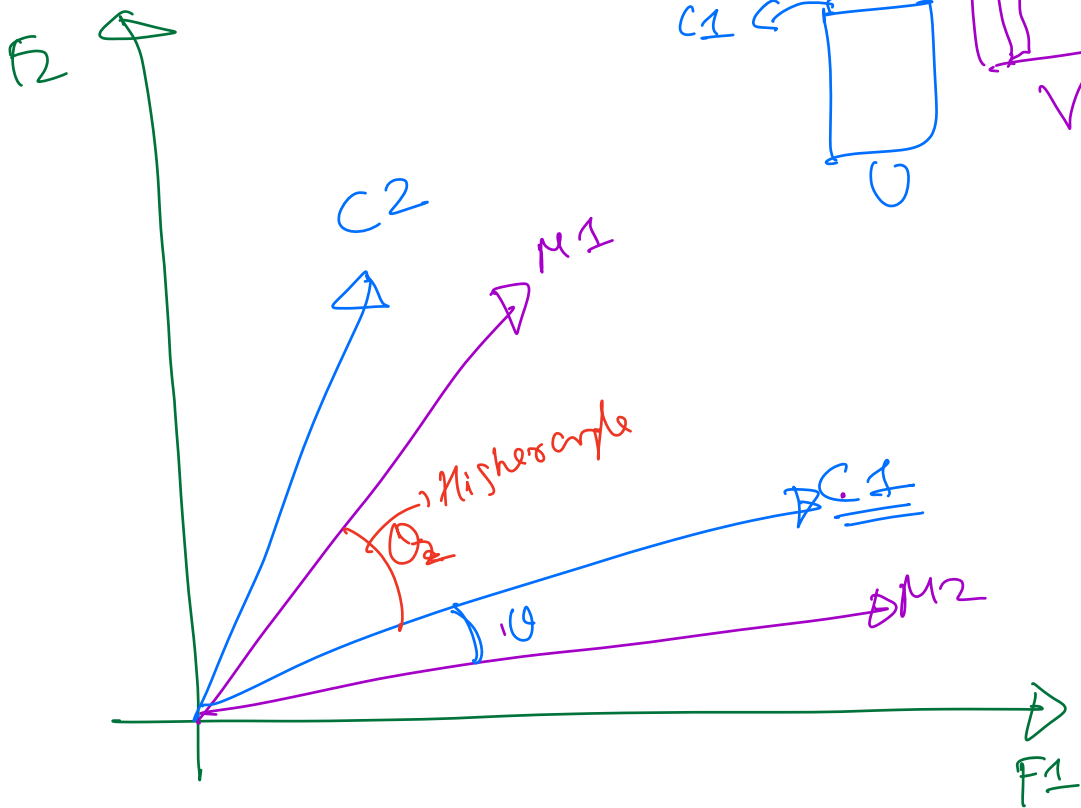
Flexible Formulation - Can add regularization!

$$\min_{U, V} \frac{1}{2} \|X - UV\|_F^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2$$

Minimizing overfitting

λ - Hyper-parameters

$\lambda = 0 \Rightarrow$ Back to matrix completion! $\lambda \uparrow \Rightarrow$ Reduce overfitting further
 $\lambda \rightarrow \infty \Rightarrow U, V \Rightarrow 0$



Matrix Factorization Extensions

Flexible Formulation - Can add regularization!

$$\min_{U, V} \frac{1}{2} \|X - UV\|_F^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2$$

Reduced rank matrix factorization

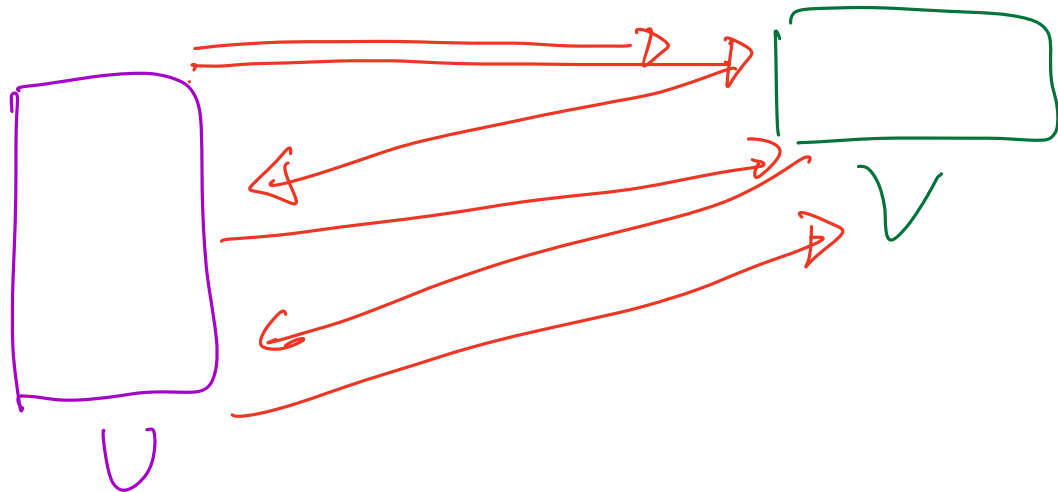
Above formulation with the regularization is known to reduce the dimensionality of the factors obtained.

