

Scaling ML to huge datasets

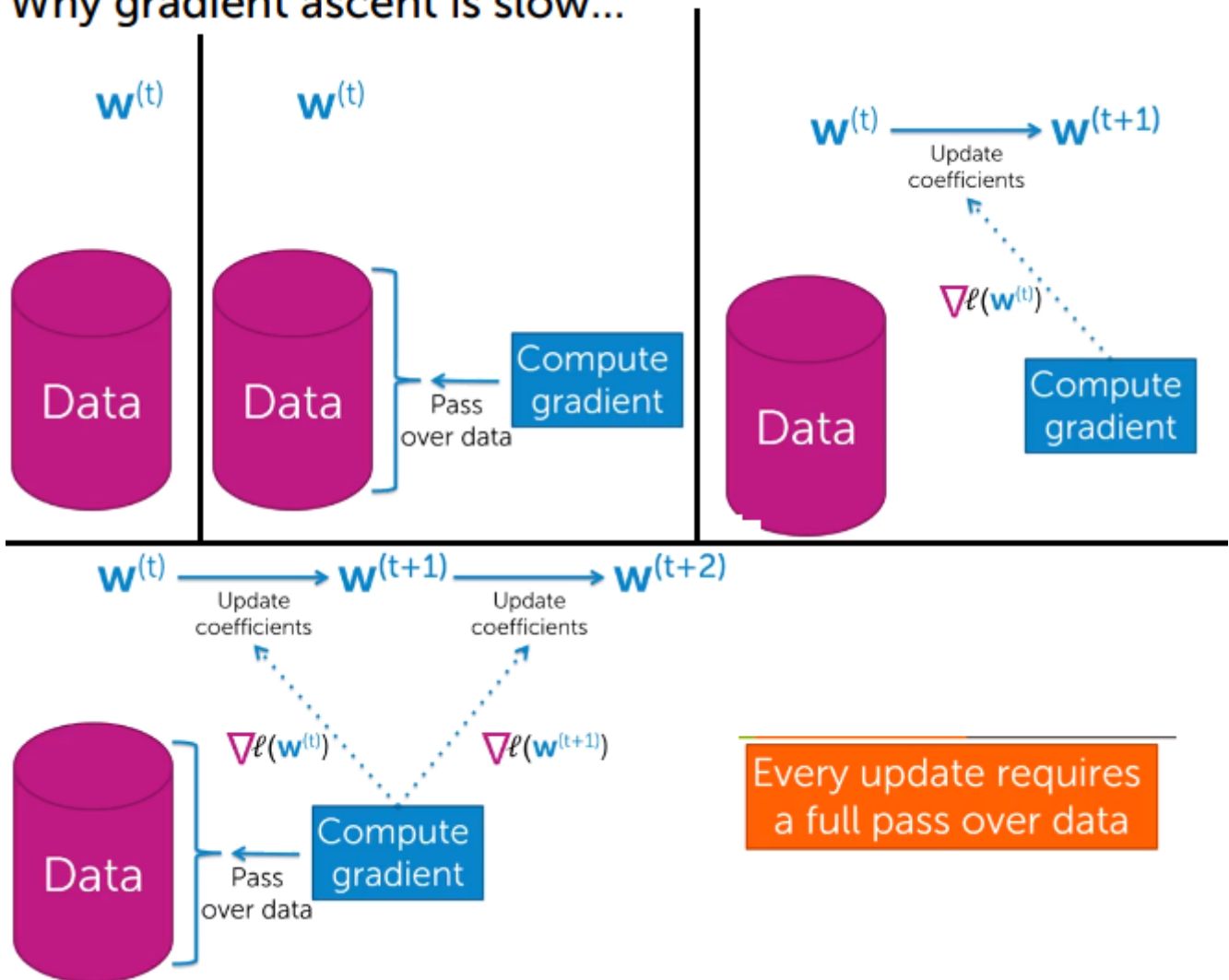
Gradient Ascent will not scale to today's huge datasets

Gradient Ascent

- Consider we have a large dataset (DATA) and a set of coefficients ($\mathbf{w}^{(t)}$) to be updated.
- Using gradient ascent compute the gradient on the dataset. This requires to make a pass or a scan over the data, computing the contribution of each of these data points to the gradient.
- Then compute the gradient and update the coefficients / parameters and get ($\mathbf{w}^{(t+1)}$).
- Then return to the dataset and make another pass over each data-points and compute a new gradient and update the parameters and coefficients ($\mathbf{w}^{(t+2)}$).

There every time a coefficients needs to be updated, then we need to perform a full scan or a full pass over the entire dataset. This process can be 'very slow' if 'dataset is very big'.

Why gradient ascent is slow...



Datasets are getting huge, and we need them!

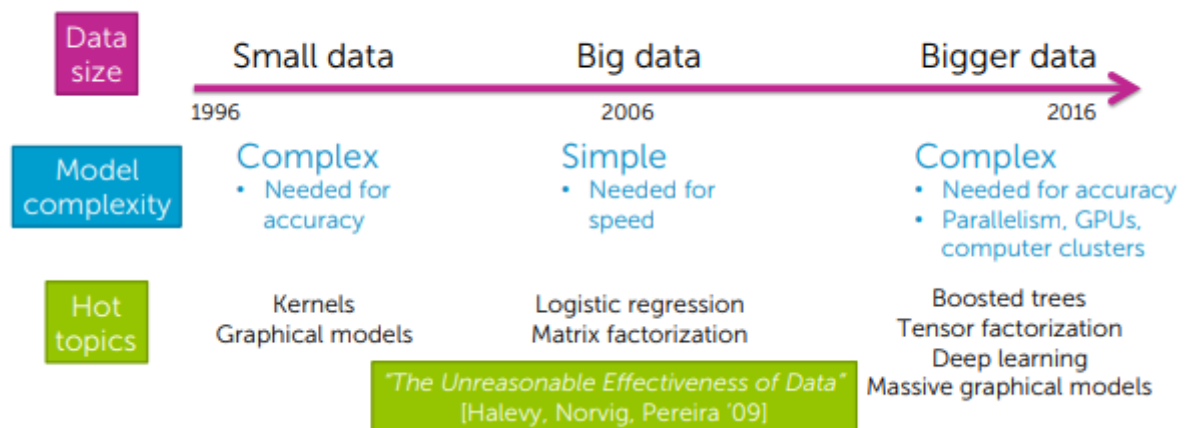
- Internet -> WWW - 4.8B webpages
- twitter -> 500M Tweets/day
- Internet of Things -> Sensors everywhere
- Youtube -> 300 hours uploaded/min

- 1B users
- Ad revenue
- 5B views/day
- (Need ML algorithm to learn from billions of video views every day, & to recommend ads within milliseconds).

