COMS 4771 Lecture 24

- 1. Collaborative filtering
- 2. Stochastic gradient descent

Collaborative filtering

Netflix problem: 480000 users, 18000 movies.

Netflix problem: 480000 users, 18000 movies.

▶ Each user rates some subset of the movies with a score in $\{1, 2, \dots, 5\}$.

On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated $>\!17000$ movies).

Netflix problem: 480000 users, 18000 movies.

- ► Each user rates some subset of the movies with a score in {1,2,...,5}.
 On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- ▶ Goal is to predict how users would rate movies they haven't seen.

Netflix problem: 480000 users, 18000 movies.

- ► Each user rates some subset of the movies with a score in {1,2,...,5}.
 On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- ▶ Goal is to predict how users would rate movies they haven't seen.
- ▶ Common to reduce $\{1, 2, ..., 5\}$ to $\{-1, +1\}$ (e.g., $\{4, 5\} \mapsto +1$).

Netflix problem: 480000 users, 18000 movies.

- ▶ Each user rates some subset of the movies with a score in $\{1, 2, ..., 5\}$. On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- Goal is to predict how users would rate movies they haven't seen.
- ▶ Common to reduce $\{1,2,\ldots,5\}$ to $\{-1,+1\}$ (e.g., $\{4,5\}\mapsto +1$).

Common supervised learning approach:

▶ Get *features* for each user i and movie j (e.g., $x_i \in \mathbb{R}^d$ and $y_j \in \mathbb{R}^d$); goal is to predict rating as function of (x_i, y_j) .

Netflix problem: 480000 users, 18000 movies.

- ▶ Each user rates some subset of the movies with a score in $\{1, 2, ..., 5\}$. On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- ▶ Goal is to predict how users would rate movies they haven't seen.
- ▶ Common to reduce $\{1,2,\ldots,5\}$ to $\{-1,+1\}$ (e.g., $\{4,5\}\mapsto +1$).

Common supervised learning approach:

- ▶ Get *features* for each user i and movie j (e.g., $x_i \in \mathbb{R}^d$ and $y_j \in \mathbb{R}^d$); goal is to predict rating as function of (x_i, y_j) .
- ▶ Linear function: $x_i y_j^\top \mapsto \langle W, x_i y_j^\top \rangle = x_i^\top W y_j$ for $W \in \mathbb{R}^{d \times d}$. Can use SVM, logistic regression, boosting, etc.

Netflix problem: 480000 users, 18000 movies.

- ▶ Each user rates some subset of the movies with a score in $\{1, 2, \ldots, 5\}$. On average, each user rates around 200 movies, though the variance is high (e.g., some user has rated >17000 movies).
- ▶ Goal is to predict how users would rate movies they haven't seen.
- ▶ Common to reduce $\{1,2,\ldots,5\}$ to $\{-1,+1\}$ (e.g., $\{4,5\}\mapsto +1$).

Common supervised learning approach:

- ▶ Get *features* for each user i and movie j (e.g., $x_i \in \mathbb{R}^d$ and $y_j \in \mathbb{R}^d$); goal is to predict rating as function of (x_i, y_j) .
- ▶ Linear function: $x_i y_j^\top \mapsto \langle W, x_i y_j^\top \rangle = x_i^\top W y_j$ for $W \in \mathbb{R}^{d \times d}$. Can use SVM, logistic regression, boosting, etc.

What if you don't have any features?

COLLABORATIVE FILTERING

Collaborative filtering:

- Users who rate the same movie similarly are likely to be similar.
- ▶ Movies that are rated similarly by the same user are likely to be similar.

COLLABORATIVE FILTERING

Collaborative filtering:

- ▶ Users who rate the same movie similarly are likely to be similar.
- ▶ Movies that are rated similarly by the same user are likely to be similar.

Upshot: there may be a wealth of information just in the ratings themselves . . .

Collaborative filtering

what a genius idea??? get the information from ratings themselves

Collaborative filtering:

- ▶ Users who rate the same movie similarly are likely to be similar.
- ▶ Movies that are rated similarly by the same user are likely to be similar.

Upshot: there may be a wealth of information just in the ratings themselves . . .

How can we effectively formalize the intuition stated above?

Simple latent factor model:

Simple latent factor model:

lacktriangle Each user i has some hidden feature vector $oldsymbol{u}_i \in \mathbb{R}^k$.

