

$$B = \{n \in \mathbb{R}^d \mid \|n\|_2 \leq 1\}$$

(ii)

$$\text{Let } \Pi_B(n) = \min_{z \in B} \|z - n\|_2$$

if $\|n\|_2 \leq 1$ i.e. $\forall z \in B$

$$\text{then, } \|z\|_2 \leq \|z - n\|_2 + \|n\|_2 \leq \|z - n\|_2 + 1$$

$$\Pi_B(n) = \|z - n\|_2 \geq 0$$

and the minimum is achieved

$$\|z - n\|_2 = 0 \quad \text{when } z = n \Rightarrow \Pi_B(n) = n$$

if $\|n\|_2 \geq 1$

$$\text{Let } \Pi_B(n) = z + \frac{x}{\|n\|_2}$$

using the property,

$$\langle y' - \Pi_D(y'), \Pi_D(y') - y \rangle \geq 0$$

where D is a convex set,

$$y' \notin D, \quad z \notin y \in D$$

~~Let $y' = x$ s.t. $\|x\|_2 > 0$ and $D = B$~~
~~and $\Pi_B(x) = \Pi_B$~~

let $D = B$, $y' = x \in B$

and $\Pi_D(y') = \Pi_B(x) = z \neq \frac{x}{\|x\|_2}$

\Rightarrow If $\Pi_B(x) = z$ is correct projection then

$$\langle x - z, z - y' \rangle \geq 0$$

$$\text{or } \langle x - z, y' - z \rangle \leq 0 \quad \forall y' \in B$$

$$\begin{aligned} \text{LHS} &\rightarrow \langle x, y' \rangle - \langle x - z, y' \rangle - \langle z - y', y' \rangle + \langle z, z \rangle \geq 0 \\ &\approx \end{aligned}$$

$$\text{Let } y = \frac{x}{\|x\|_2}$$

$$\Rightarrow \langle x, y \rangle - \langle x - z, y \rangle - \langle z - y, y \rangle + \|z\|_2^2$$

$$= \frac{\langle x, x \rangle}{\|x\|_2} - \frac{\langle x - z, x \rangle}{\|x\|_2} - \frac{\langle z - y, x \rangle}{\|x\|_2} + \|z\|_2^2$$

$$= \frac{\|x\|_2^2 + \langle x, z \rangle}{\|x\|_2} \left(\frac{1}{\|x\|_2} - 1 \right) + \|z\|_2^2$$

$$= \frac{\|x\|_2^2 + \|z\|_2^2}{\|x\|_2} + \frac{\langle x, z \rangle}{\|x\|_2} \left(\frac{1}{\|x\|_2} - 1 \right) - 0$$

Now given $\|u\|_2 \geq 1$

$$\langle u, z \rangle \left(\frac{1}{\|u\|_2} - 1 \right) \leq \|u\|_2 \|z\|_2 \left(\frac{1}{\|u\|_2} - 1 \right)$$

Substituting in ①

~~LHS ≥ 0~~ $\Rightarrow \text{LHS} \geq \|u\|_2 + \|z\|_2^2 + \|u\|_2 \|z\|_2 \left(\frac{1}{\|u\|_2} - 1 \right)$

$$\geq \|u\|_2 + \|z\|_2^2 + \|z\|_2 - \|u\|_2 \|z\|_2$$

$$\text{given } z \in B \Rightarrow \|z\|_2 \leq 1$$

$$\Rightarrow \text{LHS} \geq \|u\|_2 + \|z\|_2^2 + \|z\|_2 - \|u\|_2$$

$$\text{LHS} \geq \|z\|_2^2 + \|z\|_2 \geq 0$$

\Rightarrow Our assumption that $y, z \neq x$ is not true.

$$\|y\|_2 \geq \|u\|_2$$

for at least $y = \underline{x}$

as it violates, $\langle u, y - z \rangle \leq 0$

$$\langle u - z, y - z \rangle \geq 0$$

or $\langle u - z, z - y \rangle \leq 0$ for at least $y = \underline{x}$

But does $z = \underline{x}$ satisfies this. $\forall y \in B$

$$\langle \underline{x} - z, z - y \rangle = \langle \underline{x} - \frac{\underline{x}}{\|\underline{x}\|_2}, \frac{\underline{x}}{\|\underline{x}\|_2} - y \rangle$$

$$= \langle \underline{x}, \frac{\underline{x}}{\|\underline{x}\|_2} \rangle - \langle \underline{x}, y \rangle - \langle \frac{\underline{x}}{\|\underline{x}\|_2}, y \rangle$$

$$+ \langle \frac{\underline{x}}{\|\underline{x}\|_2}, y \rangle$$

$$= \|\underline{x}\|_2 + \langle \underline{x}, y \rangle \left(\frac{1}{\|\underline{x}\|_2} - 1 \right) - 1$$

$$= (\|\underline{x}\|_2 - 1) + \langle \underline{x}, y \rangle \left(\frac{1}{\|\underline{x}\|_2} - 1 \right)$$

given $\|\underline{x}\|_2 > 1$

$$\text{and } \langle \underline{x}, y \rangle \left(\frac{1}{\|\underline{x}\|_2} - 1 \right) \geq \|\underline{x}\|_2 \|y\|_2 \left(\frac{1}{\|\underline{x}\|_2} - 1 \right)$$

we can get

$$\langle \underline{x} - z, z - y \rangle \geq \|\underline{y}\|_2^2 + (\|\underline{x}\|_2 - 1)^2$$

$$\langle \underline{x} - z, z - y \rangle \geq (\|\underline{x}\|_2 - 1)^2 + \|\underline{y}\|_2^2$$

$$\geq (\|\underline{x}\|_2 - 1) (\|\underline{x}\|_2 - 1 + \|\underline{y}\|_2)$$

given $\|\underline{x}\|_2 > 1$ and $\|\underline{y}\|_2 \leq 1$ as $y \in B$

$$\Rightarrow \langle x - z, z - y \rangle \geq (||x||_2 - 1) \cancel{(||y||_2)} (1 - ||y||_2) \geq 0$$

$\Rightarrow \pi_B(x) = z = \frac{x}{||x||_2}$ is the only

projection which satisfies the given condition.

$$\Rightarrow \pi_B(x) = \begin{cases} x, & \text{if } ||x||_2^2 \leq 1 \\ \frac{x}{||x||_2}, & \text{if } ||x||_2^2 > 1 \end{cases}$$

1)

1.2)

$$S = \{x \in \mathbb{R}^d \mid U^T x = b\}$$

$$\text{s.t. } U^T U = I$$

$$b \in \mathbb{R}^k$$

$$\pi_S(y) = \underset{x \in S}{\operatorname{argmin}} \|x - y\|_2 = \underset{x \in S}{\operatorname{argmin}} \frac{1}{2} \|x - y\|_2^2$$

$$P(UU^T + \nu I)$$

Lagrangian for the constraint minimization,

$$f(x) = \frac{1}{2} \|x - y\|_2^2$$

$$U^T x - b = 0$$

$$L(x, \alpha) = f(x) + \alpha^T (U^T x - b) \quad \{x \in \mathbb{R}^{k+1}\}$$

Using KKT conditions:-

$$0 + \nabla_x P(UU^T + \nu I) \rightarrow p = -U^T(\alpha - y)$$

$$-\text{Stationarity} \quad \nabla_x L(x, \alpha) = \frac{1}{2} U^T (Ux - b) + U^T \alpha = 0$$

$$UU^T x + \nu I - U^T b = U^T \alpha$$

$$\Rightarrow x = -\nu^{-1} U^T \alpha - U^{-1} b \quad (1)$$

$$\Rightarrow x = y - U^T \alpha \quad (2)$$

- Primal feasibility

$$x \in S$$

$$U^T x - b = 0$$

$$\Rightarrow U^T (y - U^T \alpha) - b = 0$$

$$\Rightarrow U^T y - U^T U^T \alpha - b = 0$$

using $U^T U = I$

$$d = U^T y \cancel{\neq} -b \quad (2)$$

Put (2) in (1)

$$x = y \cancel{-} U(-b + U^T y)$$

$$x = y + Ub \cancel{-} UU^T y$$

using ~~$U^T U = I$~~

$$\boxed{x = Ub}$$

$$\Rightarrow \boxed{T_S(y) = Ub}$$

$$\Rightarrow x = \cancel{y} + (I - UU^T)y + Ub$$

$$\Rightarrow \boxed{T_S(y) = (I - UU^T)y + Ub}$$

$$\Delta = \left\{ x \in \mathbb{R}^d \mid \forall i \in [d], x_i \geq 0, \sum_{i \in [d]} x_i = 1 \right\}$$

$$\Pi_\Delta(y) = \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \|x - y\|_2^2 = \underset{x \in \Delta}{\operatorname{argmin}} \frac{1}{2} \|x - y\|_2^2$$

$$f(x) = \frac{1}{2} \|x - y\|_2^2$$

$$\forall i \in [d] -x_i \leq 0$$

$$\sum_{i=1}^d x_i = 1$$

Assume the dimensions of y are arranged in sorted increasing order.

Lagrangian

$$L(x, u, v) = \frac{1}{2} \|x - y\|_2^2 - \sum_{i=1}^d u_i x_i + v \left(\sum_{i=1}^d x_i - 1 \right)$$

Using KKT conditions

~~$\forall i \in [d]$~~

Stationarity:

$$\forall i \in [d] \quad \nabla_{x_i} L(x, u, v) = (x_i - y_i) - u_i + v = 0 \quad \text{--- (1)}$$

complimentary slackness

$$\forall i \in [d] \quad (x_i - u_i)x_i = 0 \quad \text{--- (2)}$$

$$\text{Dual feasibility } \forall i \in [d] \quad u_i \geq 0 \quad \text{--- (3)}$$

Primal feasibility: $\sum_{i=1}^d x_i = 1 \Rightarrow u_i = 1 - x_i \geq 0 \forall i \in [d]$

$$\sum_{i=1}^d x_i = 1 \quad (1) \quad u_i \geq 0 \quad \forall i \in [d] \\ \text{from } (1) \quad x_i = y_i + u_i - v \quad (2)$$

$$x_i = y_i + u_i - v$$

using (2) we get 2 cases,

Case 1 if $x_i > 0 \Rightarrow$ from (2) $y_i + u_i - v = 0$

(Let there be ~~from~~ \exists such i) $\Rightarrow x_i = y_i - v$

Case 2 ~~or~~ $x_i = 0 \Rightarrow u_i > 0$

Now using (1)

$$\sum_{i=1}^d x_i = 1$$

$$\sum_{\substack{i=p+1 \\ i=d-p+1}}^d x_i = \sum_{i=d-p+1}^d y_i - p v = 1$$

$$\Rightarrow v = \frac{1}{p} \left(\sum_{i=d-p+1}^d y_i - 1 \right)$$

$$\Rightarrow u_i = y_i - \frac{1}{p} \left(\sum_{i=d-p+1}^d y_i - 1 \right) \quad i \in \{d-p, \dots, d\}$$

$$\Rightarrow \forall i \in \{1, d-p-1\}$$

$$\Rightarrow \begin{cases} \forall i \in \{1, \dots, d-p\} \quad x_i = 0 \\ \forall i \in \{d-p+1, \dots, d\} \quad x_i = y_i - \frac{1}{p} \left(\sum_{i=d-p}^d y_i - 1 \right) \end{cases}$$

Bonus:- For every

The function $\frac{\sum_{i=d-p}^d y_i - 1}{p}$ is monotonically increasing for wrt. p .