Timeline of scalable machine learning and 'Stochastic Gradient'

Machine Learning improves significantly with bigger datasets

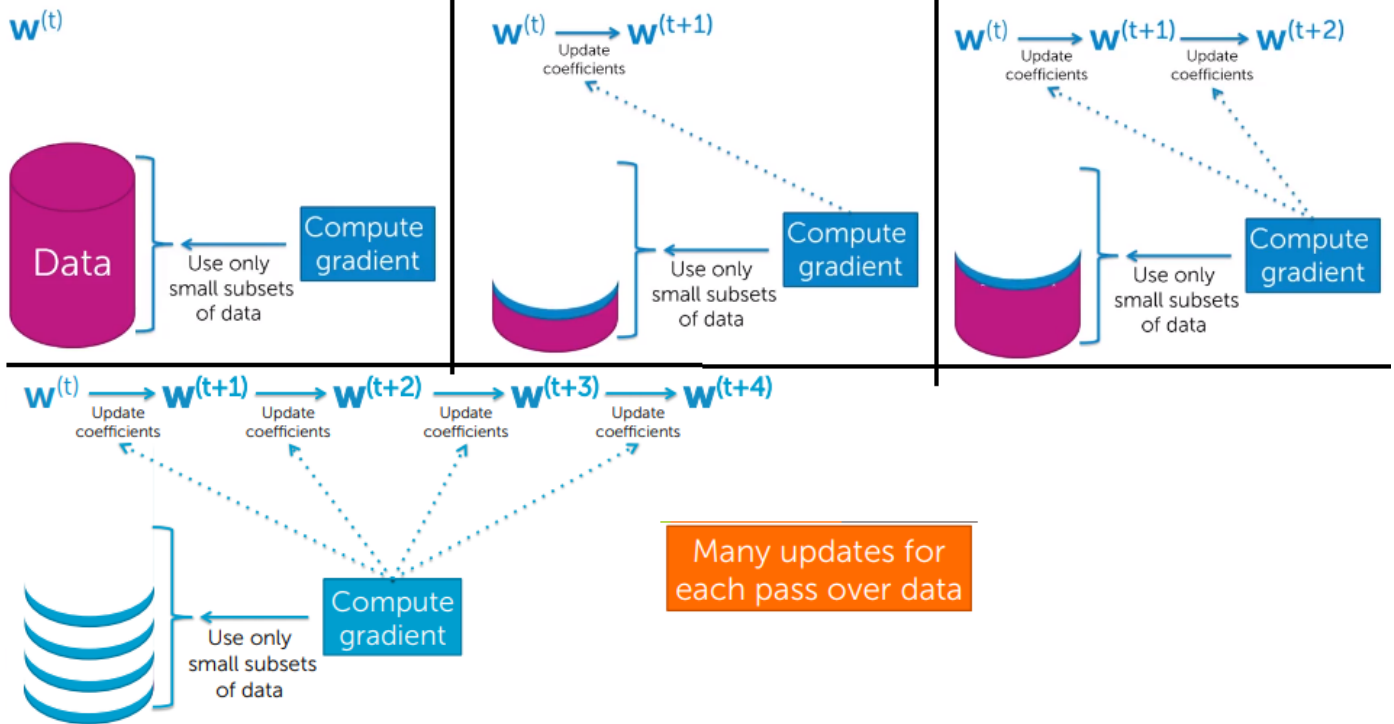
- ML has evolved with the size of the datasets over the decades.



Stochastic Gradient Ascent

- Small change to the Gradient ascent algorithm.
- Takes a massive dataset and current parameters ($\mathbf{w}(\mathbf{t})$).
- Computes the gradient for only a small subset of data rather than the entire dataset. Then it updates the coefficients.

Stochastic gradient ascent



Scaling ML with stochastic gradient

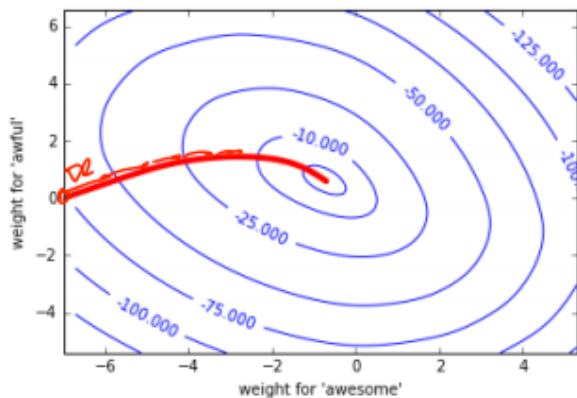
Why gradient ascent won't scale?

- Learning one data point at a time.
- Start at some point, and compute the gradient descent of the likelihood function and take a step in that space.
- Then compute the gradient again until the convergence is met, or stop condition.

Gradient descent expenses

- Sum over datapoint the contribution of each datapoint $x_i y_i$ to gradient.

Gradient ascent



Algorithm:

```
while not converged
     $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta \nabla \ell(\mathbf{w}^{(t)})$ 
```

How expensive is gradient ascent?