Alternating Minimization Method for Matrix Completion

Skip to next lecture

$$\min_U \min_V \frac{1}{2} \|X - UV\|_F^2 + \dots$$

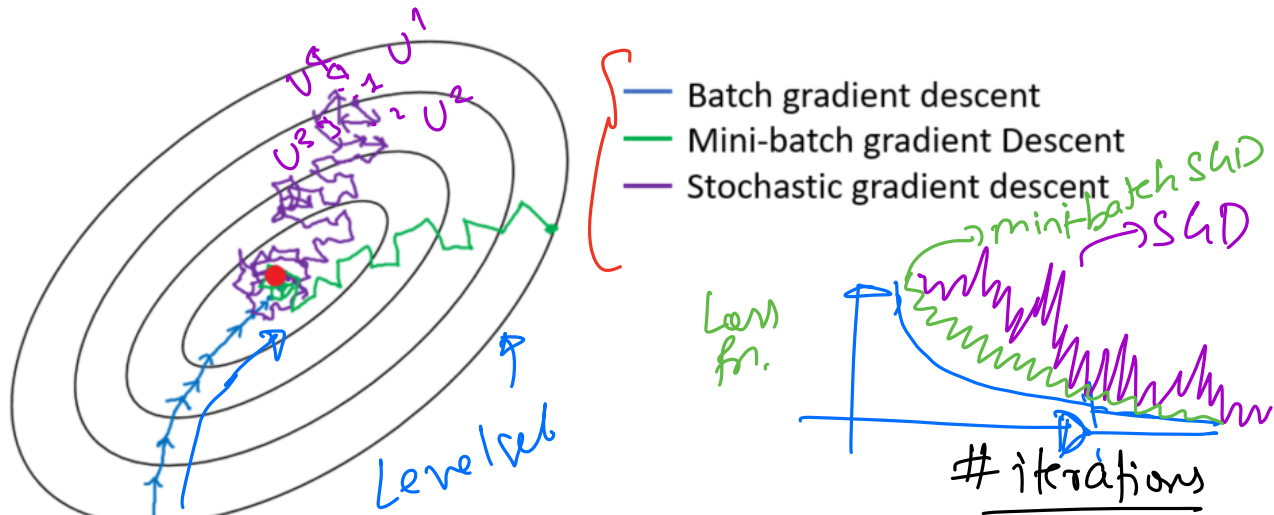
Two sets of parameters:- U, V



Algorithmic foundations to Machine Learning

Underlying Engine behind ML Training

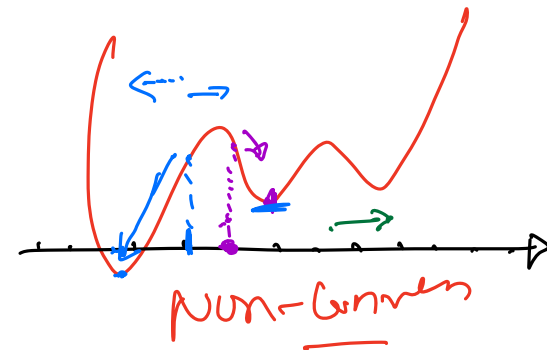
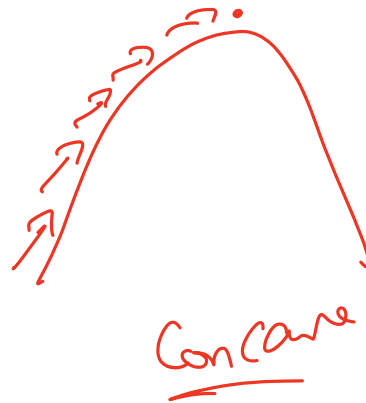
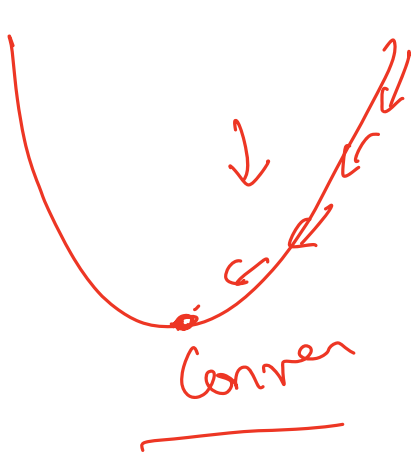
(Mini-batch) Stochastic Gradient Descent Almost every model and problem-space in ML uses SGD of some kind - Clustering, Regression, Deep Learning, Computer Vision and NLP to name a few. Almost every algorithm in every library - Scikit-learn, Keras, Pytorch, etc uses **mini-batch SGD** under the hood.