\rightarrow Sort x_i with merge sort $= O(d \log d)$

- For every $p \in \{1, 2, \dots, d\} = O(d)$

$$S_p \text{ calculate } \sum_{i=d-p}^d y_i = \sum_{i=d-(p-1)}^d y_i + y_{d-p}$$

$$S_p = \sum_{i=d-p}^d y_i = S_{p-1} + y_{d-p}$$

- Do binary search on $p \leftarrow O(\log d)$

+ (i) calculate $\sum_{i=d-p}^d y_i = O(d)$ time as $\sum_{i=d-p}^d y_i$ is known

Check if $\sum_{i=1}^d x_i = 1$ and $x_i \geq 0 \forall i = O(d)$
if $x \in \Delta$; return x .

\Rightarrow Total $O(d \log d)$ time complexity.

$$2.1) \quad \mathbb{E}[\nabla f(x_t)] = \nabla f(x) \quad \text{and} \quad \mathbb{E}[||\nabla f(x)||_2^2] \leq G$$

update rule :- $x_{t+1} = x_t - \eta_t \nabla f(x_t)$

Starting from the basic mirror descent lemma for SGD.

$$f(x_t) \leq f(x) + \mathbb{E} \left[\frac{1}{2\eta_t} \left(||x_t - x||_2^2 - ||x_{t+1} - x||_2^2 + (||x_t - x_{t+1}||_2^2) \right) \right]$$

replace η with η_t

$$f(x_t) \leq f(x) + \mathbb{E} \left[\frac{1}{2\eta_t} \left(||x_t - x||_2^2 - ||x_{t+1} - x||_2^2 + (||x_t - x_{t+1}||_2^2) \right) \right]$$

$$\Rightarrow \eta_t f(x_t) \leq \eta_t f(x) + \mathbb{E} \left[\frac{1}{2} \left(||x_t - x||_2^2 - ||x_{t+1} - x||_2^2 + (||x_t - x_{t+1}||_2^2) \right) \right]$$

Taking telescopic sum

$$\sum_{t=0}^{T-1} \eta_t \mathbb{E}[f(x_t)] \leq f(x) \sum_{t=0}^{T-1} \eta_t + \mathbb{E} \left[\frac{1}{2} \left(||x_0 - x||_2^2 + \sum_{t=0}^{T-1} ||x_t - x_{t+1}||_2^2 \right) \right]$$

$\cancel{\sum_{t=0}^{T-1}}$

$$\frac{1}{\sum_{s=0}^{T-1} \eta_s} \sum_{t=0}^{T-1} \mathbb{E}[f(x_t)] \leq f(x) + \mathbb{E}\left[\frac{\|x_0 - x\|_2^2 + \mathbb{E}}{2\eta}\right]$$

$$\left(\frac{1}{\sum_{s=0}^{T-1} \eta_s}\right) \left(\sum_{t=0}^{T-1} \eta_t \mathbb{E}[f(x_t)]\right) \leq f(x) + \|x_0 - x\|_2^2 + \mathbb{E}\left[\sum_{t=0}^{T-1} \|x_t - x_{t+1}\|_2^2\right]$$

using the update rule

$$x_t - x_{t+1} = -\eta_t \tilde{\nabla} f(x_t)$$

$$\Rightarrow \left(\frac{1}{\sum_{s=0}^{T-1} \eta_s}\right) \left(\sum_{t=0}^{T-1} \eta_t \mathbb{E}[f(x_t)]\right) \leq f(x) + \|x_0 - x\|_2^2 + \mathbb{E}\left[\frac{\|\tilde{\nabla} f(x_t)\|_2^2}{2\eta_t}\right]$$

$$\text{given } \mathbb{E}[\|\tilde{\nabla} f(x_t)\|_2^2] \leq G$$

$$\Rightarrow \left(\frac{1}{\sum_{s=0}^{T-1} \eta_s}\right) \left(\sum_{t=0}^{T-1} \eta_t \mathbb{E}[f(x_t)]\right) \leq f(x) + \|x_0 - x\|_2^2 + G + \frac{\sum_{t=0}^{T-1} \eta_t^2}{2 \sum_{t=0}^{T-1} \eta_t}$$

(15)

2.2) Starting from the second last step of Q. (2.1)

$$\text{and setting } \eta_t = \gamma \Rightarrow \sum_{t=0}^{T-1} \gamma_t = \gamma T$$

$$\frac{1}{\gamma T} \left(\mathbb{E} \left[\sum_{t=0}^{T-1} f(x_t) \right] \right) \leq f(x) + \|x_0 - x\|_2^2 + \gamma^2 \mathbb{E} \left[\sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \right]$$

$\leq f(x) + \|x_0 - x\|_2^2 + 2\gamma T$

$$\frac{1}{T} \left(\mathbb{E} \left[\sum_{t=0}^{T-1} f(x_t) \right] \right) \leq f(x) + \left(\|x_0 - x\|_2^2 + 2\gamma \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \right)$$

$2\gamma T$

now f being a convex function,

$$\mathbb{E} \left[\sum_{t=0}^{T-1} f(x_t) \right] \geq \frac{1}{T} \sum_{t=0}^{T-1} f(\bar{x}_T)$$

$$\Rightarrow \mathbb{E} \left[f(\bar{x}_T) \right] \leq f(x) + \|x_0 - x\|_2^2 + 2\gamma \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

- (1)

(1) Using twice quadratic upper bound

$$f(x_{t+1}) \leq f(x_t) + \langle \nabla f(x_t), (x_{t+1} - x_t) \rangle + \frac{L}{2} \|x_{t+1} - x_t\|_2^2$$

using the update rule
 $x_{t+1} = x_t - \gamma \nabla f(x_t)$

$$f(x_{t+1}) \leq f(x_t) - \gamma \langle \nabla f(x_t), \nabla f(x_t) \rangle + \frac{L}{2} \gamma^2 \|\nabla f(x_t)\|_2^2$$

Taking expectation

$$\mathbb{E}[f(x_{t+1})] \leq \mathbb{E}[f(x_t)] - \gamma \mathbb{E}[\nabla f(x_t)]^\top \mathbb{E}[\nabla \tilde{f}(x_t)] + \frac{L}{2} \gamma^2 \mathbb{E}[\|\nabla \tilde{f}(x_t)\|_2^2]$$

given,

$$\mathbb{E}[\nabla \tilde{f}(x_t)] = \nabla f(x_t)$$

$$\text{and } \mathbb{E}[\|\nabla \tilde{f}(x_t)\|_2^2] \leq 2 \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow \mathbb{E}[f(x_{t+1})] \leq \mathbb{E}[f(x_t)] - \gamma \|\nabla f(x_t)\|_2^2 + \frac{L}{2} \gamma^2 \|\nabla f(x_t)\|_2^2$$

for $\gamma \leq \frac{1}{2L}$ $\Rightarrow L \leq \frac{1}{2\gamma}$

$$\mathbb{E}[f(x_{t+1})] \leq \mathbb{E}[f(x_t)] - \frac{L}{2} \|\nabla f(x_t)\|_2^2 - (2)$$

Taking telescopic sum,

$$\sum_{t=0}^{T-1} \mathbb{E}[f(x_{t+1})] \leq \sum_{t=0}^{T-1} \mathbb{E}[f(x_t)] - \frac{\eta}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\mathbb{E}[f(x_T)] - \mathbb{E}[f(x_0)] \leq -\frac{\eta}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq \left(\frac{2}{\eta} \right) (f(x_0) - \mathbb{E}[f(x_T)])$$

now $f(x^*) \leq f(x_T)$ { x^* is the minimizer}
 $\Rightarrow \mathbb{E}[f(x^*)] = f(x^*) \leq \mathbb{E}[f(x_T)]$

$$\Rightarrow \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq \frac{2}{\eta} (f(x_0) - f(x^*)) \quad \text{--- (3)}$$

~~Now using Lipschitz smooth~~

Now using upper quadratic bound on x_0 and x^*

$$f(x_0) \leq f(x^*) + \langle \nabla f(x^*), x_0 - x^* \rangle +$$

$$+ \frac{L}{2} \|x_0 - x^*\|_2^2$$

given x^* is the minimizer

$$\Rightarrow \nabla f(x^*) = 0$$

$$\Rightarrow f(x_0) - f(x^*) \leq \frac{L}{2} \|x_0 - x^*\|_2^2$$

$$\Rightarrow \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq \frac{L}{\eta} \|x_0 - x^*\|_2^2 \quad \text{--- (4)}$$

using ① in ①

$$\Rightarrow \frac{1}{T} \mathbb{E}[f(\bar{x}_T)]$$

$$[\mathbb{E}[f(\bar{x}_T)]] \leq f(x_0) + \frac{\|x_0 - x\|_2^2}{2\gamma T} + \frac{2\gamma}{T} \leq \|x_0 - x^*\|_2^2$$

given $\gamma \leq \frac{1}{2L}$

$$\leq f(x_0) + \frac{\|x_0 - x\|_2^2}{2\gamma T} + \frac{\|x_0 - x^*\|_2^2}{\gamma T}$$

if x_0 is $x \in \mathbb{R}^d$ is farther away from x_0 than x^*

$$\Rightarrow \|x_0 - x\|_2^2 \geq \|x_0 - x^*\|_2^2$$

$$\Rightarrow \boxed{\mathbb{E}[f(\bar{x}_T)] \leq f(x_0) + \frac{3}{2} \frac{\|x_0 - x^*\|_2^2}{\gamma T}}$$

⇒ The convergence

$$\Rightarrow \boxed{\mathbb{E}[f(\bar{x}_T)] \leq f(x_0) + O\left(\frac{\|x_0 - x^*\|_2^2}{\gamma T}\right)}$$

~~Exhibit~~

$$2.3) \quad f(w) = -\frac{1}{N} \sum_{i=1}^N \log (1 + \exp(-y_i \cdot w, x_i))$$

$$\nabla_w f(w) = \frac{1}{N} \sum_{i=1}^N \frac{\exp(-y_i \cdot w, x_i)}{1 + \exp(-y_i \cdot w, x_i)} \nabla_w (-y_i \cdot w, x_i)$$

$$\boxed{\nabla_w f(w) = \frac{1}{N} \sum_{i=1}^N \frac{\exp(-y_i \cdot w, x_i)}{1 + \exp(-y_i \cdot w, x_i)} (-y_i \cdot x_i)}$$

Observations :-

- Lower batch size gives a less smooth graph as the gradients have high variance and noisy for lower batch size.
- For a given batch size as we increase the learning rate the loss drops more quickly compared to lower ones.
- As the batch size increases, the number of operations required to calculate gradients will increase and thus there is an ~~slight~~ increase in training time for higher batch sizes.

sgd_mnist

March 6, 2020

```
[31]: import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt
import os
import argparse
import pdb
import time

[32]: def parse_args():
    parser = argparse.ArgumentParser()

    parser.add_argument('--lr', dest='lr', type=float, default=1e-3, help="learning rate")
    parser.add_argument('--batch_size', dest='momentum', type=float, default=0.9, help="batch_size")

    return parser.parse_args()

[33]: def load_dataset(data_path="."):
    image_size = 28 # width and length of mnist image
    num_labels = 10 # i.e. 0, 1, 2, 3, ..., 9
    image_pixels = image_size * image_size
    train_data = np.loadtxt(os.path.join(data_path, "mnist_train.csv"), delimiter=",")
    test_data = np.loadtxt(os.path.join(data_path, "mnist_train.csv"), delimiter=",")
    return {"train_data":train_data,
            "test_data": test_data,
            }

[12]: def process_data(raw_data, labels_req=[0,1]):
    train_data = raw_data["train_data"]
    test_data = raw_data["test_data"]
```

```

# rescale image from 0-255 to 0-1
fac = 1.0 / 255
train_imgs = np.asarray(train_data[:, 1:])
test_imgs = np.asarray(test_data[:, 1:])
train_labels = np.asarray(train_data[:, :1])
test_labels = np.asarray(test_data[:, :1])

train_imgs = np.divide(train_imgs, np.linalg.norm(train_imgs, axis=1, u
↪keepdims=True))
test_imgs = np.divide(test_imgs, np.linalg.norm(test_imgs, axis=1, u
↪keepdims=True))

train_mask = np.isin(train_labels[:,0],labels_req)
test_mask = np.isin(test_labels[:,0],labels_req)

dataset = { "X_train": train_imgs[train_mask],
            "Y_train": train_labels[train_mask]*2.0 - 1.0,
            "X_test": test_imgs[test_mask],
            "Y_test": test_labels[test_mask]*2.0 - 1.0,
        }

return dataset

```

[34]: raw_data = load_dataset(data_path=".")