Simple latent factor model:

- ▶ Each user i has some *hidden* feature vector $u_i \in \mathbb{R}^k$.
- ▶ Each movie j has some *hidden* feature vector $v_j \in \mathbb{R}^k$.

Simple latent factor model:

- ▶ Each user i has some *hidden* feature vector $u_i \in \mathbb{R}^k$.
- **Each** movie j has some *hidden* feature vector $v_j \in \mathbb{R}^k$.
- ▶ Rating $a_{i,j}$ assigned by user i to movie j is, in expectation, $\langle u_i, v_j \rangle$.

$$\mathbb{E} \left\{ \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \right\} = \begin{bmatrix} - & \boldsymbol{u}_1^\top & - \\ - & \boldsymbol{u}_2^\top & - \\ \vdots & \vdots & - \\ - & \boldsymbol{u}_m^\top & - \end{bmatrix} \begin{bmatrix} \boldsymbol{\downarrow} & \boldsymbol{\downarrow} & \boldsymbol{\downarrow} \\ \boldsymbol{v}_1 & \boldsymbol{v}_2 & \cdots & \boldsymbol{v}_n \\ \boldsymbol{\downarrow} & \boldsymbol{\downarrow} & \boldsymbol{\downarrow} \end{bmatrix}$$

Simple latent factor model:

- ▶ Each user i has some *hidden* feature vector $u_i \in \mathbb{R}^k$.
- lacktriangle Each movie j has some *hidden* feature vector $oldsymbol{v}_j \in \mathbb{R}^k$.
- ▶ Rating $a_{i,j}$ assigned by user i to movie j is, in expectation, $\langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle$.

$$\mathbb{E} \left\{ \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \right\} = \begin{bmatrix} - & \boldsymbol{u}_1^\top & - \\ - & \boldsymbol{u}_2^\top & - \\ \vdots & \vdots & - \\ - & \boldsymbol{u}_m^\top & - \end{bmatrix} \begin{bmatrix} | & | & | & | \\ \boldsymbol{v}_1 & \boldsymbol{v}_2 & \cdots & \boldsymbol{v}_n \\ | & | & | & | \end{bmatrix}$$
 great way to understand

We'll assume $k \leq \min\{m,n\}$ for reasons that will become clear later matrix!!!

Complete ratings

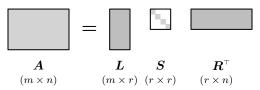
$$\mathbb{E} \left\{ \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \right\} = \begin{bmatrix} - & \boldsymbol{u}_1^\top & - \\ - & \boldsymbol{u}_2^\top & - \\ \vdots & \vdots & - \\ - & \boldsymbol{u}_m^\top & - \end{bmatrix} \begin{bmatrix} \boldsymbol{\downarrow} & \boldsymbol{\downarrow} & \boldsymbol{\downarrow} \\ \boldsymbol{v}_1 & \boldsymbol{v}_2 & \cdots & \boldsymbol{v}_n \\ \boldsymbol{\downarrow} & \boldsymbol{\downarrow} & \boldsymbol{\downarrow} \end{bmatrix}$$

Suppose every user rates every movie.

COMPLETE RATINGS

$$\mathbb{E} \left\{ \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{bmatrix} \right\} = \begin{bmatrix} \begin{matrix} \begin{matrix} & \boldsymbol{u}_1^\top & \begin{matrix} \\ \end{matrix} \\ \begin{matrix} \boldsymbol{u}_2^\top \end{matrix} \\ \begin{matrix} \vdots \\ \end{matrix} \\ \begin{matrix} \boldsymbol{u}_n^\top \end{matrix} \\ \begin{matrix} \begin{matrix} \boldsymbol{v}_1 & \boldsymbol{v}_2 & \cdots & \boldsymbol{v}_n \\ \end{matrix} \\ \begin{matrix} \begin{matrix} & \begin{matrix} \end{matrix} \\ \end{matrix} \\ \begin{matrix} \end{matrix} \end{matrix} \\ \begin{matrix} \begin{matrix} \end{matrix} \end{matrix} \end{bmatrix}$$

Suppose every user rates every movie. Can use SVD:



Use rows of $m{L}m{S}^{1/2}$ as the $m{u}_i$, and rows of $m{R}m{S}^{1/2}$ as the $m{v}_j$.

Complete ratings

Suppose every user rates every movie. Can use SVD:

Use rows of $m{L}m{S}^{1/2}$ as the $m{u}_i$, and rows of $m{R}m{S}^{1/2}$ as the $m{v}_j$.