So what is Gradient Descent?

Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .



So what is Gradient Descent?

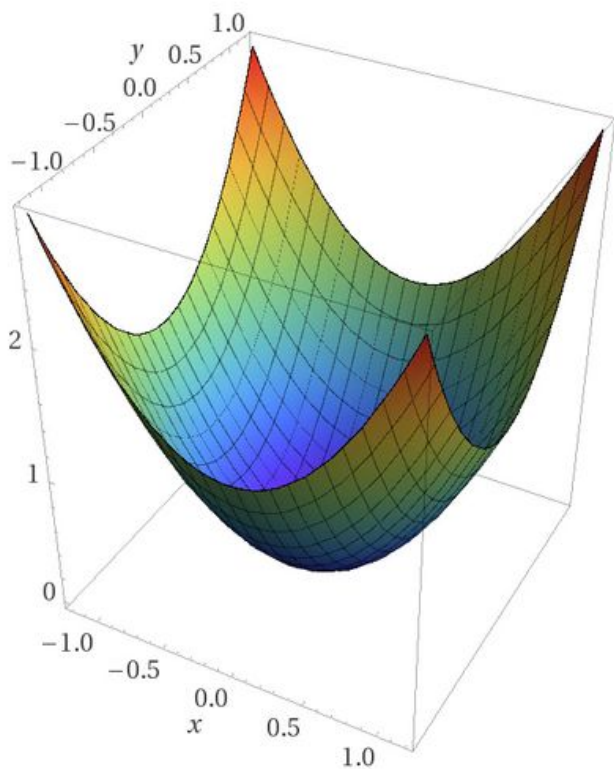
Fundamentally

Take a convex/non-convex function, f . GD allows you to find a local optimum to f .

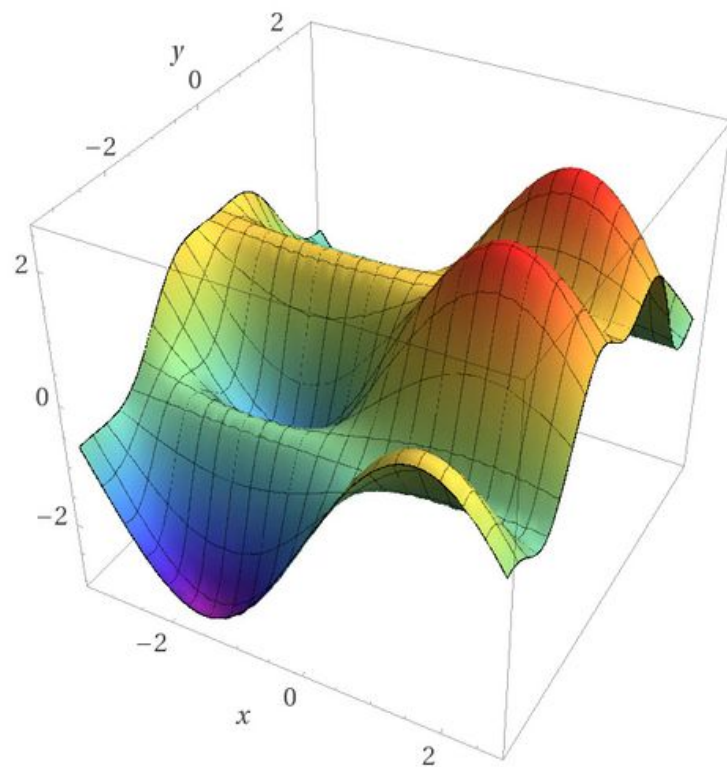
Why is this important?