Sum over data points

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N h_j(\mathbf{x}_i) \left(\mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}) \right)$$

Contribution of data point \mathbf{x}_i, y_i to gradient

$$\frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$$

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N \frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$$

Time to compute contribution of \mathbf{x}_i, y_i	# of data points (N)	Total time to compute 1 step of gradient ascent
1 millisecond	1000	1 sec
1 second	1000	16.7 min
1 millisecond	10 million	2.8 hours
1 millisecond	10 billion	115.7 days

Stochastic gradient

- Learning one data point at a time.
- Use one data-point to compute the gradient.
- Each time, while computing the gradient use a new data point i .
- Therefore, it will be computationally cheaper to perform stochastic gradient ascent when compared to gradient ascent.

gradient ascent for logistic regression

```
init  $\mathbf{w}^{(1)}=0$ ,  $t=1$ 
until converged
    for  $j=0, \dots, D$ 
        partial[j] =  $\sum_{i=1}^N h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}) \right)$ 
         $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} + \eta$  partial[j]
     $t \leftarrow t + 1$ 
```

Sum over data points

Comparing computational time per step

Gradient ascent

Stochastic gradient ascent

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N \frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$$

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} \approx \frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$$

Time to compute contribution of $\mathbf{x}_i y_i$	# of data points (N)	Total time for 1 step of gradient	Total time for 1 step of stochastic gradient
1 millisecond	1000	1 second	1 millisecond
1 second	1000	16.7 minutes	1 sec
1 millisecond	10 million	2.8 hours	1 millisecond
1 millisecond	10 billion	115.7 days	1 millisecond

Stochastic gradient ascent for logistic regression

```
init  $\mathbf{w}^{(1)}=0$ ,  $t=1$ 
until converged
    for  $i=1, \dots, N$ 
        for  $j=0, \dots, D$ 
            partial[j] =  $h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}) \right)$ 
             $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} + \eta$  partial[j]
         $t \leftarrow t + 1$ 
```

Each time, pick different data point i

Instead of all data points for gradient, use 1 data point only

Sum over data points

Gradient ascent: $\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N \frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$

Stochastic gradient ascent: $\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} \approx \frac{\partial \ell_i(\mathbf{w})}{\partial \mathbf{w}_j}$

Each time, pick different data point i

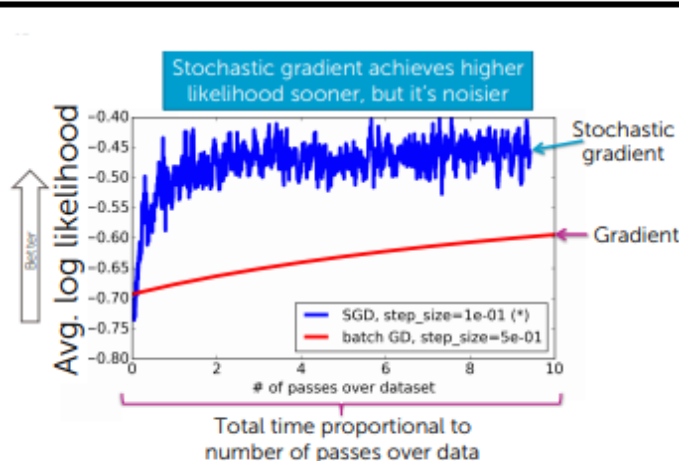
Compare gradient to stochastic gradient

- Stochastic gradient is computationally cheaper for large datasets.
- In practice stochastic gradient is often faster than gradient.
- But stochastic gradient is **'very sensitive to parameters'**. Choice of step size, etc.
- A graph is drawn between x-axis (# passes over dataset) & y-axis (average likelihood);
 - higher the likelihood, better it is.

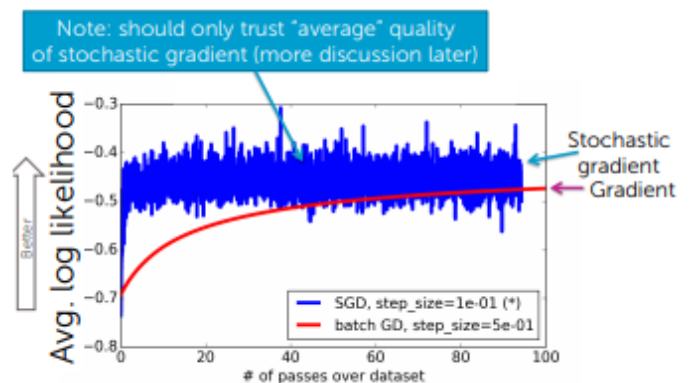
- From the graph, after 10 passes, gradient has much lower likelihood when compared to stochastic gradient.
- stochastic gradient achieves higher likelihood sooner, but it is noisier.
- stochastic gradient converges faster but oscillates more.
- eventually the gradient descent catches up. It is smooth.

Comparing gradient to stochastic gradient

Algorithm	Time per iteration	Total time to convergence for large data		Sensitivity to parameters
		In theory	In practice	
Gradient	Slow for large data	Slower	Often slower	Moderate
Stochastic gradient	Always fast	Faster	Often faster	Very high



Eventually, gradient catches up

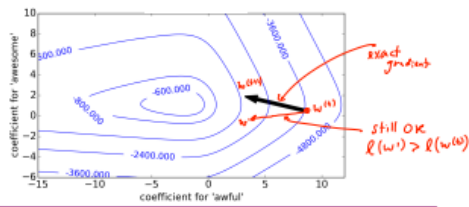


Understanding why stochastic gradient works

Why would stochastic gradient ever work?

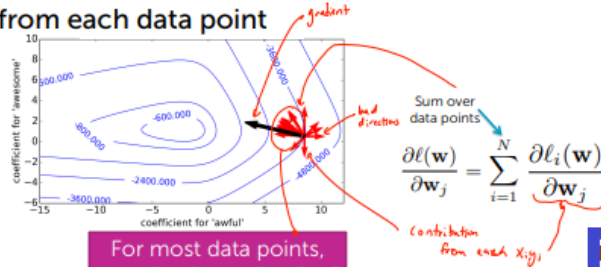
- The Gradient is the "best" direction, but any direction that goes "up" would be useful.
- In machine learning the steepest direction is the sum of "little directions" from each data-point. For most data-points contributions points "up" - ascent.
- In case of the gradient - the gradient is the likelihood of sum over all data-points.
- In case of the stochastic gradient - pick a data-point and move in direction.
- In case stochastic gradient - total likelihood will increase.