Unclear what use this is (e.g., if a new movie comes along, can we predict how users will rate it?).

Known user features

Suppose the users' features were not hidden.

Known user features

Suppose the users' features were not hidden.

Could try to *estimate* each movie vector v_j using linear regression:

Known user features

Suppose the users' features were not hidden.

Could try to *estimate* each movie vector v_j using linear regression:

The catch: we don't observe all of the entries of A.

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed. Regard these $a_{i,j}$ as the training data.

Let $\Omega \subseteq m] \times [n]$ be the entries for which $a_{i,j}$ is observed. Regard these $a_{i,j}$ as the training data.

Training obj.: find $U = [u_1| \dots | u_m] \in \mathbb{R}^{k \times m}$ and $V = [v_1| \dots | v_n] \in \mathbb{R}^{k \times n}$ with as small squared error as possible:

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j)\in\Omega} (a_{i,j} - \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle)^2.$$

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed. Regard these $a_{i,j}$ as the training data.

Training obj.: find $U = [u_1| \dots | u_m] \in \mathbb{R}^{k \times m}$ and $V = [v_1| \dots | v_n] \in \mathbb{R}^{k \times n}$ with as small squared error as possible:

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j)\in\Omega} (a_{i,j} - \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle)^2.$$

Often called (low rank) matrix completion, because this is equivalent to:

$$\label{eq:min_problem} \begin{split} \min_{\boldsymbol{X} \in \mathbb{R}^{m \times n}} & & \sum_{(i,j) \in \Omega} (a_{i,j} - X_{i,j})^2 \\ \text{s.t.} & & \operatorname{rank}(\boldsymbol{X}) \leq k. \end{split}$$

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed. Regard these $a_{i,j}$ as the training data.

Training obj.: find $U = [u_1| \dots | u_m] \in \mathbb{R}^{k \times m}$ and $V = [v_1| \dots | v_n] \in \mathbb{R}^{k \times n}$ with as small squared error as possible:

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j)\in\Omega} (a_{i,j} - \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle)^2.$$

Often called (low rank) matrix completion, because this is equivalent to:

$$\min_{\boldsymbol{X} \in \mathbb{R}^{m \times n}} \qquad \sum_{(i,j) \in \Omega} (a_{i,j} - X_{i,j})^2$$
s.t.
$$\operatorname{rank}(\boldsymbol{X}) \le k.$$

Unfortunately, objective is not convex due to rank constraint.

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed. Regard these $a_{i,j}$ as the training data.

Training obj.: find $U = [u_1| \dots | u_m] \in \mathbb{R}^{k \times m}$ and $V = [v_1| \dots | v_n] \in \mathbb{R}^{k \times n}$ with as small squared error as possible:

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j)\in\Omega} (a_{i,j} - \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle)^2.$$

Often called (low rank) matrix completion, because this is equivalent to:

$$\min_{\boldsymbol{X} \in \mathbb{R}^{m \times n}} \qquad \sum_{(i,j) \in \Omega} (a_{i,j} - X_{i,j})^2$$
s.t.
$$\operatorname{rank}(\boldsymbol{X}) \le k.$$

Unfortunately, objective is not convex due to rank constraint.

Nevertheless, can derive an alternating minimization algorithm (similar to k-means and dictionary learning) to iteratively update U and V.

ALTERNATING MINIMIZATION

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

- ▶ Somehow initialize $u_i \in \mathbb{R}^k$ for each $i \in [m]$ and $v_j \in \mathbb{R}^k$ for each $j \in [n]$.
- ► Repeat until convergence:
 - ▶ For each user $i \in [m]$,

$$oldsymbol{u}_i \coloneqq rg\min_{oldsymbol{u}_i \in \mathbb{R}^k} f(oldsymbol{U}, oldsymbol{V})$$

▶ For each movie $j \in [n]$,

$$oldsymbol{v}_j \coloneqq rg\min_{oldsymbol{v}_j \in \mathbb{R}^k} f(oldsymbol{U}, oldsymbol{V})$$

ALTERNATING MINIMIZATION

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

- ▶ Somehow initialize $u_i \in \mathbb{R}^k$ for each $i \in [m]$ and $v_j \in \mathbb{R}^k$ for each $j \in [n]$.
- ► Repeat until convergence:
 - ▶ For each user $i \in [m]$,

$$oldsymbol{u}_i \coloneqq \left(\sum_{j \in [n]: (i,j) \in \Omega} oldsymbol{v}_j oldsymbol{v}_j^ op
ight)^{-1} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} oldsymbol{v}_j$$