Consider the Linear Regression problem. \hat{w} is a local optimum to the function $f(w) = \frac{1}{2} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$

Negative Gradient helps you view the direction of descent



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

$$\underline{u}^3 \leftarrow \underline{u}^2 + \alpha (-\nabla l(\underline{u}^2, \underline{v}^2))$$

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(\hat{w})$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)

Gradient Descent

Batch Gradient Descent

Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

① **Initialize** $w = w_0$ (maybe randomize)

② **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$

\leftarrow step size = Learning Rate

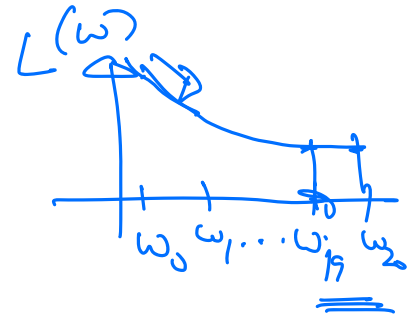
Gradient Descent

Batch Gradient Descent

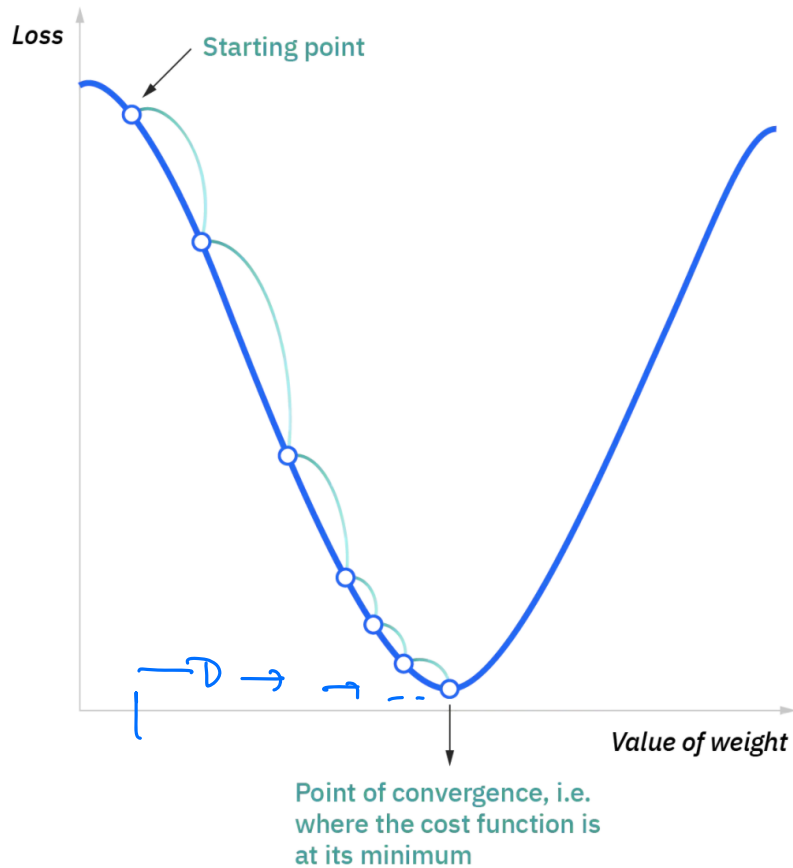
Let us say we want to minimize $L(w)$ - Loss Function and find the best \hat{w} that does that.

- 1 **Initialize** $w = w_0$ (maybe randomize)
- 2 **Gradient Descent** $w \leftarrow w - lr * \nabla L(w)$
- 3 **Iterate** Repeat step 2 until w converges, i.e.

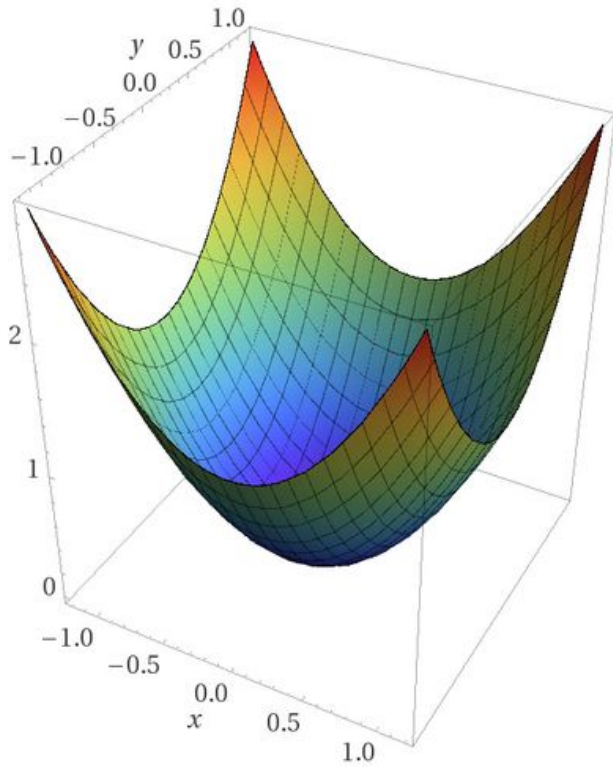
$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$