[35]: dataset = process_data(raw_data.copy(),labels_req=[0,1])

```

[36]: def plot_props(data_arr,prop_names, figname, xlabel, x_data=None):
    fig = plt.figure(figsize=(16,9))

    for i in range(len(data_arr)):
        if(x_data):
            print(x_data.shape, data_arr.shape)
            plt.plot(x_data[i], data_arr[i], label=prop_names[i])
        else:
            plt.plot(data_arr[i], label=prop_names[i])
    plt.ylabel("train_losses")
    plt.xlabel(xlabel)
    plt.legend()
    plt.title(figname)
    plt.savefig("./{}.pdf".format(figname))
    # plt.show()

```

[37]: def get_loss_grad(W, X, y_true, require_grad=True):
 ...
 W: weight vector (n,)

```

X: input batch (batch_size, n)
'''

dot_prod = np.matmul(X,W)
expo = np.exp(-np.multiply(y_true, dot_prod))
loss = np.mean(np.log(1 + expo))

if require_grad:
    grad = np.divide(expo , (1+expo))
    grad = np.multiply(grad, -1.0*np.multiply(y_true,X))
    grad = np.mean(grad, 0)
    return loss, grad
return loss

```

```
[38]: def test_train_data(W, train_data, train_labels):
    train_loss = get_loss_grad(W, train_data, train_labels, require_grad=False)
    return train_loss
```

```

[39]: def main(lr,batch_size):
    data_path = "./"
    num_iters = 500
    X_train, Y_train = dataset["X_train"], dataset["Y_train"]
    X_test, Y_test = dataset["X_test"], dataset["Y_test"]

    in_dim = X_train.shape[1]
    W = np.zeros(shape=(in_dim,1))

    train_data_loss_arr = []
    train_time_arr = []
    start_time = time.time()
    eval_time = 0.
    loss_calc_time = 0.
    for i in range(num_iters):
        train_data_loss = 0.0
        idxs = np.random.choice(X_train.shape[0], batch_size, replace=True)
        X = X_train[idxs]
        y_true = Y_train[idxs]

        loss, grad = get_loss_grad(W, X, y_true)
        train_time_arr.append(time.time() - start_time - loss_calc_time)

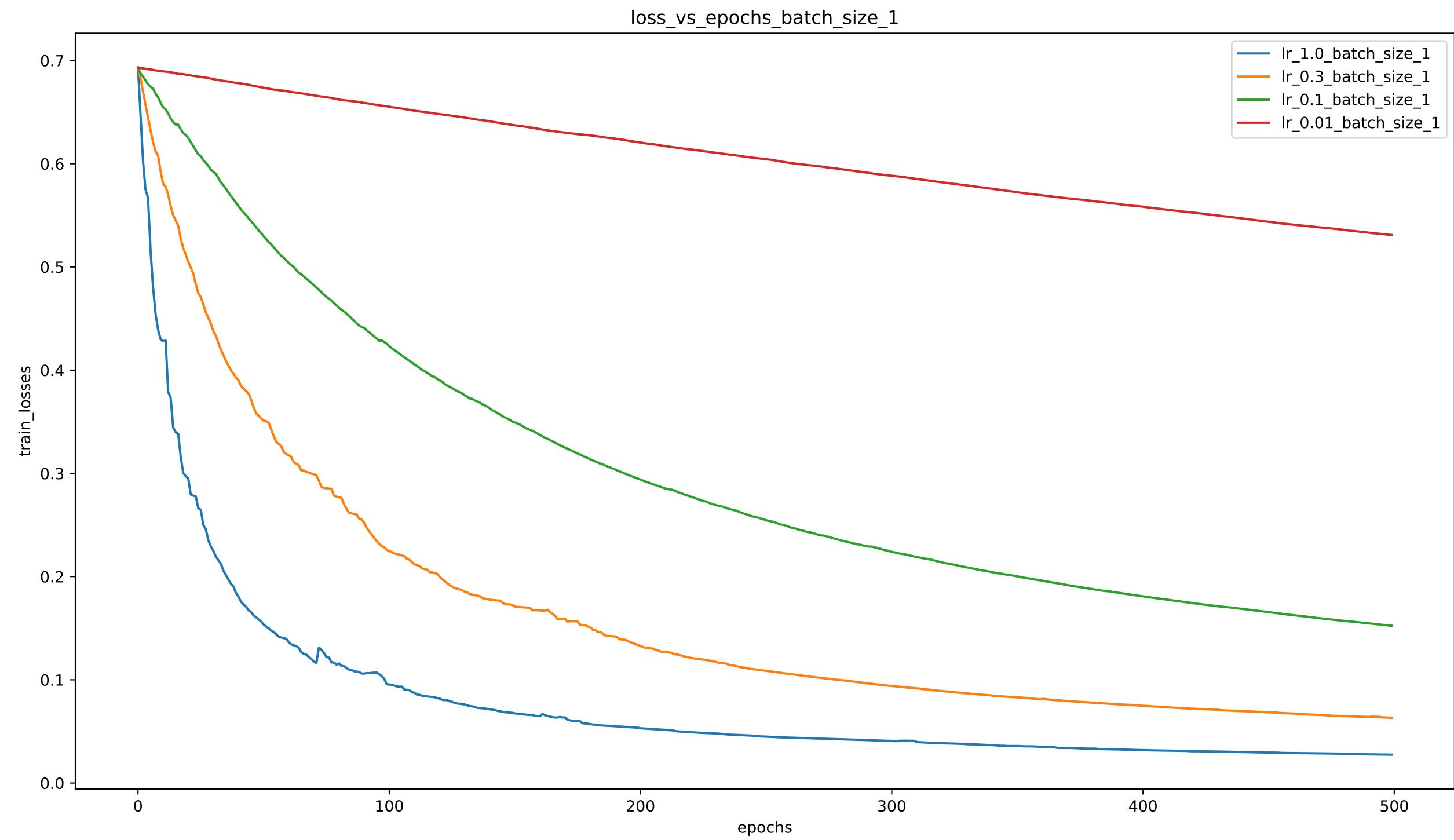
        temp = time.time()
        train_data_loss_arr.append(test_train_data(W, X_train, Y_train)*1.0)
        loss_calc_time += time.time() - temp
        W[:,0] -= lr*grad

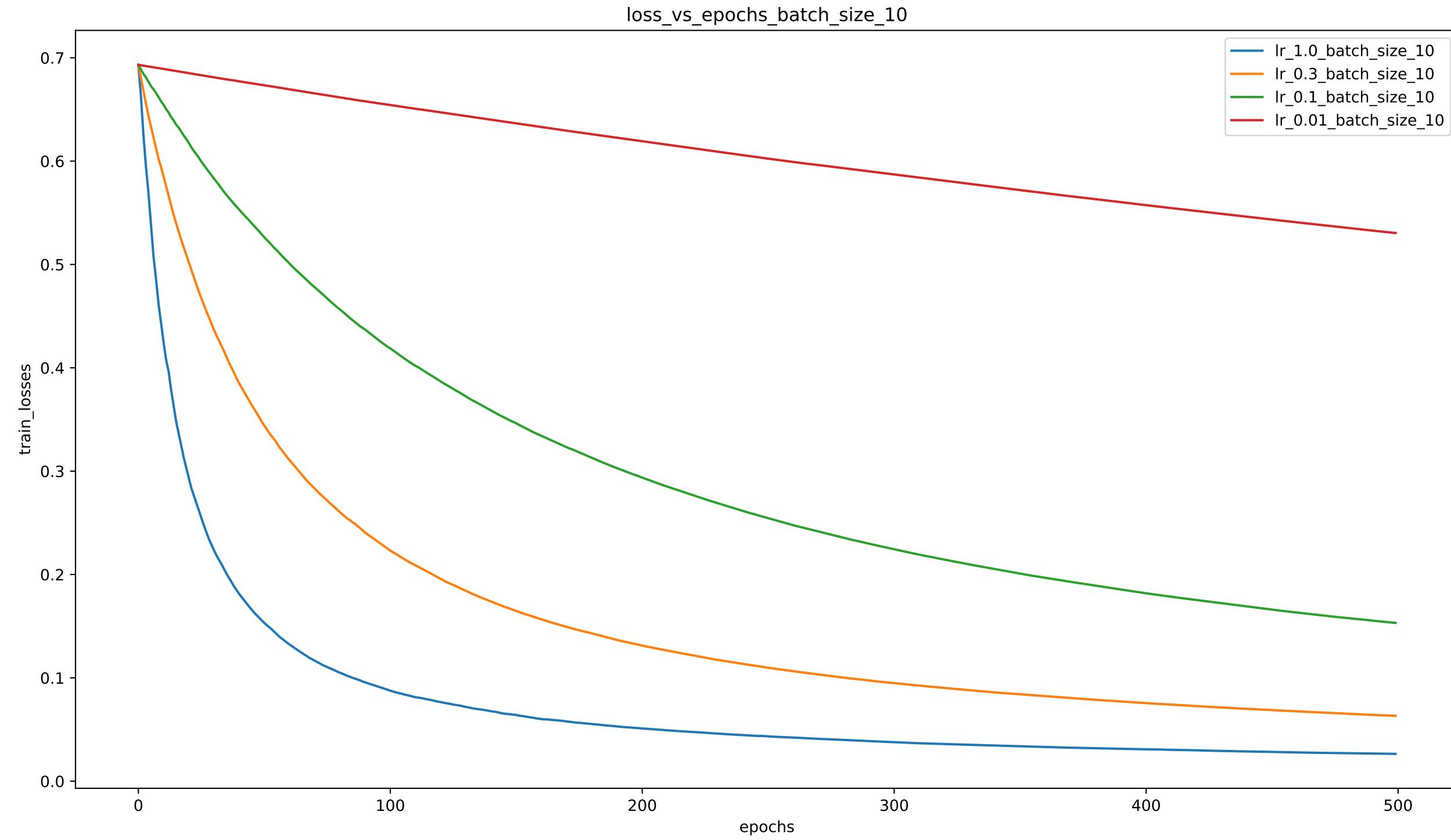
    return np.array(train_data_loss_arr), np.array(train_time_arr)

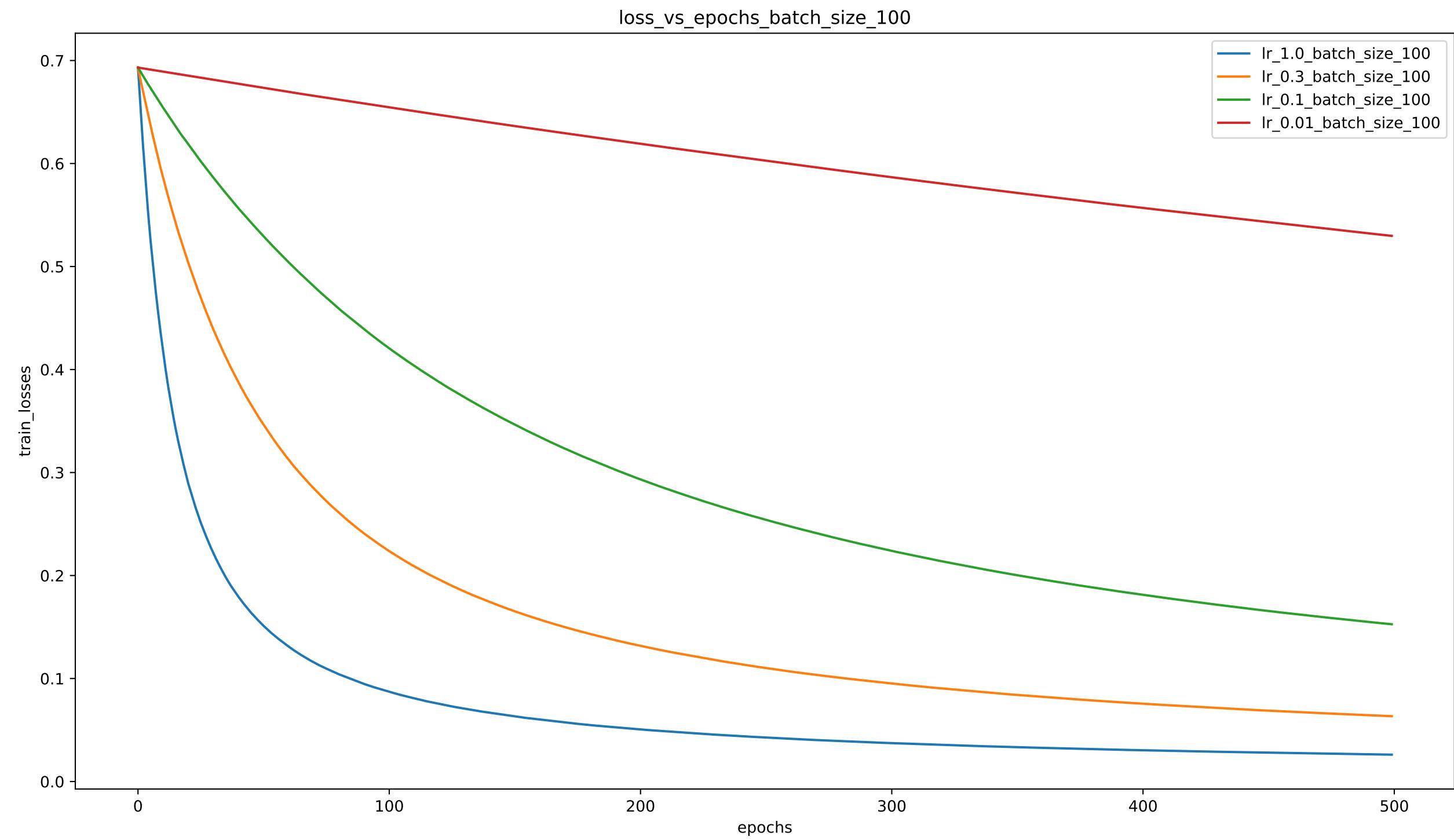
```