Gradient is direction of steepest ascent



Gradient is "best" direction, but any direction that goes "up" would be useful

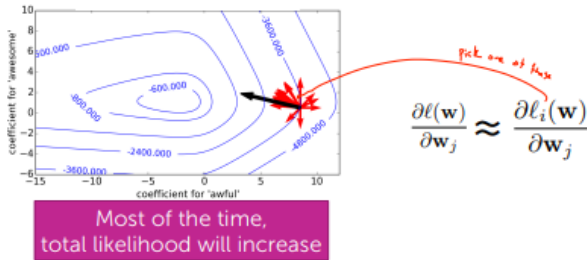
In ML, steepest direction is sum of "little directions" from each data point



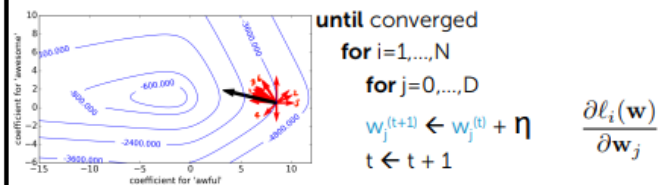
For most data points, contribution points "up"

Rather than take the sum of little data-points, in case of stochastic gradient case pick a data-point and move in direction.

Stochastic gradient: pick a data point and move in direction



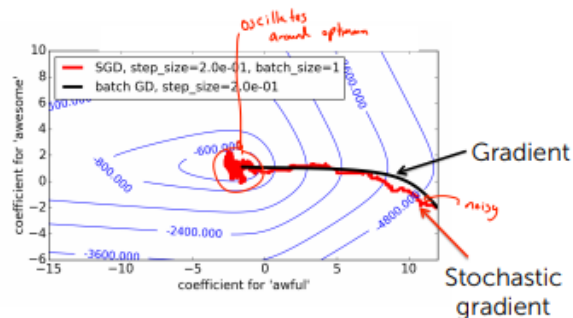
Stochastic gradient ascent: Most iterations increase likelihood, but sometimes decrease it → On average, make progress



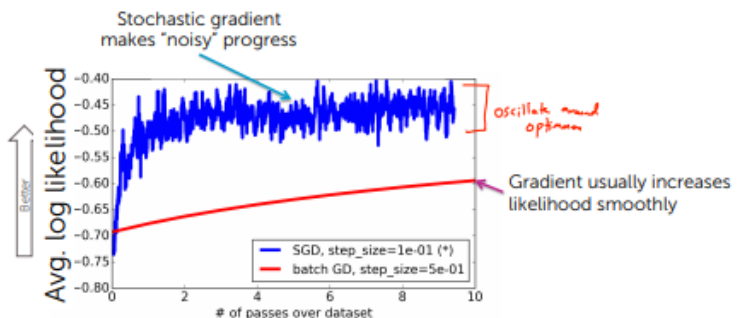
Convergence path

- Gradient finds the direction of steepest ascent. It converges smoothly.
- Gradient is the sum of contributions from each data-point.
- Stochastic gradient is the sum of contributions from each data-point.
- On an average increase likelihood and sometimes decrease.
- Stochastic gradient has a noisy convergence. It oscillates around the optimum. (Issue);

Convergence paths



Stochastic gradient convergence is "noisy"



Stochastic Gradient : Practical Tricks

Shuffle data before running stochastic gradient

- If the data is implicitly sorted, it can cause issues.
- Systematic order of data (like y is -1 comes before all the y is +1) can introduce bias.
- Therefore shuffle the dataset. No long regions in the dataset.

Algorithm:

1. Shuffle the data
2. initialize the coefficients/parameters = 0 at $t=1$
3. until convergence
 - for any dataset i
 - for the features j
 - compute the new coefficient $\rightarrow w_j(t+1) \leftarrow w_j(t) + (\text{step-size}) * \text{gradient of data-point};$
4. optimum

Stochastic gradient ascent

Shuffle data ← Before running stochastic gradient, make sure data is shuffled

```

init  $w^{(1)}=0$ ,  $t=1$ 
until converged
  for  $i=1, \dots, N$ 
    for  $j=0, \dots, D$ 
       $w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \frac{\partial \ell_i(w)}{\partial w_j}$ 
       $t \leftarrow t + 1$ 
  
```

Order of data can introduce bias

$x[1] = \text{\#awesome}$	$x[2] = \text{\#awful}$	$y = \text{sentiment}$
0	2	-1
3	3	-1
2	4	-1
0	3	-1
0	1	-1
2	1	+1
4	1	+1
1	1	+1
2	1	+1

Stochastic gradient updates parameters 1 data point at a time

Systematic order in data can introduce significant bias, e.g., all negative points first, or temporal order, younger first, or ...

Choosing a step-size

- Picking step-size for **stochastic gradient** is very similar to picking a step-size for gradient.
- But since stochastic gradient is more noiser and fluctuates more frequently, it is a lot more unstable.
- If step-size is too small, stochastic gradient slows to converge.
- If step-size is too large, stochastic gradient oscillates a lot.
- If step-size is very large, stochastic gradient goes crazy.

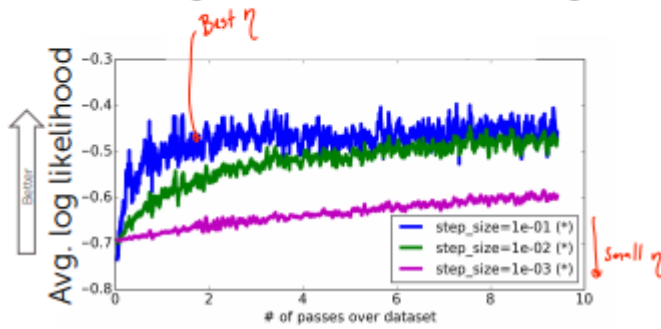
Simple rule of thumb for picking step size η similar to gradient

- It requires a lot of trial and error.
- Try several values exponentially spaced.
 - Goal: plot learning curves to
 - find one η that is too small.
 - find one η that is too large.
- Advanced tip : step-size that decreases with iterations is very important for stochastic gradient.