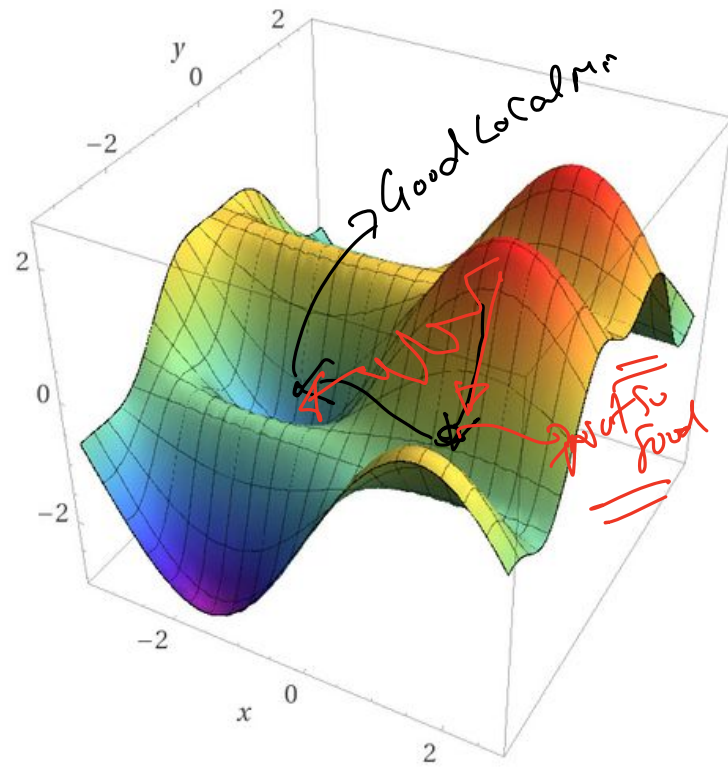
GD in one dimension



Loss function in 2 dimensions



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

Deep Learning and SGD!

Gradient in DL ← "Back propagation"
↓
Auto-Differentiation

Computing Gradients!

$$l(u) = u^T d \quad \begin{matrix} \text{constant vector} \\ \left[\begin{array}{l} d \in \mathbb{R}^{10} \\ u \in \mathbb{R}^{10} \end{array} \right. \end{matrix}$$

$$\frac{\nabla l}{u} = d$$

$$l(u) = \underbrace{u_1 d_1 + u_2 d_2 + \dots + u_{10} d_{10}}$$

$$\frac{\nabla l}{u_1} = \frac{d}{du_1} (u_1 d_1) = d_1 \quad \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_{10} \end{bmatrix} = d$$

$$\frac{\nabla l}{u_k} = d_k$$

ICE #1

$$\|w\|_2^2 = w_1^2 + w_2^2 + \dots + w_d^2$$

$w \in \mathbb{R}^d$

Gradient of Ridge Regularizer (2 mins)

Find the gradient of the regularization function, $R(w) = \lambda \|w\|_2^2$. I.e. obtain the expression for, $\nabla_w R(w)$?

- a) $2\lambda \|w\|_2$
- b) $\lambda \|w\|_2 w$
- c) $2\lambda w$
- d) $2\lambda \|w\|_2 w$

$$\begin{bmatrix} \frac{\partial R}{\partial w_1} \\ \frac{\partial R}{\partial w_2} \\ \vdots \\ \frac{\partial R}{\partial w_d} \end{bmatrix} = \nabla_w R$$

[Poll Link](#)

$$\frac{\partial R}{\partial w_1} = \frac{\partial \lambda(w_1^2 + w_2^2 + \dots + w_d^2)}{\partial w_1}$$

$$\Rightarrow \lambda(2w_1) = 2\lambda w_1$$

$$\frac{\partial R}{\partial w_2} = 2\lambda w_2$$

$$2\lambda \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = 2\lambda \underbrace{\omega}_{\text{vector}}$$