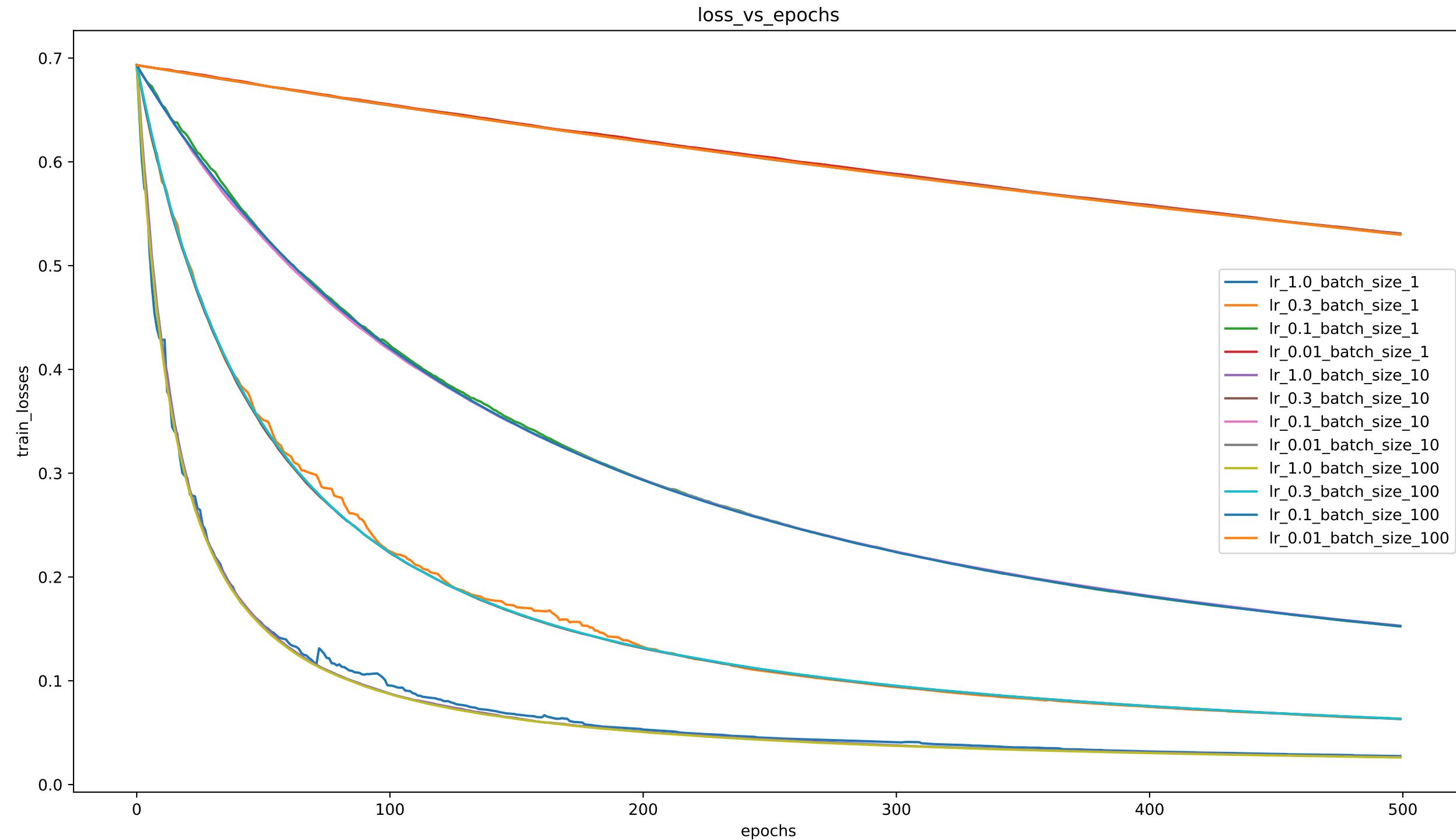
```
#     plot_props(train_data_loss_arr, "train_data_loss_lr_{}_batch_size_{}".
→format(lr, batch_size))
```

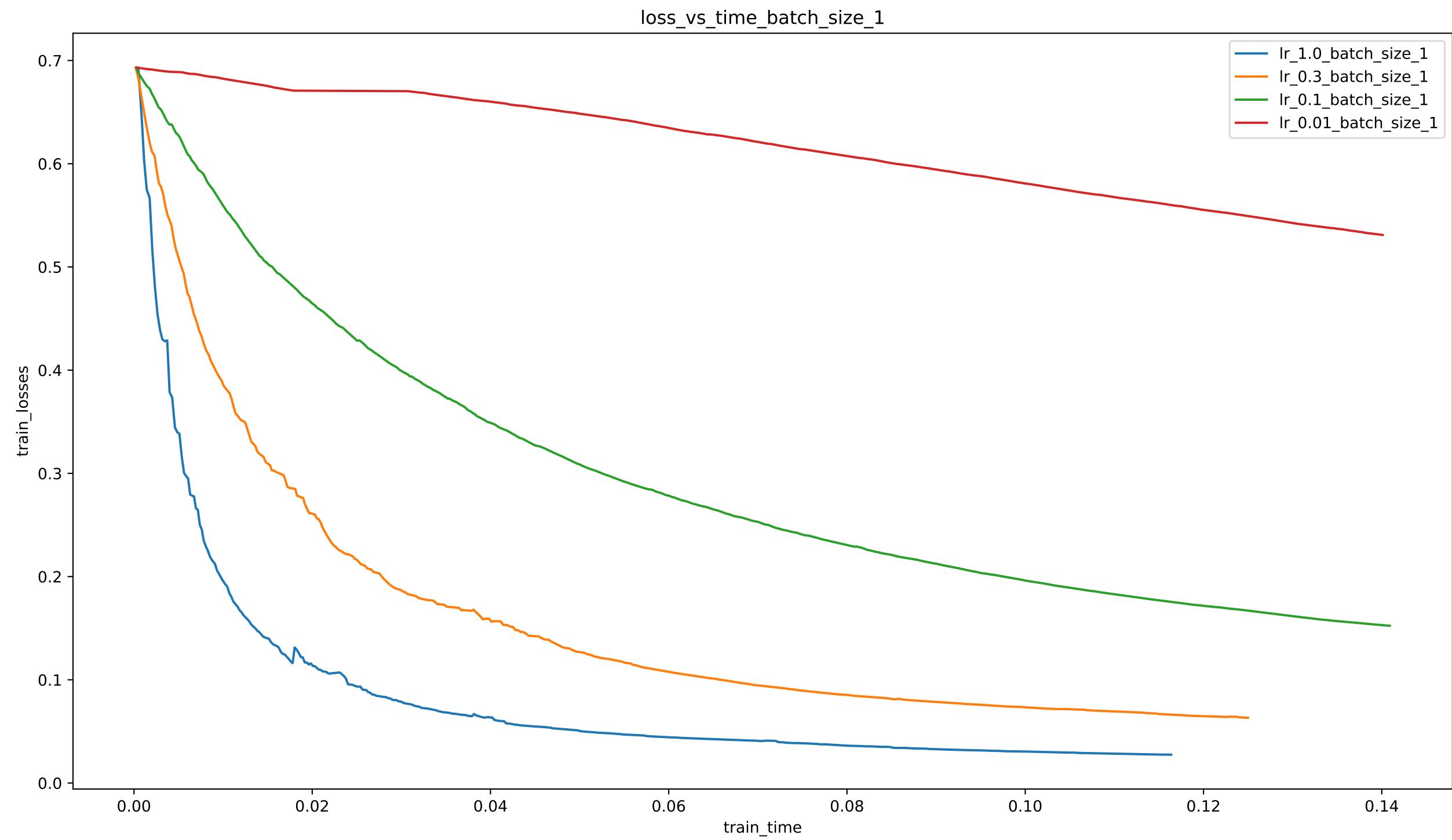
```
[ ]: lrs = [1.0, 0.3, 0.1, 0.01]
batch_sizes = [1, 10, 100]
# lrs = [1.0]
# batch_sizes = [10]
TRAIN_LOSSES = []
TRAIN_TIMES = []
PROP_NAMES = []
for batch_size in batch_sizes:
    print("batch_size", batch_size)
    prop_names = []
    train_losses = []
    train_times = []
    for lr in lrs:
        print("learning_rate", lr)
        start_time = time.time()
        train_loss, train_time = main(lr, batch_size)
        train_losses.append(train_loss * 1.0)
        train_times.append(train_time * 1.0)
    #     print("Time_taken: {:.1f}".format(time.time() - start_time))
    prop_names.append("lr_{}_batch_size_{}".format(lr, batch_size))
    plot_props(train_losses, prop_names, "loss_vs_epochs_batch_size_{}".
→format(batch_size), "epochs")
    plot_props(train_losses, prop_names, "loss_vs_time_batch_size_{}".
→format(batch_size), "train_time", train_times)
    TRAIN_LOSSES.extend(train_losses)
    TRAIN_TIMES.extend(train_times)
    PROP_NAMES.extend(prop_names)
plot_props(TRAIN_LOSSES, PROP_NAMES, "loss_vs_epochs", "epochs")
plot_props(TRAIN_LOSSES, PROP_NAMES, "loss_vs_time", "train_time", TRAIN_TIMES)
```