▶ For each movie $j \in [n]$,

$$oldsymbol{v}_j \coloneqq \left(\sum_{i \in [m]: (i,j) \in \Omega} oldsymbol{u}_i oldsymbol{u}_i^ op
ight)^{-1} \sum_{i \in [m]: (i,j) \in \Omega} a_{i,j} oldsymbol{u}_i.$$

ALTERNATING MINIMIZATION

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

- ▶ Somehow initialize $u_i \in \mathbb{R}^k$ for each $i \in [m]$ and $v_j \in \mathbb{R}^k$ for each $j \in [n]$.
- ► Repeat until convergence:
 - ▶ For each user $i \in [m]$,

$$oldsymbol{u}_i \coloneqq \left(\sum_{j \in [n]: (i,j) \in \Omega} oldsymbol{v}_j oldsymbol{v}_j^ op
ight)^{-1} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} oldsymbol{v}_j$$

▶ For each movie $j \in [n]$,

$$oldsymbol{v}_j \coloneqq \left(\sum_{i \in [m]: (i,j) \in \Omega} oldsymbol{u}_i oldsymbol{u}_i^ op
ight)^{-1} \sum_{i \in [m]: (i,j) \in \Omega} a_{i,j} oldsymbol{u}_i.$$

This assumes certain matrices involving u_i and v_j are invertible—but there's no guarantee of this!

MINIMIZE REGULARIZED TRAINING ERROR

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

Regularized training objective:

$$f(U, V) := \sum_{(i,j) \in \Omega} (a_{i,j} - \langle u_i, v_j \rangle)^2 + \lambda \left(\sum_{i=1}^m ||u_i||_2^2 + \sum_{j=1}^n ||v_j||_2^2 \right).$$

MINIMIZE REGULARIZED TRAINING ERROR

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

Regularized training objective:

$$f(U, V) := \sum_{(i,j) \in \Omega} (a_{i,j} - \langle u_i, v_j \rangle)^2 + \lambda \left(\sum_{i=1}^m ||u_i||_2^2 + \sum_{j=1}^n ||v_j||_2^2 \right).$$

New updates:

▶ For each user $i \in [m]$,

$$\boldsymbol{u}_i := \left(\sum_{j \in [n]: (i,j) \in \Omega} \boldsymbol{v}_j \boldsymbol{v}_j^\top + \lambda \boldsymbol{I}\right)^{-1} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} \boldsymbol{v}_j$$

lacktriangle (And similar for updating the $oldsymbol{v}_j$.)

MINIMIZE REGULARIZED TRAINING ERROR

Let $\Omega \subseteq [m] \times [n]$ be the entries for which $a_{i,j}$ is observed.

Regularized training objective:

$$f(U, V) := \sum_{(i,j) \in \Omega} (a_{i,j} - \langle u_i, v_j \rangle)^2 + \lambda \left(\sum_{i=1}^m ||u_i||_2^2 + \sum_{j=1}^n ||v_j||_2^2 \right).$$

New updates:

▶ For each user $i \in [m]$,

invertiable

$$oldsymbol{u}_i \coloneqq \left(\sum_{j \in [n]: (i,j) \in \Omega} oldsymbol{v}_j oldsymbol{v}_j^ op + \lambda oldsymbol{I}
ight)^{-1} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} oldsymbol{v}_j$$

• (And similar for updating the v_{i} .)

Prediction: For a movie j that user i has not rated, predict rating to be

$$\hat{a}_{i,j} := \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle.$$

OTHER THINGS

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- ▶ How to pick k or λ ?

OTHER THINGS

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- How to pick k or λ?
 Answer is obvious by now, but there are some subtle issues.

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- ▶ How to pick k or λ ?

Answer is obvious by now, but there are some subtle issues.

If λ = 0, then for updates to work, every user must have seen at least k movies, and every movie must have been rated by at least k users.

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- ▶ How to pick k or λ ?

Answer is obvious by now, but there are some subtle issues.

- If λ = 0, then for updates to work, every user must have seen at least k movies, and every movie must have been rated by at least k users.
- ▶ If $\lambda > 0$, then technically updates work for any k.

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- ▶ How to pick k or λ ?

Answer is obvious by now, but there are some subtle issues.

- If λ = 0, then for updates to work, every user must have seen at least k movies, and every movie must have been rated by at least k users.
- ▶ If $\lambda > 0$, then technically updates work for any k.
- ▶ Perhaps should have separate λ for each user/movie, and perhaps it should change as the u_i and v_j change!