Entry/Scalars

Gradient Descent Properties

- ① Gradient Descent converges to a local minimum

Gradient Descent Properties

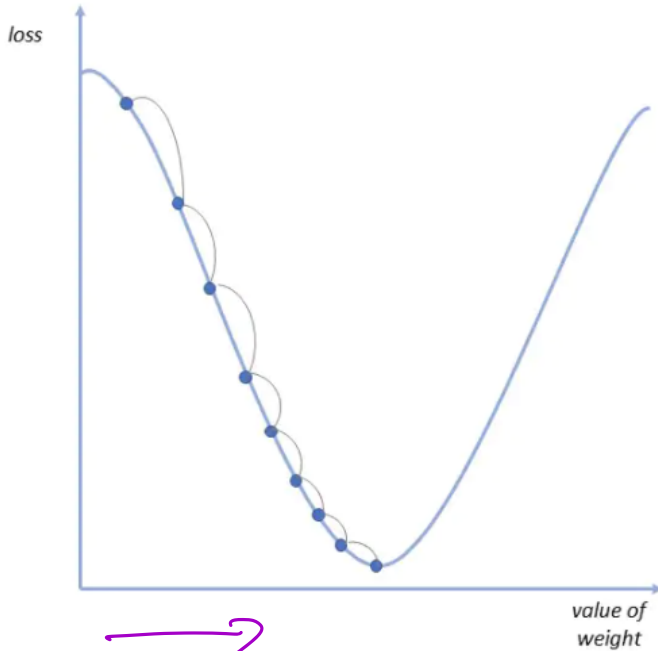
- ① Gradient Descent converges to a local minimum
- ② If L is a convex function, all local minima become a global minima!

Gradient Descent Properties

- ① Gradient Descent converges to a local minimum
- ② If L is a convex function, all local minima become a global minima!
- ③ Wherever we start, gradient descent usually finds a local minima closest to the start.

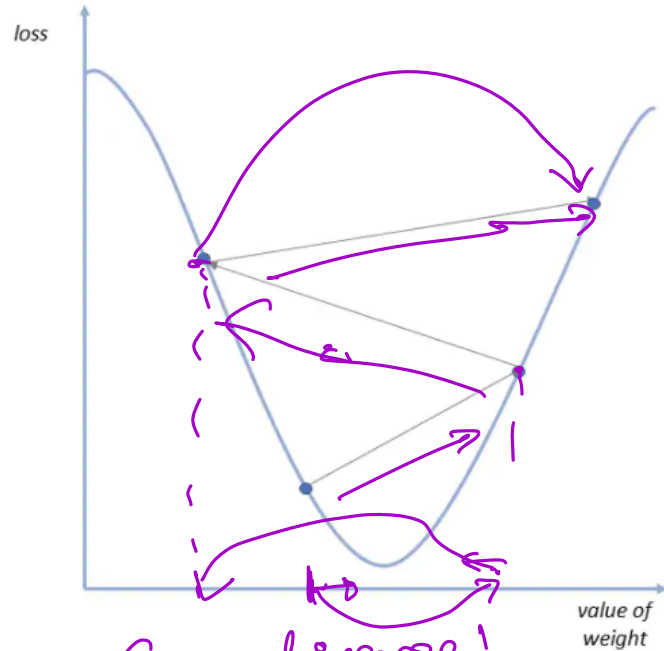
Effect of Learning Rate

Small Learning Rate



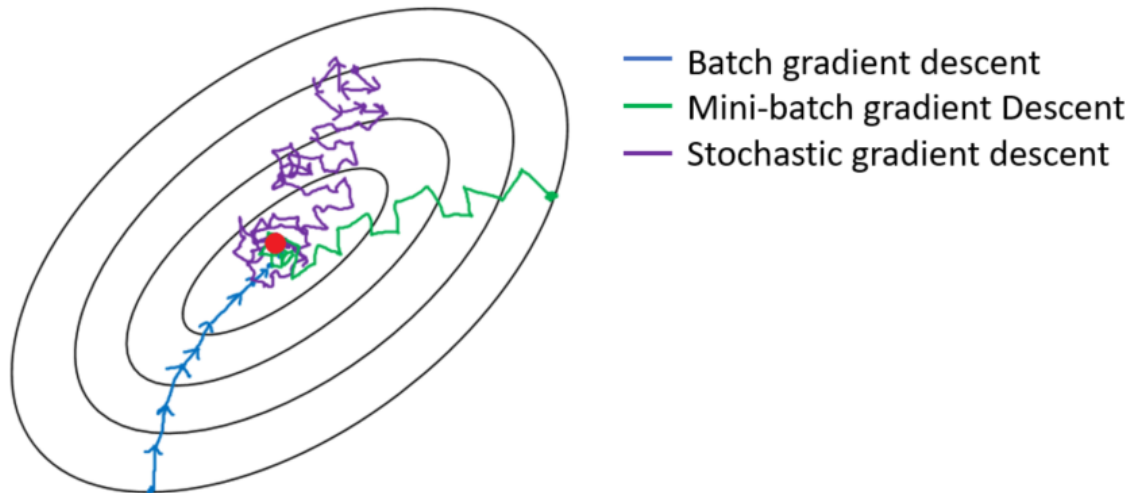
→
More time
to Converge

Large Learning Rate



Can diverge!