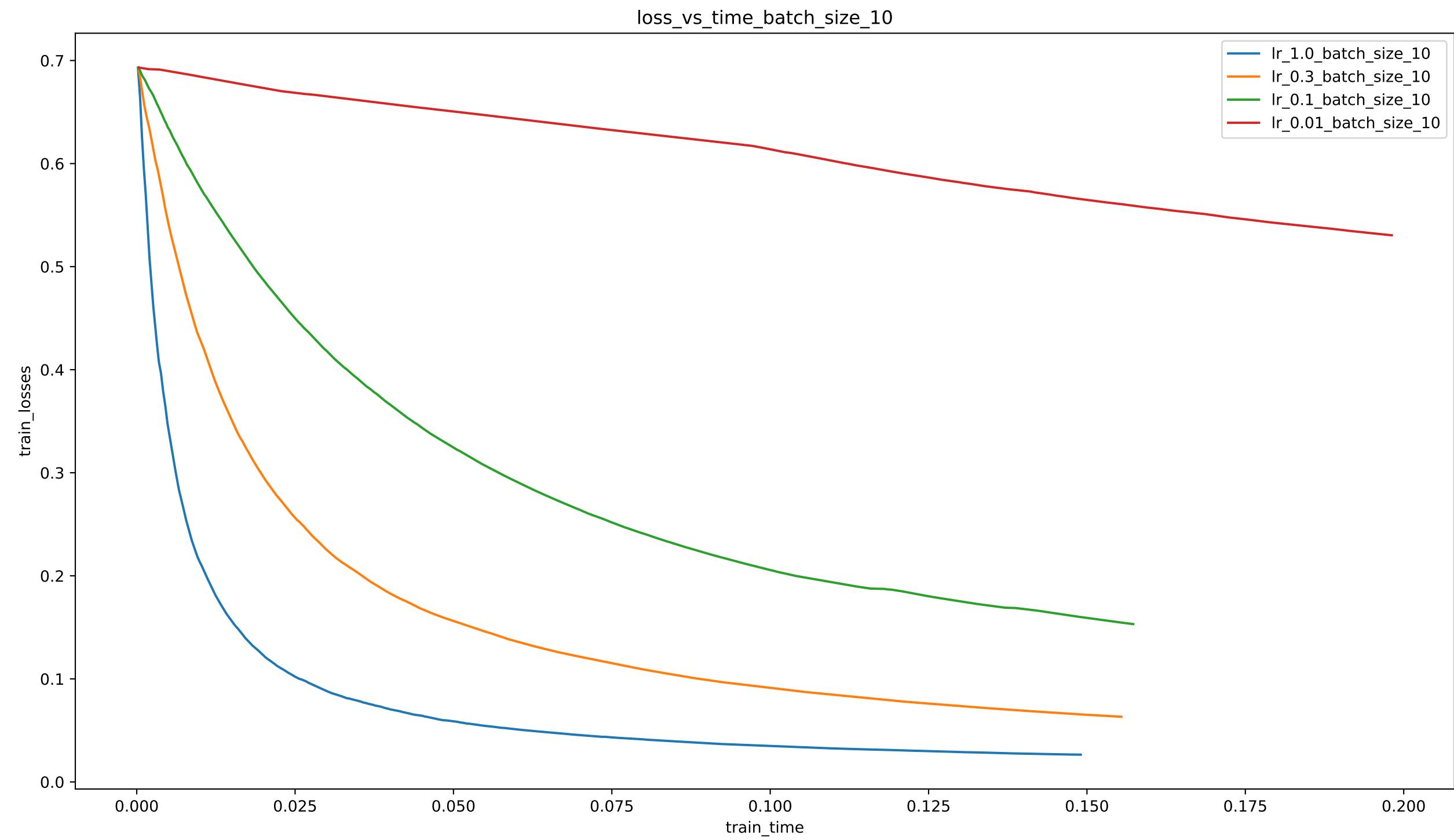


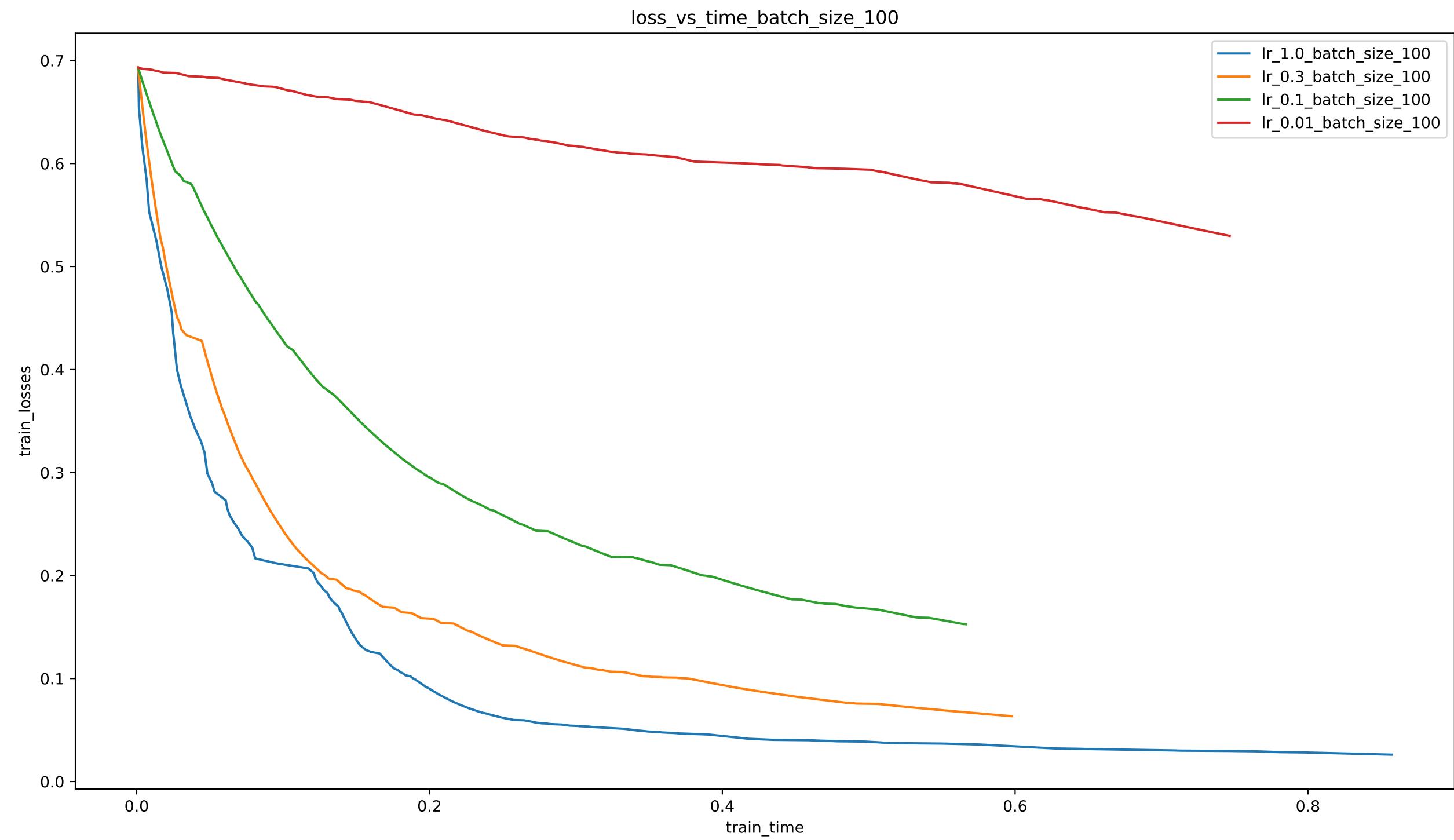


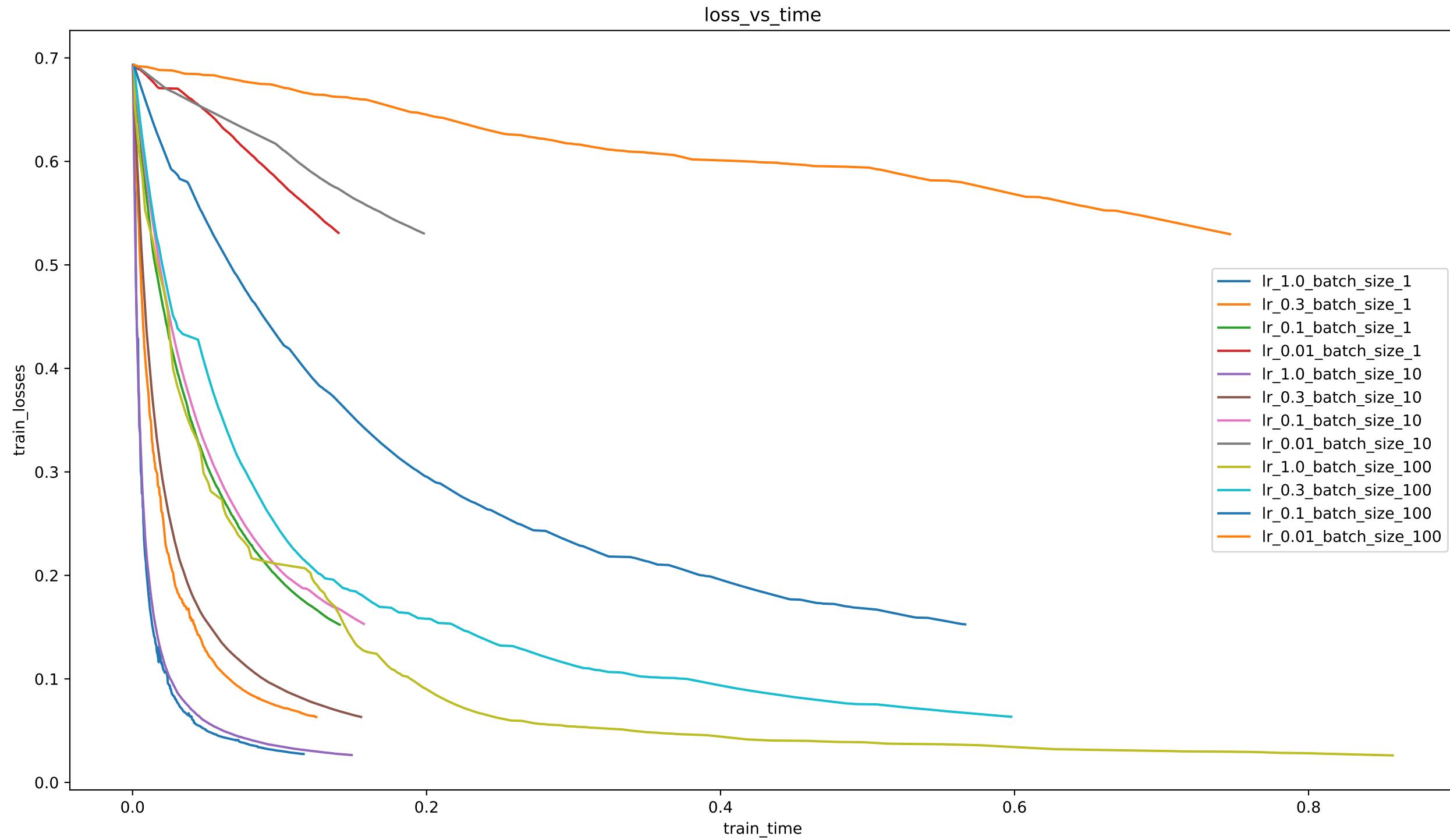












3.1)

$$f(w) = f_1(w) + f_2(w)$$

↓ visit each:

$$f_1(w) = \frac{w_2}{2}(w-1)^2 - w^2$$

$$f_2(w) = 4(w+1)^2 - w^2$$

Global minimum of f

$$\nabla_w f(w) = \nabla f_1(w) + \nabla f_2(w) = 0$$

$$= 4(w-1) - 2w + 8(w+1) - 2w = 0$$

$$(w-1) + 2w + 6w + 8 = 0$$

$$\Rightarrow w = -\frac{1+8}{8} = -\frac{9}{8}$$

Naive ADMM :- $1 = w$

$$w_j^{(t+1)} = \arg \min \left\{ f_j(w) + \lambda \|w - w^{(t)}\|_2^2 \right\}$$

Iteration 1 :-for $j=1$

$$w_1^{(1)} = \arg \min_{w_1} \left\{ \frac{1}{2}(w_1 - 1)^2 + \lambda \|w_1 - w^{(0)}\|_2^2 \right\}$$

$$\nabla_{w_1} f_1$$

$$w_1 = \frac{1}{2}w_1 + \frac{\lambda}{2}$$

$$w_1 = \frac{1}{2}w_1 + \frac{\lambda}{2}$$

$$w_1 = \frac{1}{2}w_1 + \frac{\lambda}{2}$$

iteration 1: $(w^{(0)} = 0), \quad \lambda = 1$ (1.8)

$$w_1^{(1)} = \underset{w}{\operatorname{argmin}} \left(2(w+1)^2 - w^2 + \lambda \|w - w^{(0)}\|_2^2 \right)$$

$$\text{setting } \nabla_w = 2(w+1) - 2w + 2(w-w^{(0)}) = 0$$

$$2w + 2w - 4 = 0$$

$$\Rightarrow w_1^{(1)} = 1$$

$$w = w^{(1)} + \lambda(w - (-1)) \approx$$

$$w_2^{(1)} = \underset{w}{\operatorname{argmin}} \left(2(w+1)^2 - w^2 + \|w - w^{(0)}\|_2^2 \right)$$

$$\text{setting } \nabla_w = 2(w+1) - 2w + 2(w-w^{(0)}) = 0$$

$$\Rightarrow 8w + 8 = 0$$

$$\Rightarrow w = -1$$

$$\Rightarrow w_2^{(1)} = -1$$

$$\Rightarrow w^{(1)} = \frac{1}{2} (w_1^{(1)} + w_2^{(1)}) = \frac{1}{2} (1 + -1) = 0$$

$$\Rightarrow \boxed{w^{(1)} = 0}$$

$$\text{as } w^{(1)} = w^{(0)}$$

$$\Rightarrow w_1^{(2)} = w_1^{(1)} = -1$$

$$\text{and } w_2^{(2)} = w_2^{(1)} = -1$$

$$(\Rightarrow) \quad w_2^{(2)} = \lim_{n \rightarrow \infty} w_2^{(n)} \text{ and } (\star)$$

$$\Rightarrow \exists w^{(i)} \forall i=0 \dots + i \in \mathbb{N}$$

i.e. w converges to $w=0$

\Rightarrow Clearly the naive ADMM method does not converge to the global minimum.

$$f(w) = \left(\|w - g\|_2^2 + \frac{\rho}{2} \|w - b\|_2^2 \right) \sum_{i=1}^m \frac{1}{\lambda_i} \leq f^*$$

$$0 = \rho s + (w - w)^s + w^s - b - w^b = \int_s w$$

$$(D) \vdash b + w^b - w^s \in \mathbb{R} \quad (=)$$

$$0 = \rho s + (w - w)^s + w^s - b + w^b = \int_s w$$

$$(D) \vdash b - \rho w^b = w^s \in \mathbb{R} \quad (=)$$

$$0 = \rho s + (w - w)^s + (w^s - w)^s = \int_s w$$

(Q) Now, Q is true, D is true, so we have

so $b - \rho w^b = w^s$ and $w^s - w = 0$

3.2)

$$\min_{w_j, w} \frac{1}{m} \sum_{j=1}^m (f_j(w_j) + \lambda \|w_j - w\|_2^2)$$