- ▶ Initialization: can't initialize with all $u_i = v_j = 0$. (Try random.)
- ▶ How to pick k or λ ?

Answer is obvious by now, but there are some subtle issues.

- ▶ If $\lambda = 0$, then for updates to work, every user must have seen at least k movies, and every movie must have been rated by at least k users.
- ▶ If $\lambda > 0$, then technically updates work for any k.
- ▶ Perhaps should have separate λ for each user/movie, and perhaps it should change as the u_i and v_j change!
- Incorporating biases

$$egin{aligned} f(oldsymbol{U},oldsymbol{V},oldsymbol{b},oldsymbol{c},\mu) &:= \sum_{(i,j)\in\Omega} ig(a_{i,j} - \langle oldsymbol{u}_i,oldsymbol{v}_j
angle - b_i - c_j - \muig)^2 \ &+ \lambda \left(\sum_{i=1}^m \lVert oldsymbol{u}_i
Vert_2^2 + \sum_{j=1}^n \lVert oldsymbol{v}_j
Vert_2^2
ight). \end{aligned}$$

What about different loss functions?

Different training objective: logistic loss (each $a_{i,j} \in \{-1, +1\}$)

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j)\in\Omega} \ell_{\log}(a_{i,j}\langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle).$$

No closed form minimizer of $U \mapsto f(U, V)$.

What about different loss functions?

Different training objective: logistic loss (each $a_{i,j} \in \{-1, +1\}$)

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j) \in \Omega} \ell_{\log}(a_{i,j} \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle).$$

No closed form minimizer of $U \mapsto f(U, V)$.

Alternating gradient updates

For t = 1, 2, ...:

▶ For each user $i \in [m]$,

$$u_i^{(t+1)} := u_i^{(t)} - \eta \nabla_{u_i} \{ f(U^{(t)}, V^{(t)}) \}.$$

▶ For each movie $j \in [n]$,

$$\boldsymbol{v}_{j}^{(t+1)} \ \coloneqq \ \boldsymbol{v}_{j}^{(t)} - \eta \nabla_{\boldsymbol{v}_{j}} \Big\{ f(\boldsymbol{U}^{(t+1)}, \boldsymbol{V}^{(t)}) \Big\}.$$

What about different loss functions?

Different training objective: logistic loss (each $a_{i,j} \in \{-1, +1\}$)

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j) \in \Omega} \ell_{\log}(a_{i,j} \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle).$$

No closed form minimizer of $\boldsymbol{U} \mapsto f(\boldsymbol{U}, \boldsymbol{V})$.

Alternating gradient updates

For t = 1, 2, ...:

▶ For each user $i \in [m]$,

$$\mathbf{u}_{i}^{(t+1)} := \mathbf{u}_{i}^{(t)} + \eta_{t} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} \frac{1}{1 + \exp(a_{i,j} \langle \mathbf{u}_{i}^{(t)}, \mathbf{v}_{j}^{(t)} \rangle)} \mathbf{v}_{j}^{(t)}.$$

▶ For each movie $j \in [n]$,

$$\boldsymbol{v}_{j}^{(t+1)} := \boldsymbol{v}_{j}^{(t)} + \eta_{t} \sum_{i \in [m]: (i,j) \in \Omega} a_{i,j} \frac{1}{1 + \exp(a_{i,j} \langle \boldsymbol{u}_{i}^{(t)}, \boldsymbol{v}_{j}^{(t)} \rangle)} \boldsymbol{u}_{i}^{(t)}.$$

WHAT ABOUT DIFFERENT LOSS FUNCTIONS?

Different training objective: logistic loss (each $a_{i,j} \in \{-1, +1\}$)

$$f(\boldsymbol{U}, \boldsymbol{V}) := \sum_{(i,j) \in \Omega} \ell_{\log}(a_{i,j} \langle \boldsymbol{u}_i, \boldsymbol{v}_j \rangle).$$

No closed form minimizer of $U \mapsto f(U, V)$.