GD behavior in the search space



Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 Initialize w^0 (randomize)

$(u^0, v^0) \rightarrow$ matrix completion

data pts.

Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !

Gradient descent in practice - SGD!

SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !
- 2 **Gradient Descent** $\underline{w^{k+1}} \leftarrow w - \underline{lr} * \underline{\nabla L_i(w^k)}$

Gradient descent in practice - SGD!

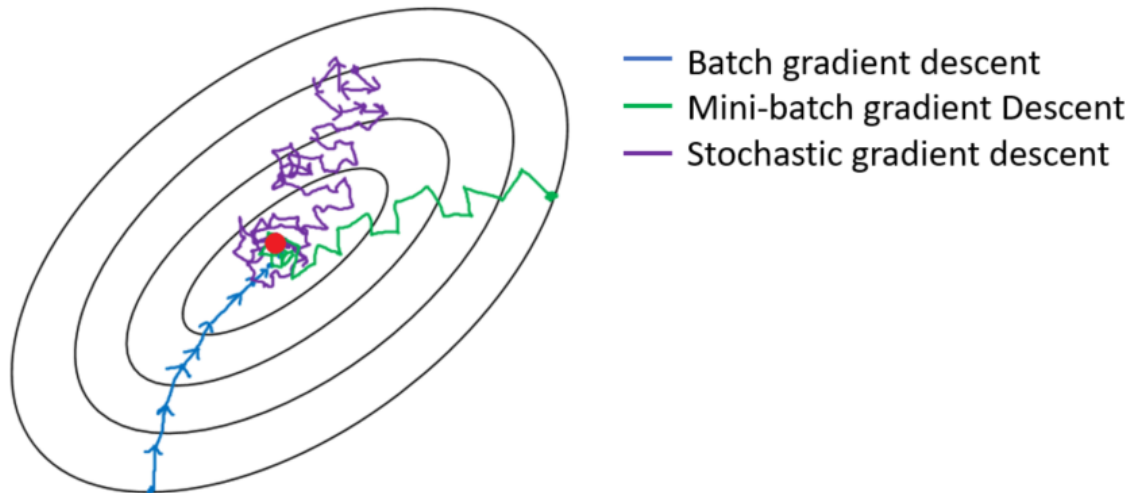
SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w .

- 1 **Initialize** w^0 (randomize) Pick index i at random between 1 and N !
- 2 **Gradient Descent** $w^{k+1} \leftarrow w - lr * \nabla L_i(w^k)$ }
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

SGD behavior in search space



SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

$\equiv N/k$

256

- 1 **Initialize** $w = w_0$ (randomize)

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : i_1, i_2, \dots, i_k ↗

SGD in practice - mini-batch SGD!

mini-batch SGD

Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

① **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : i_1, i_2, \dots, i_k !

② **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$

*pool info from k data pts.
(SGD: $k=1$)*

SGD in practice - mini-batch SGD!

mini-batch SGD

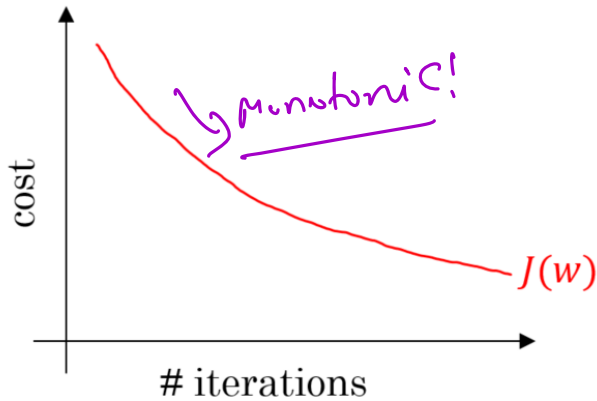
Let $L(w) = \sum_{i=1}^N L_i(w)$ where L_i is a function of only the i th data point (x_i, y_i) and parameter w . Let B be the number of batches and k be the batch size.

- 1 **Initialize** $w = w_0$ (randomize) Pick a batch of k data points at random between 1 and N : i_1, i_2, \dots, i_k !
- 2 **Gradient Descent** $w^{k+1} \leftarrow w^k - lr * \sum_{j=1}^k \nabla_w L_{i_j}(w^k)$
- 3 **Iterate** Repeat step 2 and 3 until w converges, i.e.