$$\text{s.t. } w_j = w \quad \forall j$$

Dual problem (Maximize the dual function)

$$L(\alpha, w_j, w) = \frac{1}{m} \sum_{j=1}^m \left(f_j(w_j) + \lambda \|w_j - w\|_2^2 \right) + \alpha_j (w_j - w)$$

$$\nabla_{w_1} L = 4w_1 - 4 - 2w_1 + 2(w_1 - w) + \alpha_1 = 0$$

$$\Rightarrow \alpha_1 = 2w - 4w_1 + 4 \quad \textcircled{1}$$

$$\nabla_{w_2} L = 8w_2 + 8 - 2w_2 + 2(w_2 - w) + \alpha_2 = 0$$

$$\Rightarrow \alpha_2 = 2w_2 - 8w_2 - 8 \quad \textcircled{2}$$

$$\nabla_w L = 2(w_1 - w) + 2(w_2 - w) - \alpha_1 - \alpha_2 = 0 \quad \textcircled{3}$$

Put $w_1 = w_2 = w$ in (1), ~~(2)~~, and (3)

$$\Rightarrow \alpha_1 = -2w + 4 \quad \textcircled{4}$$

$$\alpha_2 = -6w - 8 \quad \textcircled{5}$$

$$\alpha_1 = -\alpha_2 \quad \textcircled{6}$$

using ④ and ⑤ in ⑥

$$-2w + 4 = 6w + 8$$

$$\Rightarrow 8w = -4$$

$$\Rightarrow w = -\frac{1}{2}$$

and $\alpha_1 = 5$ and $\alpha_2 = -5$

Clearly the optimal dual solution $w = -\frac{1}{2}$ is

the same as the primal optimal solution obtained in (3.1) when $f(w)$ was directly minimized.

admm

March 6, 2020

```
[1]: import numpy as np  
import matplotlib.pyplot as plt
```

```
[17]: def plot_props(data,prop_name, figname, xlabel):  
    fig = plt.figure(figsize=(7,5))  
    plt.plot(data, label=prop_name)  
    plt.ylabel(prop_name)  
    plt.xlabel(xlabel)  
    plt.legend()  
    plt.savefig("./{}{}.pdf".format(figname))  
#     plt.show()
```

```
[19]: def admm(lambda_, w1_init, w2_init, W_init, alpha_init):  
    # learning rate  
    lr = 2*lambda_  
    # initialize weights  
    w1, w2, W = w1_init, w2_init, W_init  
  
    # initialize alpha  
    alpha = alpha_init  
  
    alpha1_arr = []  
    alpha2_arr = []  
  
    w1_arr = []  
    w2_arr = []  
  
    W_arr = []  
  
    del_W = 1e5  
  
    for i in range(50):  
        # inner minimization  
        w1 = (2*W +4 - alpha[0]) / 4.0  
        w2 = (2*W -8 - alpha[1]) / 8.0
```

```

# update W
del_W = W*1.0
W = (alpha.sum()/(4.0*lambda_)) + 0.5*(w1+w2)

del_W = abs(del_W - W)
#print(del_W)
#update alpha
alpha = alpha - lr*np.array([W-w1, W-w2])
alpha1_arr.append(alpha[0]*1.0)
alpha2_arr.append(alpha[1]*1.0)
w1_arr.append(w1*1.0)
w2_arr.append(w2*1.0)
W_arr.append(W*1.0)

## plotting

plot_props(alpha1_arr, "alpha_1", "alpha_1_vs_iters","iterations")
plot_props(alpha2_arr, "alpha_2", "alpha_2_vs_iters","iterations")
plot_props(w1_arr, "w1", "w1_vs_iters","iterations")
plot_props(w2_arr, "w2", "w2_vs_iters","iterations")
plot_props(W_arr, "W", "W_vs_iters","iterations")