Alternating gradient updates

For t = 1, 2, ...:

▶ For each user $i \in [m]$,

$$\boldsymbol{u}_{i}^{(t+1)} := \boldsymbol{u}_{i}^{(t)} + \eta_{t} \sum_{j \in [n]: (i,j) \in \Omega} a_{i,j} \frac{1}{1 + \exp(a_{i,j} \langle \boldsymbol{u}_{i}^{(t)}, \boldsymbol{v}_{j}^{(t)} \rangle)} \boldsymbol{v}_{j}^{(t)}.$$

 $\blacktriangleright \ \, \text{For each movie} \,\, j \in [n]\text{,}$

$$\boldsymbol{v}_{j}^{(t+1)} \ := \ \boldsymbol{v}_{j}^{(t)} + \eta_{t} \sum_{i \in [m]: (i,j) \in \Omega} a_{i,j} \frac{1}{1 + \exp(a_{i,j} \langle \boldsymbol{u}_{i}^{(t)}, \boldsymbol{v}_{j}^{(t)} \rangle)} \boldsymbol{u}_{i}^{(t)}.$$

(Ambiguous whether to use $m{u}_i^{(t)}$ or $m{u}_i^{(t+1)}$ when updating $m{v}_j^{(t)} o m{v}_j^{(t+1)}$.)

Stochastic gradient descent

Unconstrained convex optimization

Unconstrained convex optimization problem

$$\min_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{x})$$

(f is the convex objective function).

Unconstrained convex optimization

Unconstrained convex optimization problem

$$\min_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{x})$$

(f is the convex objective function).

Gradient descent for differentiable objectives

- Start with some initial $oldsymbol{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For t = 1, 2, ... until some stopping condition is satisfied.
 - ▶ Compute gradient of f at $w^{(t)}$:

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)}).$$

Update:

$$\boldsymbol{w}^{(t+1)} := \boldsymbol{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

Unconstrained convex optimization

Unconstrained convex optimization problem

$$\min_{\boldsymbol{w} \in \mathbb{R}^d} f(\boldsymbol{x})$$

(f is the convex objective function).

Gradient descent for differentiable objectives

- Start with some initial $oldsymbol{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For t = 1, 2, ... until some stopping condition is satisfied.
 - ▶ Compute gradient of f at $w^{(t)}$:

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)}).$$

Update:

$$\boldsymbol{w}^{(t+1)} := \boldsymbol{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

If f is "nice", GD converges to global minimizer w_{\star} , where $\nabla f(w_{\star}) = 0$.

In machine learning, we frequently deal with objectives of the form

$$f(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{w})$$

where $f_1, f_2, \dots, f_n \colon \mathbb{R}^d \to \mathbb{R}$ are convex functions (e.g., loss on example i).

In machine learning, we frequently deal with objectives of the form

$$f(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{w})$$

where $f_1, f_2, \dots, f_n \colon \mathbb{R}^d \to \mathbb{R}$ are convex functions (e.g., loss on example i).

Example: Gradient descent algorithm for logistic regression

- ▶ Start with some initial $w^{(1)} \in \mathbb{R}^d$.
- ▶ For t = 1, 2, ...

$$egin{align*} oldsymbol{w}^{(t+1)} &:= oldsymbol{w}^{(t)} - \eta_t
abla f(oldsymbol{w}^{(t)}) & ext{training loss} \ &= oldsymbol{w}^{(t)} + \eta_t rac{1}{|S|} \sum_{(oldsymbol{x}, y) \in S} rac{1}{1 + e^{y \langle oldsymbol{w}^{(t)}, oldsymbol{x}
angle}} y oldsymbol{x}. \end{split}$$

In fact, we often really care about

$$f(\boldsymbol{w}) = \mathbb{E}[f(\boldsymbol{w}; Z)]$$

where Z is a random example, and $f(\cdot;Z)$ represents the (convex) loss on this random example.

test loss

In fact, we often really care about

$$f(\boldsymbol{w}) = \mathbb{E}[f(\boldsymbol{w}; Z)]$$

where Z is a random example, and $f(\cdot;Z)$ represents the (convex) loss on this random example.

Let us continue with this for now . . .

Suppose at the push of a red button, we get an independent draw from the distribution ${\cal P}$ over labeled examples that we care about.



Suppose at the push of a red button, we get an independent draw from the distribution P over labeled examples that we care about.



To make things even simpler, before the t-th button push, we must specify our current weight vector $\boldsymbol{w}^{(t)}$, and all we get back is

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)}; z_t)$$

where z_t is a independent random realization of Z.

Suppose at the push of a red button, we get an independent draw from the distribution P over labeled examples that we care about.



To make things even simpler, before the t-th button push, we must specify our current weight vector $\boldsymbol{w}^{(t)}$, and all we get back is

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)}; z_t)$$

where z_t is a independent random realization of Z.