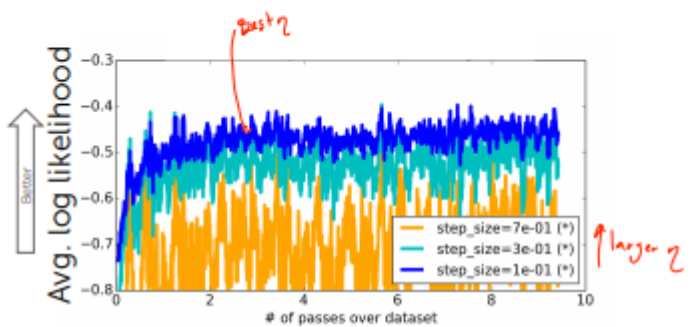
$\eta(t) = \eta_0 / t$; ($\eta_0 \rightarrow \text{constant}$; $t \rightarrow \text{iteration \#}$; $\eta(t) \rightarrow \text{step-size for that iteration}$);

Choosing the step size η

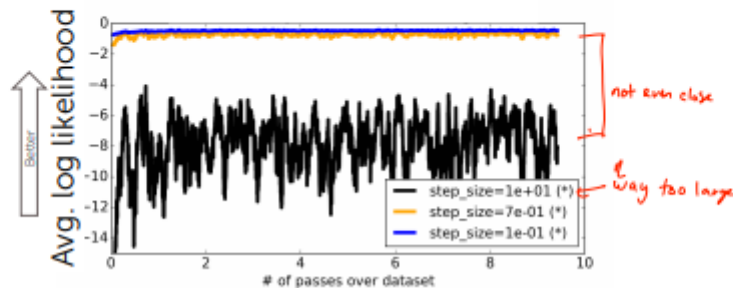
If step size is too small,
stochastic gradient slow to converge



If step size is too large,
stochastic gradient oscillates



If step size is very large,
stochastic gradient goes crazy ☹

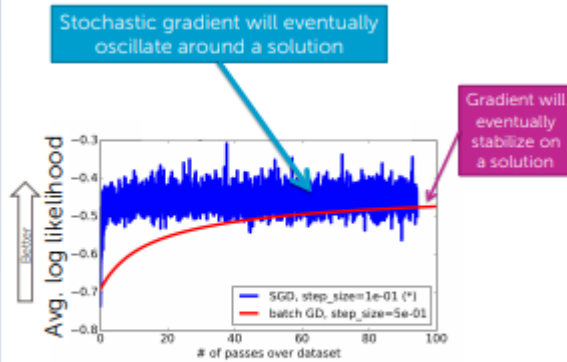


Stochastic Gradient convergence

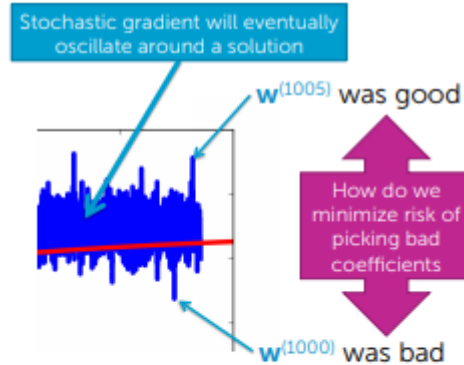
- Stochastic gradient oscillates vastly over the minimum. Therefore cannot trust the last parameter/coefficient.
- The gradient will eventually stabilize must the stochastic gradient oscillates around the solution.
- Stochastic gradient return the average coefficients.

$$\hat{\mathbf{w}} = (1/T) \sum_{t=1, \dots, T} \mathbf{w}(t);$$

Stochastic gradient never fully "converges"



The last coefficients may be really good or really bad!! ☹



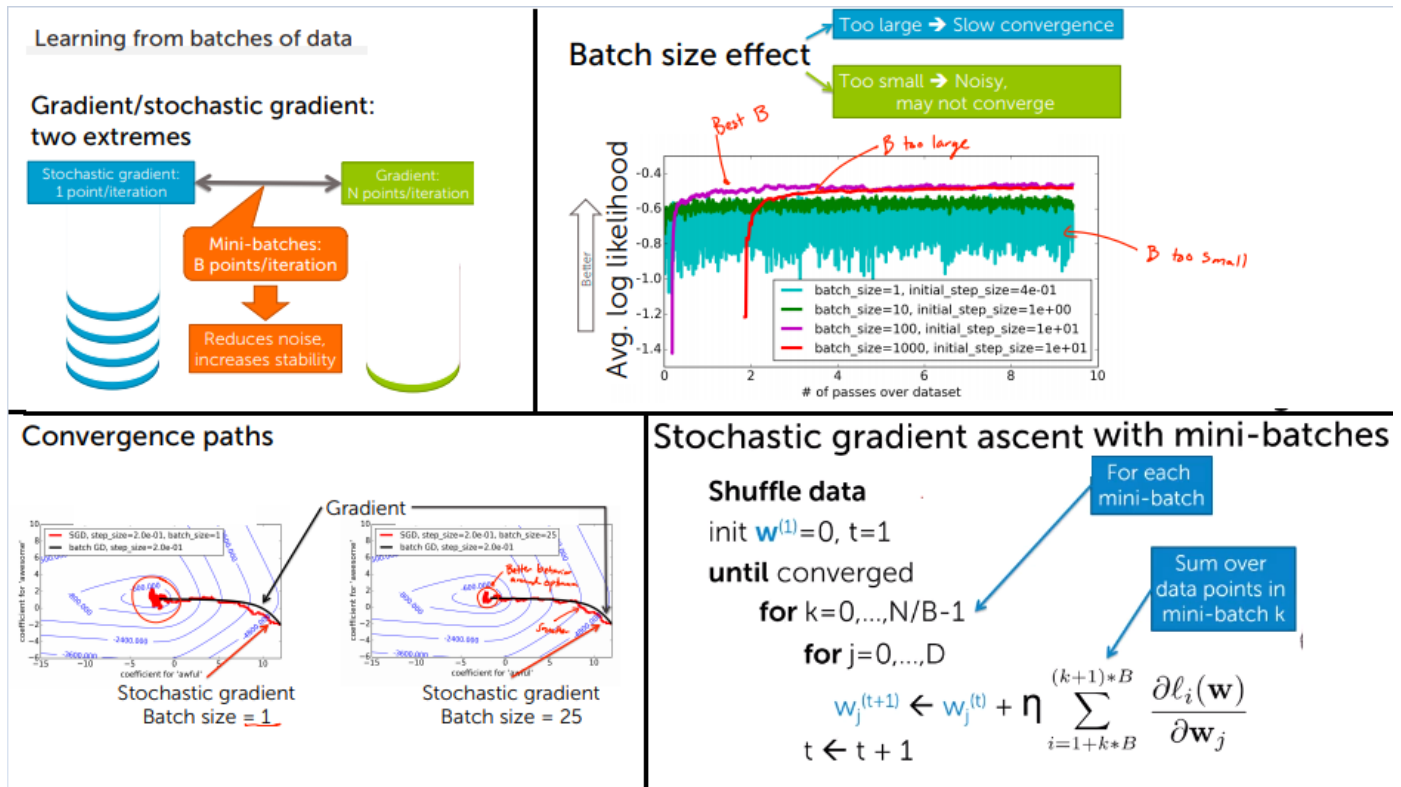
Stochastic gradient returns average coefficients

$$\hat{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

- Minimize noise: don't return last learned coefficients
- Instead, output average:

Learning from batches of data

- Learning from just one data-point is too noisy, therefore prefer mini-batches.
- The gradient descent -> pass over all data-points and stochastic gradient -> uses only one data-point are at the extremes. Require something in between.
- **Mini-batches:** B points/iteration -> reduces noise, increases stability.
- Stochastic gradient with batch-size of 25 when compared to a batch-size of 1 has a smoother curve towards convergence and oscillates less at the optimum.
- Batch-size must not be too-large (gradient - slow to converge) or too small (batch-size -1 -> noisy);



Measuring Convergence

- Stochastic gradient gets to the optimum before a single pass over the data. While the gradient takes 100 or more passes over the data.
- In order to get one point on the plot (#passes over dataset v/s avg likelihood) then if we had to compute the whole likelihood over the entire dataset, then it would require to compute the product over all the dataset, makes the process much slower.
- Therefore measure the log-likelihood based on the probability.

Estimating the log-likelihood with sliding window

- For every iteration t can compute the likelihood of a particular data point.
- We can't use this value to measure the conversion because it would do well on one data point classified perfectly but not for others -> leading to noise.
- In order to compute the progress after certain iterations. Eg $t = 75$, report the average of the last k values (Take the likelihood for the last few datapoints, average it and create a smoother curve).

For every timestamp need to keep an average of the last few likelihoods in order to measure the convergence.

- Minibatches on size = 100. Converges faster than gradient.
- In order to draw the blue line -> need to average the likelihood over the last 30 data-points.

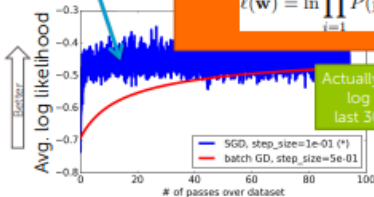
Measuring convergence

Need to compute log likelihood of data at every iteration???

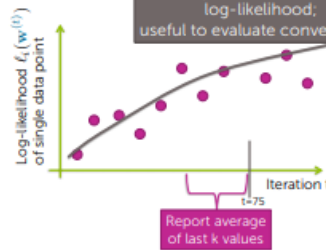
→ Really really slow, product over all data points!

$$\ell(\mathbf{w}) = \ln \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$$

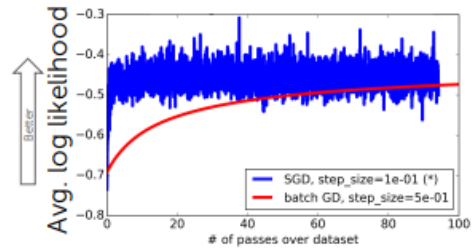
Actually, plotting average log likelihood over last 30 mini-batches...



Estimate log-likelihood with sliding window



That's what average log-likelihood meant... 😊
(In this case, over last $k=30$ mini-batches, with batch-size $B=100$)



Computing log-likelihood during run of stochastic gradient ascent

init $\mathbf{w}^{(1)}=0, t=1$
until converged

for $i=1, \dots, N$

for $j=0, \dots, D$

partial[j] = $h_j(\mathbf{x}_i) (\mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}))$

$\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} + \eta \text{partial}[j]$

$t \leftarrow t + 1$

Log-likelihood of data point i is simply:

$$\ell_i(\mathbf{w}^{(t)}) = \begin{cases} \ln P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}), & \text{if } y_i = +1 \\ \ln (1 - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)})), & \text{if } y_i = -1 \end{cases}$$

Adding regularization

- Stochastic gradient -> take a contribution for a single data point and add the contributions we get the gradient. Sum of the stochastic gradients is equal to gradient.
- Total cost for the algorithm
 - Total quality = measure of fit - measure of magnitude of coefficients
 - Total quality = log data likelihood - L2 penalty
 - Total quality = $l(\mathbf{w}) + (\lambda/2) \|\mathbf{w}\|^2$
 - λ -> tuning/regularization parameter;
- Gradient - Derivative of the Total cost
 - Total derivative = sum of the derivative of likelihood - $\lambda \mathbf{w}$;
- Stochastic gradient -> gradient over a data-point or a batch
 - Need to introduce the regularization parameter while computing stochastic gradient.
 - Consider each datapoint contributes $1/N$ to regularization.
 - Adding all the stochastic gradients will give the gradient descent.

Adding regularization

Consider specific total cost

max w

Total quality =
measure of fit - measure of magnitude of coefficients

$\ell(w)$
log data likelihood

$\|w\|_2^2$
L2 penalty

Gradient of L₂ regularized log-likelihood

Total quality =
measure of fit - measure of magnitude of coefficients

$\ell(w)$

$\lambda \|w\|_2^2$

Total derivative = $\sum_{i=1}^N \frac{\partial \ell_i(w)}{\partial w_j} - 2 \lambda w_j$

Stochastic gradient for regularized objective

Total derivative = $\sum_{i=1}^N \frac{\partial \ell_i(w)}{\partial w_j} - 2 \lambda w_j$

• What about regularization term?

Each time, pick different data point i

Stochastic gradient ascent

Total derivative $\approx \frac{\partial \ell_i(w)}{\partial w_j} - \frac{2}{N} \lambda w_j$

Each data point contributes 1/N to regularization

Stochastic gradient ascent with regularization

Shuffle data

init $w^{(1)}=0$, $t=1$

until converged

for $i=1, \dots, N$

for $j=0, \dots, D$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \left(\frac{\partial \ell_i(w)}{\partial w_j} - \frac{2}{N} \lambda w_j \right)$$

$t \leftarrow t + 1$

For a mini-batch B

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \left(\frac{\partial \ell_i(w)}{\partial w_j} - \frac{2}{N_B} \lambda w_j \right)$$

$t \leftarrow t + 1$

Online Learning

The online learning task

- Batch learning -> All data is available at the start of training time.
- Online learning -> Data arrives (streams in) over time.
 - Must train the model as data arrives.
 - Compute the coefficients as data is feed to the model.