```

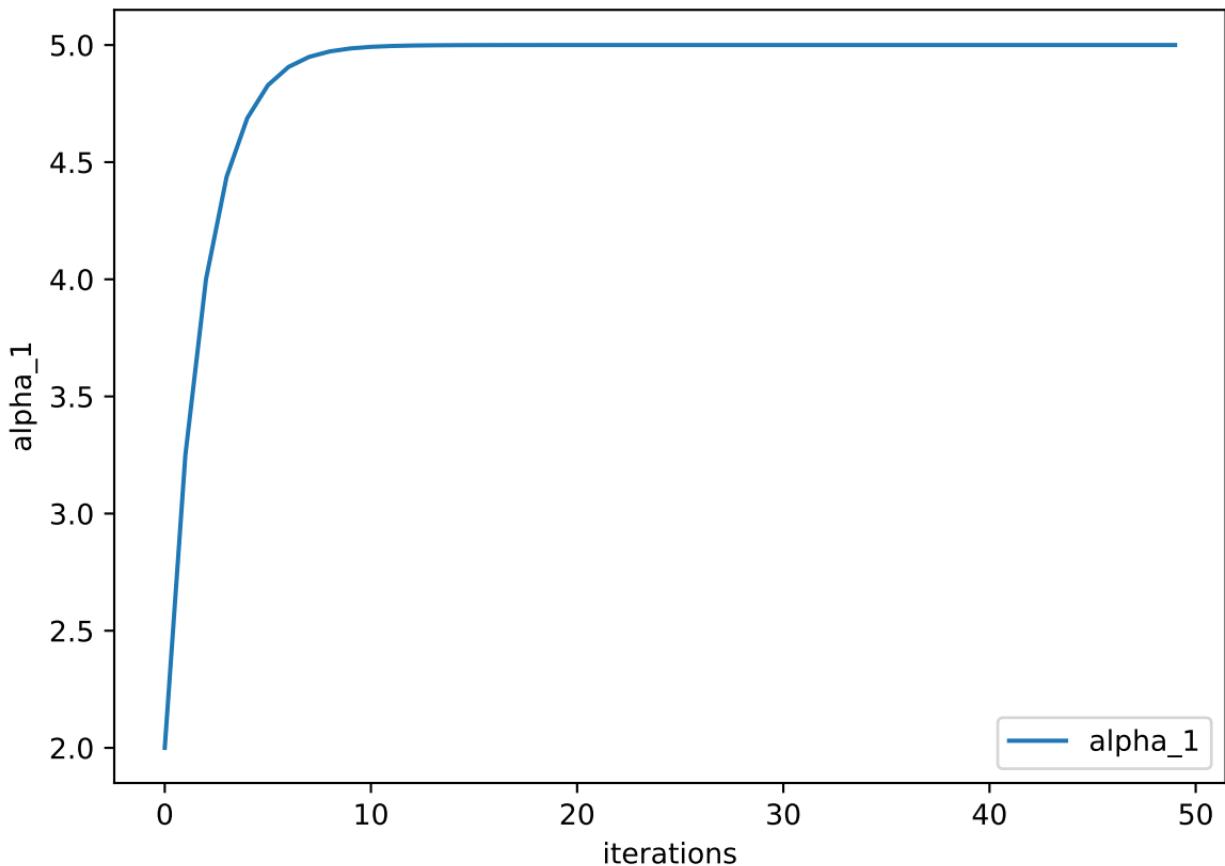
[20]:

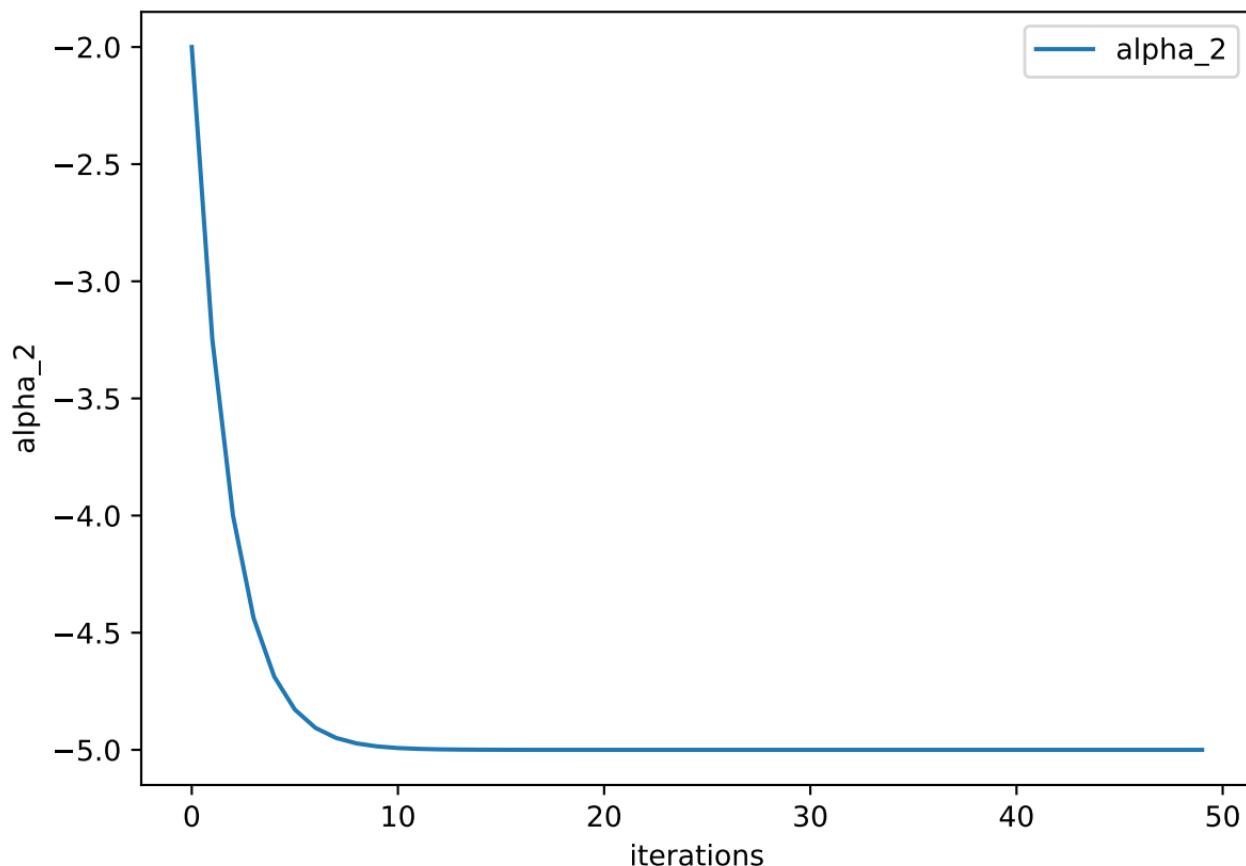
```

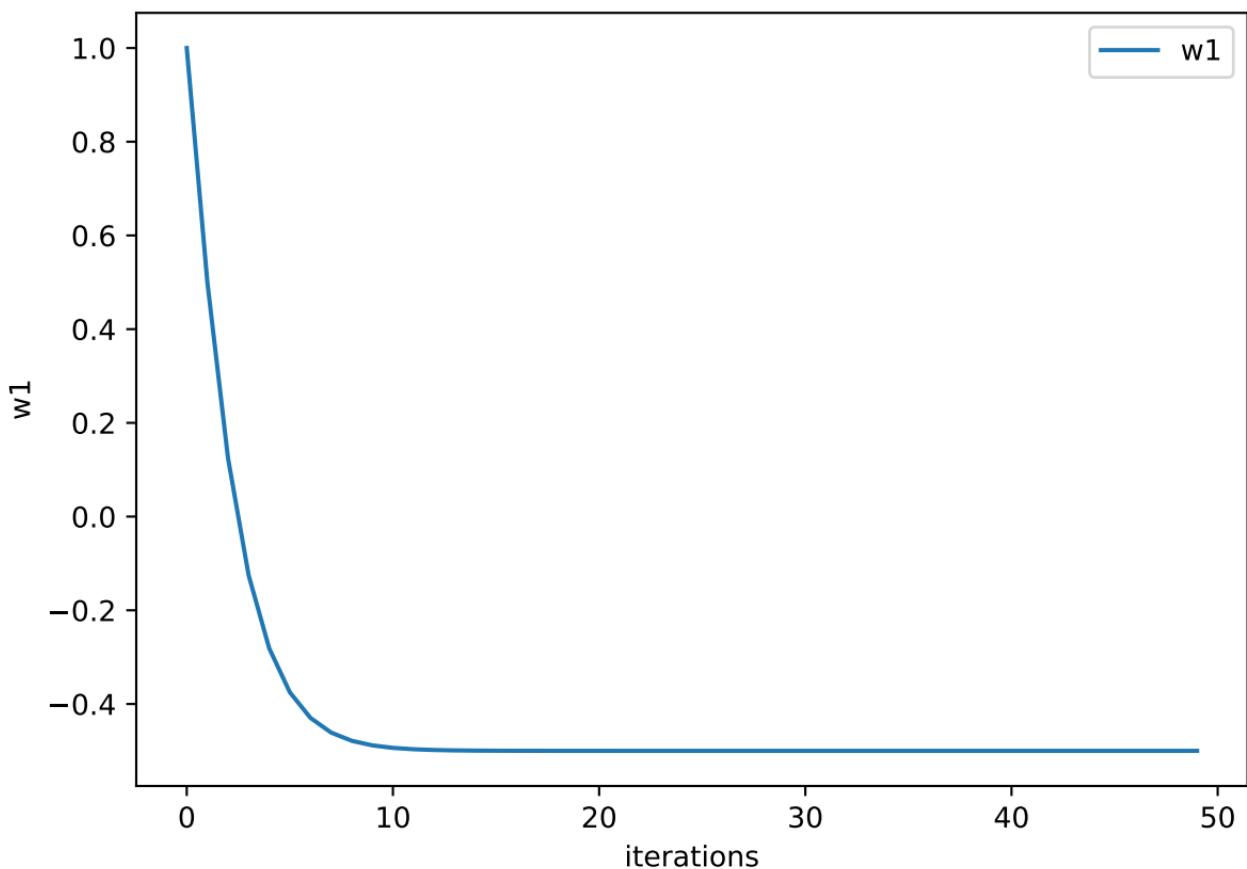
lambda_ = 1.
w1_init = 0.0
w2_init = 0.0
W_init = 0.0
alpha_init = np.zeros(shape=(2,1))

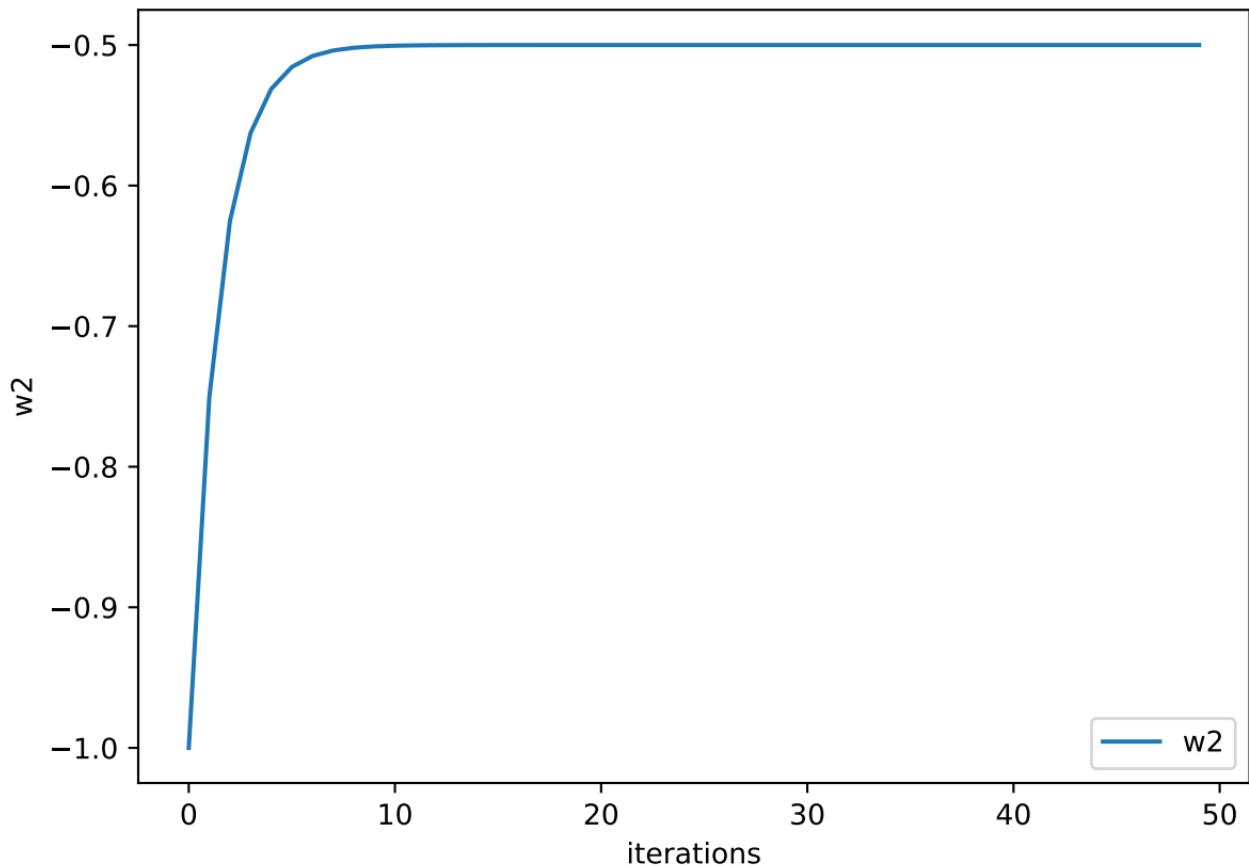
admm(lambda_, w1_init, w2_init, W_init, alpha_init)

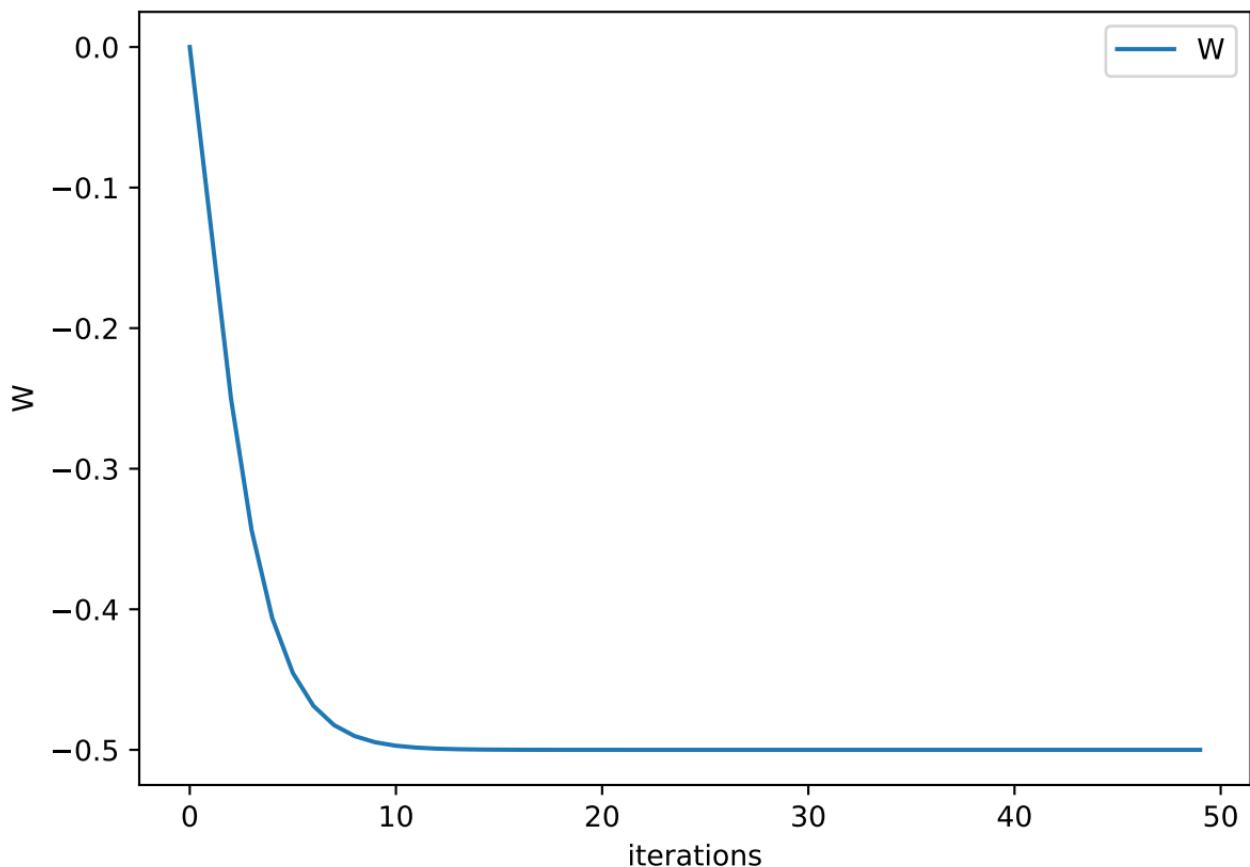
```











$$4.1) \quad h(x) = \lambda|x|, \quad \lambda > 0$$

$$\text{prox}(x) = \arg\min \left\{ h(z) + \frac{1}{2} \|z - x\|_2^2 \right\}$$

$$= \arg\min \left\{ \lambda|z| + \frac{1}{2} \|z - x\|_2^2 \right\}$$

$$\text{let } f(z) = \lambda|z| + \frac{1}{2} \|z - x\|_2^2$$

$$\Rightarrow \nabla_z f(z) \text{ for } z \neq 0 \text{ npf } \dots$$

$$\nabla_z f(z) = \lambda \text{sgn}(z) + (z - x) = 0$$

$$\Rightarrow z = x - \lambda \text{sgn}(z)$$

$$\underline{\text{Case 1}} \quad x > \lambda$$

$$\Rightarrow z = x - \lambda \text{sgn}(z) > \lambda - \lambda \text{sgn}(z)$$

$$\text{if } z > 0 \Rightarrow z > \lambda - \lambda \Rightarrow \boxed{z > 0}$$

$$\text{or } \text{sgn}(z) = \text{sgn}(x)$$

$$\Rightarrow z > \lambda + \lambda \Rightarrow z > 2\lambda > 0$$

but z has to be < 0 (contradiction)

$$\Rightarrow \text{for } x > \lambda, z > 0$$

$$\Rightarrow \text{sgn}(x) \neq \text{sgn}(z)$$

$$\Rightarrow \text{sgn}(z) = \text{sgn}(x)$$

case 2 if $x < -\lambda$ $0 < z \leq \lambda f(z) = (x)^{\alpha}$ (1.1)

$$\Rightarrow z = x - \operatorname{sgn}(z) < -\lambda - \operatorname{sgn}(z)$$

{ if $z > 0$ if $z \leq 0$ minfns = (x)^{\alpha}

$$\Rightarrow z < -\lambda - \lambda \Rightarrow z < -2\lambda < 0$$

contradiction

{ if $z \geq 0$ minfns =

$$\Rightarrow z < -\lambda + \lambda \Rightarrow z < 0$$

{ if $z < 0$ $\Rightarrow |z|^{\alpha} = (z)^{\alpha}$

$$\Rightarrow \text{for } x < -\lambda, z < 0$$

$$\Rightarrow \operatorname{sgn}(z) = \operatorname{sgn}(x)$$

case 3 when $-\lambda \leq x \leq \lambda$ $f(z) = (z)^{\alpha}$

$$(-\lambda - \lambda \operatorname{sgn}(z)) \leq z = x - \lambda \operatorname{sgn}(z) \leq \lambda - \lambda \operatorname{sgn}(z)$$

(s) if $z > 0$ $-2\lambda \leq z \leq 0$

$$\Rightarrow -2\lambda \leq z \leq 0 \quad (\text{contradiction})$$

$0 \leq z \leq 2\lambda$

if $z < 0$
 $\Rightarrow 0 \leq z \leq 2\lambda \quad (\text{contradiction})$

\Rightarrow There is no valid z which can be obtained

when $x \in [-\lambda, \lambda]$ when we set $\nabla_z f(z) = 0$

As both $z < 0$ and $z > 0 \Rightarrow z$ can't be minimizer
of $f(z)$ when $z \in [-\lambda, \lambda]$

\Rightarrow ~~$z = 0$~~ has to be the minimizer
when $z \in [-\lambda, \lambda]$

$$\therefore \text{prox}_h(x) = \begin{cases} x - \lambda \operatorname{sgn}(x), & x > |\lambda| \\ 0 & \text{otherwise} \end{cases}$$

$$4.2) \quad h(x) = \lambda \|x\|_1, \quad x \in \mathbb{R}^d$$

update rule $\rightarrow x_{t+1} = \text{prox}_{\lambda h}(x_t)$

Given From Lecture 10, slide 15,

$$\forall i \in [d], \quad [x_{t+1}]_i = \begin{cases} [x_t]_i - \lambda \text{sgn}([x_t]_i), & |[x_t]_i| > \lambda \\ 0 & \text{otherwise} \end{cases}$$

all until $\Rightarrow \|x_t\|_1 = \|x_{t+1}\|_1$

$$\Rightarrow \forall i \in [d]$$

$$\text{if } |[x_t]_i| > \lambda, \Rightarrow [x_{t+1}]_i = [x_t]_i - \lambda \text{sgn}([x_t]_i)$$

$$\Rightarrow |[x_{t+1}]_i| = |[x_t]_i|$$

and as soon as $|[x_t]_i| \leq \lambda$

$$|[x_{t+1}]_i| = 0$$

so starting from some $|[x_0]_i|$, in m number of steps where $m = \lceil \frac{|[x_0]_i|}{\lambda} \rceil$ and $\lceil \cdot \rceil$ is the ceiling function,

$$|[x_{m+1}]_i| \leq \lambda$$

$$\Rightarrow |[x_{m+1}]_i| = 0$$

every component is zero after $m+1$ steps.

and once $\|x_{t+1}\|_2 = 0$

\Rightarrow It will take $\left\lceil \frac{\max_i |x_{0,i}|}{\eta_d} - 1 \right\rceil$ steps to
take ~~reduce~~ all dimensions to exact 0.

\Rightarrow It will take $O\left(\frac{\max_i |x_{0,i}|}{\eta_d}\right)$ steps.

now, $\max_i |x_{0,i}| = \|x_0\|_\infty \leq \|x_0\|_2$

\Rightarrow It will take $O\left(\frac{\|x_0\|_2}{\eta_d}\right)$ steps
to converge to exact 0.

using gradient descent

update

In this case the gradient is still the same for $|x_t(i)| > \eta_d$ but when $|x_t(i)| \leq \eta_d$

the next $(x_{t+1})_i$ is not set to 0 but rather follows the same update.

In this case the dimension can overshoot the exact minimizer and keeps on oscillating.

e.g:- Let $x \in \mathbb{R}^2$, $x_0 = [2, -4]$

and let $\gamma \lambda = 5/3$

\Rightarrow

i=1

$$[x]_1 = 2 - 5/3 = \frac{1}{3}$$

$$[x]_2 = -4 + 5/3 = -\frac{7}{3}$$

i=2 $[x]_1 = 1/3 - 5/3 = -\frac{4}{3}$

$$[x]_2 = -7/3 + 5/3 = -2/3$$

i=3 $[x]_1 = -4/3 + 5/3 = \frac{1}{3}$

$$[x]_2 = -2/3 + 5/3 = 1$$

i=4 $[x]_1 = 1/3 - 5/3 = -4/3$

$$[x]_2 = 1 - 5/3 = -2/3$$

Clearly both the dimensions have started oscillating around 0 without converging.