(Gradient of the loss on a simple random example at current weight vector.)

Suppose at the push of a red button, we get an independent draw from the distribution P over labeled examples that we care about.



To make things even simpler, before the t-th button push, we must specify our current weight vector $w^{(t)}$, and all we get back is

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)} z_t)$$

not overall sample

where z_t is a independent random realization of Z.

(Gradient of the loss on a simple random example at current weight vector.)

Can we still find a good weight vector?

STOCHASTIC GRADIENT DESCENT (SGD)

randomly pick one and minimize its cost

- Start with some initial $oldsymbol{w}^{(1)} \in \mathbb{R}^d$.
- ▶ For t = 1, 2, ...
 - ▶ Push the red button and get

$$\boldsymbol{\lambda}^{(t)} := \nabla f(\boldsymbol{w}^{(t)}; z_t).$$

Update:

$$\boldsymbol{w}^{(t+1)} \coloneqq \boldsymbol{w}^{(t)} - \eta_t \boldsymbol{\lambda}^{(t)}.$$

just for one sample point each time

WHY DOES THIS WORK?

We care about

$$f(\boldsymbol{w}) = \mathbb{E}[f(\boldsymbol{w}; Z)].$$

We care about

$$f(\boldsymbol{w}) = \mathbb{E}[f(\boldsymbol{w}; Z)].$$

We update $oldsymbol{w}^{(t)}$ using $oldsymbol{\lambda}^{(t)}$, where

$$\mathbb{E}\left[\boldsymbol{\lambda}^{(t)}\right] = \mathbb{E}\left[\nabla f(\boldsymbol{w}^{(t)}; z_t)\right] = \nabla \left\{\mathbb{E}\left[f(\boldsymbol{w}^{(t)}; z_t)\right]\right\} = \nabla f(\boldsymbol{w}^{(t)})$$

- **Expectation** is only w.r.t. z_t , an independent realization of Z.
- lacktriangle Gradient is with respect to $m{w}$ argument, and is evaluated at $m{w} = m{w}^{(t)}$.

We care about

$$f(m{w}) = \mathbb{E} ig[f(m{w}; Z) ig].$$
 average. not a random example

We update $oldsymbol{w}^{(t)}$ using $oldsymbol{\lambda}^{(t)}$, where

$$\mathbb{E}\left[\boldsymbol{\lambda}^{(t)}\right] = \mathbb{E}\left[\nabla f(\boldsymbol{w}^{(t)}; z_t)\right] = \nabla \left\{\mathbb{E}\left[f(\boldsymbol{w}^{(t)}; z_t)\right]\right\} = \nabla f(\boldsymbol{w}^{(t)})$$

- \blacktriangleright Expectation is only w.r.t. z_t , an independent realization of Z.
- lacktriangle Gradient is with respect to w argument, and is evaluated at $w=w^{(t)}$.

In expectation, we move in the correct direction w.r.t. the objective we actually care about!

online perceptron use this idea

We care about

$$f(\boldsymbol{w}) = \mathbb{E}[f(\boldsymbol{w}; Z)].$$

We update $oldsymbol{w}^{(t)}$ using $oldsymbol{\lambda}^{(t)}$, where

$$\mathbb{E}\left[\boldsymbol{\lambda}^{(t)}\right] = \mathbb{E}\left[\nabla f(\boldsymbol{w}^{(t)}; z_t)\right] = \nabla \left\{\mathbb{E}\left[f(\boldsymbol{w}^{(t)}; z_t)\right]\right\} = \nabla f(\boldsymbol{w}^{(t)})$$

- \blacktriangleright Expectation is only w.r.t. z_t , an independent realization of Z.
- lacktriangle Gradient is with respect to w argument, and is evaluated at $w=w^{(t)}$.

In expectation, we move in the correct direction w.r.t. the objective we actually care about!

Convergence is slower than Gradient Descent due to stochasticity.

We care about

the expectation of Vambda $f(\boldsymbol{w}) \ = \ \mathbb{E}[f(\boldsymbol{w};Z)].$ (rather than directly calculate)

We update $oldsymbol{w}^{(t)}$ using $oldsymbol{\lambda}^{(t)}$, where

$$\mathbb{E}\left[\boldsymbol{\lambda}^{(t)}\right] = \mathbb{E}\left[\nabla f(\boldsymbol{w}^{(t)}; z_t)\right] = \nabla \left\{\mathbb{E}\left[f(\boldsymbol{w}^{(t)}; z_t)\right]\right\} = \nabla f(\boldsymbol{w}^{(t)})$$

- **Expectation** is only w.r.t. z_t , an independent realization of Z.
- lacktriangle Gradient is with respect to w argument, and is evaluated at $w=w^{(t)}$.

In expectation, we move in the correct direction w.r.t. the objective we actually care about!

Convergence is slower than Gradient Descent due to stochasticity. But it is w.r.t. the objective we actually care about!

Assumption: training data is i.i.d. sample from distribution we care about.

- Assumption: training data is i.i.d. sample from distribution we care about.
- ► Randomly permute order of the data.

- Assumption: training data is i.i.d. sample from distribution we care about.
- ► Randomly permute order of the data.
- ▶ Pushing the red button: take the next training example z_t (according to the random order) and use it to compute gradient of $f(\cdot; z_t)$ at $w^{(t)}$.

- Assumption: training data is i.i.d. sample from distribution we care about.
- ► Randomly permute order of the data.
- ▶ Pushing the red button: take the next training example z_t (according to the random order) and use it to compute gradient of $f(\cdot; z_t)$ at $w^{(t)}$.
- After making a pass through the data, theory says to stop (afaik)—repeating examples cannot really be regarded as independent.

- Assumption: training data is i.i.d. sample from distribution we care about.
- ► Randomly permute order of the data.
- ▶ Pushing the red button: take the next training example z_t (according to the random order) and use it to compute gradient of $f(\cdot; z_t)$ at $w^{(t)}$.
- ► After making a pass through the data, theory says to stop (afaik)—repeating examples cannot really be regarded as independent.
- ▶ In practice, can help to cycle through the data a few times.

there is no risk of overfitting!!!

- very resistant to overfitting.
- 2. voted-perceptron

SGD WITH NON-CONVEX OBJECTIVES

Can also apply SGD to non-convex summation objectives, e.g.,

$$f(\boldsymbol{U}, \boldsymbol{V}) := \mathbb{E}[\ell_{\log}(a_{I,J}\langle \boldsymbol{u}_I, \boldsymbol{v}_J \rangle)]$$

(for random user I and movie J).

SGD WITH NON-CONVEX OBJECTIVES

Can also apply SGD to non-convex summation objectives, e.g.,

$$f(\boldsymbol{U}, \boldsymbol{V}) := \mathbb{E}[\ell_{\log}(a_{I,J}\langle \boldsymbol{u}_I, \boldsymbol{v}_J \rangle)]$$

(for random user I and movie J).

Given ratings $\{a_{i,j}:(i,j)\in\Omega\}$, first randomly permute their order in Ω .

SGD WITH NON-CONVEX OBJECTIVES

Can also apply SGD to non-convex summation objectives, e.g.,

$$f(\boldsymbol{U}, \boldsymbol{V}) := \mathbb{E} [\ell_{\log}(a_{I,J}\langle \boldsymbol{u}_I, \boldsymbol{v}_J \rangle)]$$

(for random user I and movie J).

Given ratings $\{a_{i,j}:(i,j)\in\Omega\}$, first randomly permute their order in Ω .

For t = 1, 2, ...:

- ▶ Get next (i, j) and $a_{i,j}$.
- ► Updates:

$$egin{aligned} m{u}_i^{(t+1)} &:= m{u}_i^{(t)} + \eta_t a_{i,j} rac{1}{1 + \exp(a_{i,j} \langle m{u}_i^{(t)}, m{v}_j^{(t)}
angle)} m{v}_j^{(t)} \ m{v}_j^{(t+1)} &:= m{v}_j^{(t)} + \eta_t a_{i,j} rac{1}{1 + \exp(a_{i,j} \langle m{u}_i^{(t)}, m{v}_j^{(t)}
angle)} m{u}_i^{(t)}. \end{aligned}$$

 $\blacktriangleright \text{ (For all other } (i',j') \text{, no change: } \boldsymbol{u}_{i'}^{(t+1)} = \boldsymbol{u}_{i'}^{(t)} \text{ and } \boldsymbol{v}_{j'}^{(t+1)} = \boldsymbol{v}_{j'}^{(t)}.)$

RECAP

- Stochastic gradient descent (not technically a descent method): for convex problems, converges (somewhat slowly) to the optimum of the problem we actually care about!
- ► Can also be applied to non-convex problems.
- ▶ A particular non-convex problem: matrix completion, which is often used for collaborative filtering / recommender systems.