Online Learning Example:

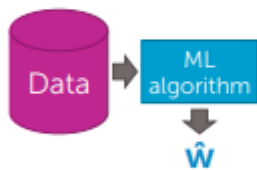
Ad targeting

- Consider visiting a websites, that shows a lot of ads.
 - The input -> user info -age, previous websites, etc -> who visits the pages is fed to the machine learning algorithm.
 - This machine learning algorithm is gonna use some set of coefficients $\hat{w}(t)$ and predict the best ads to show.
 - The ads shown are -> **y-hat (suggested ads)**.
 - In case the user clicks any ad shown, then machine learning algorithm figure that the user chose a particular ad.
 - The ML assign y_t , true label for the particular chosen ad in the website.
 - The ML algorithm updates the coefficient $w(t) \rightarrow w(t+1)$;

Batch vs online learning

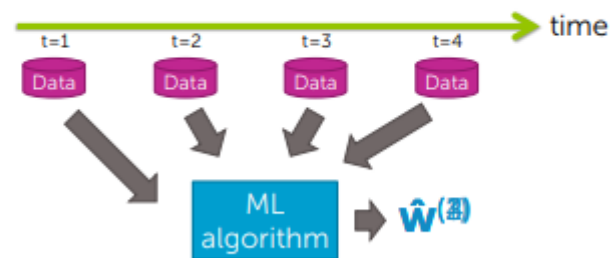
Batch learning

- All data is available at start of training time

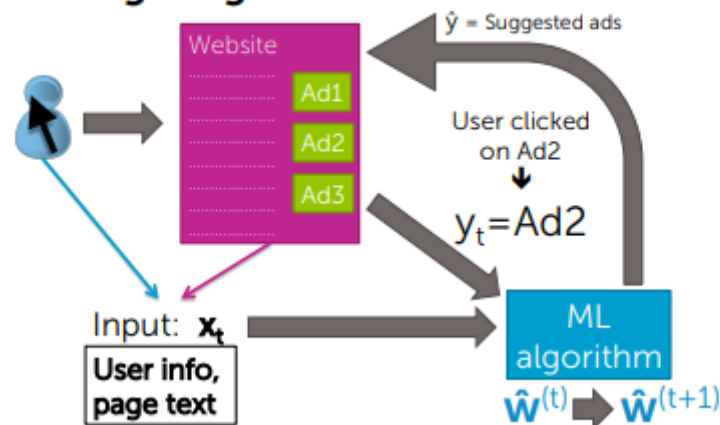


Online learning

- Data arrives (streams in) over time
– Must train model as data arrives!



Ad targeting Online learning example:



Using stochastic-gradient for online learning

Online learning problem

- Data arrives over each time step t :
 - **Observe input x_t**
 - Info of user, text of webpage, etc;
 - **Make a prediction $\hat{y}(t)$**
 - Which ad to show
 - **Observe true output $y(t)$**
 - Which ad user clicked on
 - Feed the information **ML Algorithm** to update coefficients each time.

Stochastic gradient ascent can be used for online learning.

As Stochastic gradient uses a small number of data-points to compute coefficients it can be used for online learning.

- init $w(1) = 0$, $t = 1$
- Each time step t :
 - Observe input $x(t)$
 - Make a prediction $\hat{y}(t)$
 - Observe true output $y(t)$
 - Update coefficients

- for $j = 0, \dots, D$ $w_j(t+1) \leftarrow w_j(t) + \eta * (\text{gradient likelihood})$;

Stochastic gradient ascent can be used for online learning!!!

- init $\mathbf{w}^{(1)}=0$, $t=1$
- Each time step t :
 - Observe input \mathbf{x}_t
 - Make a prediction \hat{y}_t
 - Observe true output y_t
 - Update coefficients:

for $j=0, \dots, D$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \frac{\partial \ell_t(\mathbf{w})}{\partial w_j}$$

RS

Updating coefficients immediately: Pros and Cons

Pros

- Model always up to date → Often more accurate
- Lower computational cost
- Don't need to store all data, but often do anyway

Cons

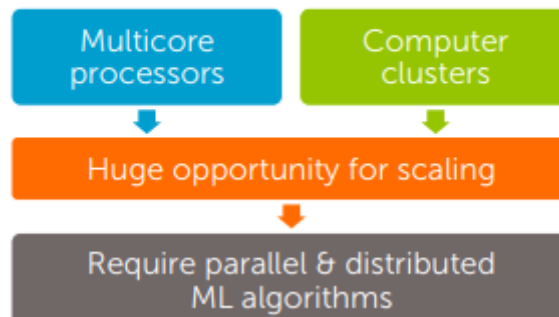
- Overall system is *much* more complex
 - Bad real-world cost in terms of \$\$\$ to build & maintain

Most companies opt for systems that save data and update coefficients every night, or hour, week,...

Scaling huge datasets & Online learning

- Stochastic gradient is a means to scale large huge datasets and online learning.
- Scaling through parallelism

Scaling through parallelism



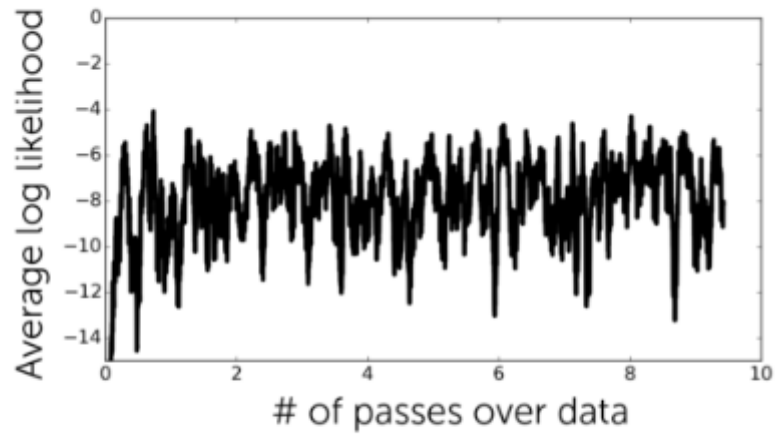
Quiz

1. (True/False) Stochastic gradient ascent often requires fewer passes over the dataset than batch gradient ascent to achieve a similar log likelihood.
 - ☒ True
 - ☐ False
2. (True/False) Choosing a large batch size results in less noisy gradients
 - ☒ True
 - ☐ False

3. (True/False) The set of coefficients obtained at the last iteration represents the best coefficients found so far.