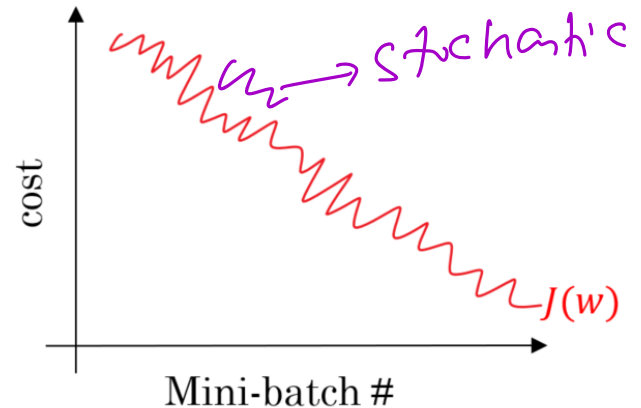
$$\|w^{k+1} - w^k\| / \|w^k\| \leq 10^{-3}$$

GD vs Mini-batch convergence behavior

Batch gradient descent



Mini-batch gradient descent



GD vs mini-batch SGD

Factor	GD	Mini-batch SGD
<u>Data</u>	<u>All per iteration</u>	<u>Mini-batch</u> (usually <u>128</u> or <u>256</u>)
Randomness	Deterministic	Stochastic
Error reduction	Monotonic	Stochastic
Computation	High	Low
Memory big data	<u>Intractable</u>	Tractable
Convergence	Low relative error	Few "passes" on data
Local Minima traps	<u>Yes</u>	<u>No</u>

Stopping Criteria

DL

Programming Assignment 2 Part 2

- 1 Implement SGD yourself, from scratch for recommender systems!

Programming Assignment 2 Part 2

- 1 Implement SGD yourself, from scratch for recommender systems!
- 2 Apply it with/without regularizers!

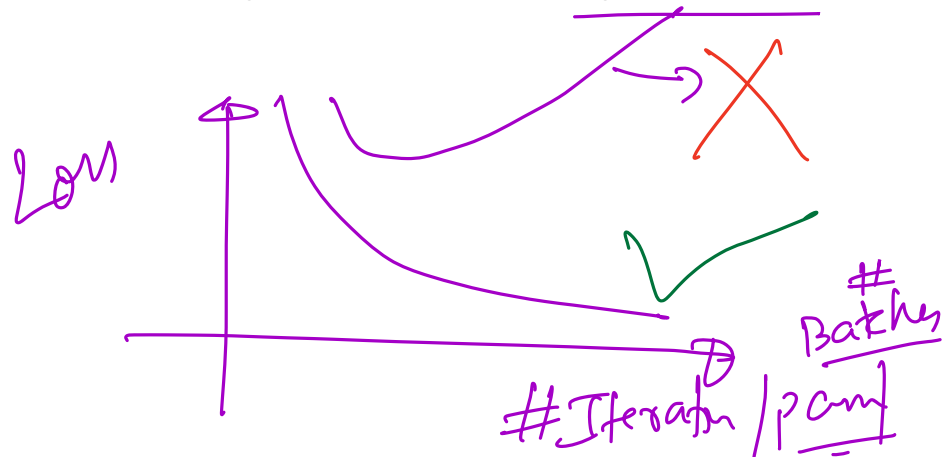
Programming Assignment 2 Part 2

MovieLens Data set look rating

- 1 Implement SGD yourself, from scratch for recommender systems!
- 2 Apply it with/without regularizers!
- 3 Use “auto-differentiation” part of MXNet/Pytorch DL frameworks to compute gradients automatically!

Programming Assignment 2 Part 2

- 1 Implement SGD yourself, from scratch for recommender systems!
- 2 Apply it with/without regularizers!
- 3 Use “auto-differentiation” part of MXNet/Pytorch DL frameworks to compute gradients automatically!
- 4 Test that you have implemented it right as the training error will make it obvious to you!



Programming Assignment 2 Part 2

- 1 Implement SGD yourself, from scratch for recommender systems!
- 2 Apply it with/without regularizers!
- 3 Use “auto-differentiation” part of MXNet/Pytorch DL frameworks to compute gradients automatically!
- 4 Test that you have implemented it right as the training error will make it obvious to you!
- 5 In practice, SGD isn't implemented - But this is the fundamental building block behind any ML/DL algorithm/framework such as Scikit learn, Pytorch, Keras, etc.