☐ True
☒ False

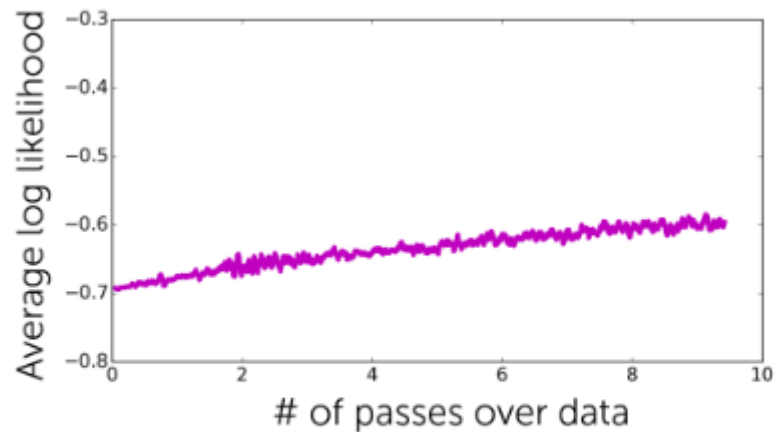
4. Suppose you obtained the plot of log likelihood below after running stochastic gradient ascent.



Which of the following actions would help the most to improve the rate of convergence?

☐ Increase step size
☒ Decrease step size
☐ Decrease batch size

5. Suppose you obtained the plot of log likelihood below after running stochastic gradient ascent.



Which of the following actions would help to improve the rate of convergence?

☐ Increase batch size
☒ Increase step size
☐ Decrease step size

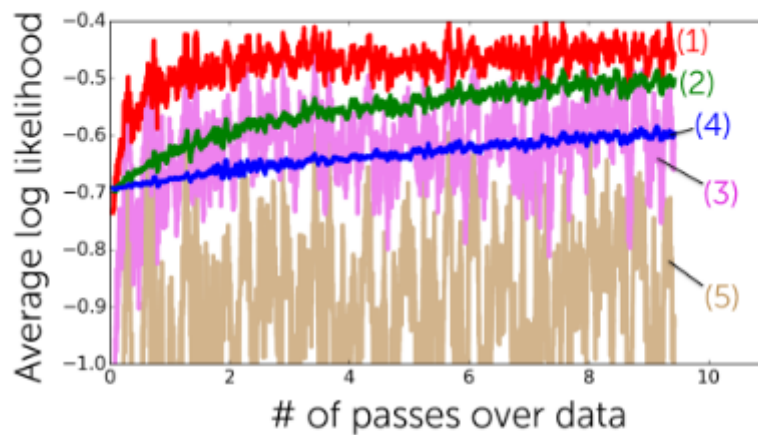
- purple line = step size is too small, takes too long to converge.
- to improve rate of convergence : increase step size.

6. Suppose it takes about 1 milliseconds to compute a gradient for a single example. You run an online advertising company and would like to do online learning via mini-batch stochastic gradient ascent. If you aim to update the coefficients once every 5 minutes, how many examples can you cover in each update? Overhead and other operations take up 2 minutes, so you only have 3 minutes for the coefficient update.

180000

- 3 minutes for coefficient update x 60 seconds x 1 millisecond to compute gradient for a single example = 180,000 gradient examples can be calculated.
- t = time to contribute contribution of x , N data points.
- gradient requires $N \times t$ time to calc 1 step of gradient
- stochastic gradient requires t time to calc 1 step of gradient. (ie: not affected by number of data points)
- ANSWER = 180,000

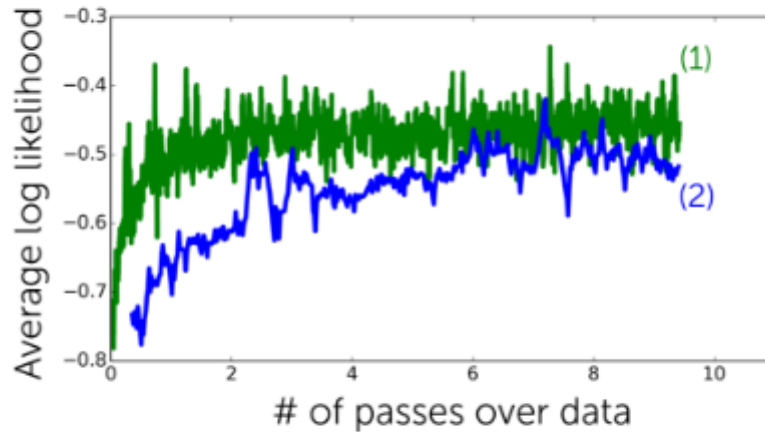
7. In search for an optimal step size, you experiment with multiple step sizes and obtain the following convergence plot.



Which line corresponds to the best step size?

- ☒ (1)
- ☐ (2)
- ☐ (3)
- ☐ (4)
- ☐ (5)

8. Suppose you run stochastic gradient ascent with two different batch sizes. Which of the two lines below corresponds to the smaller batch size (assuming both use the same step size)?



- ☒ (1)
- ☐ (2)

- too small batch size will be noisy and may not converge.

9. Which of the following is NOT a benefit of stochastic gradient ascent over batch gradient ascent? Choose all that apply.

- ☐ Each coefficient step is very fast.
- ☒ Log likelihood of data improves monotonically.
- ☐ Stochastic gradient ascent can be used for online learning.
- ☐ Stochastic gradient ascent can achieve higher likelihood than batch gradient ascent for the same amount of running time.
- ☒ Stochastic gradient ascent is highly robust with respect to parameter choices.

10. Suppose we run the stochastic gradient ascent algorithm described in the lecture with batch size of 100. To make 10 passes over a dataset consisting of 15400 examples, how many iterations does it need to run?

1540

- $15400/100 = 154$ batches.
- $10 \text{ passes} \times 154 \text{ batches} = 1540$

In []: