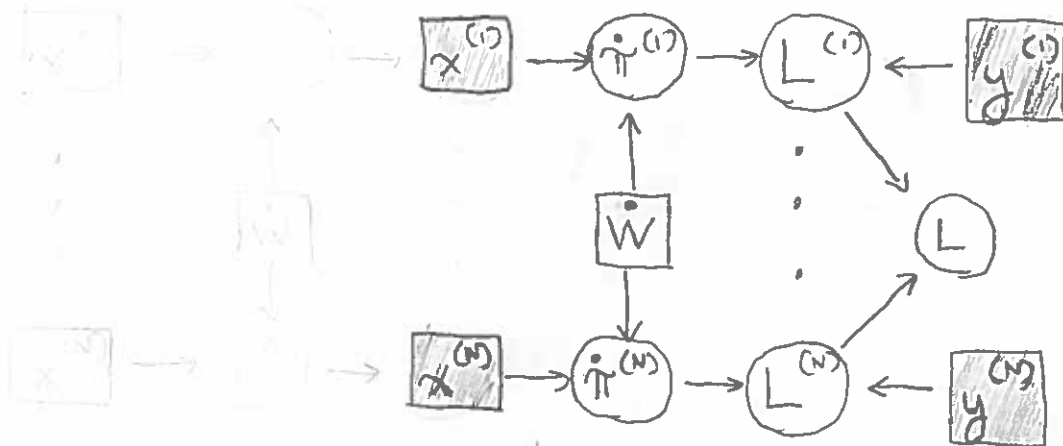


MINIBATCH GRADIENT DESCENT

- ① Consider the full causal diagram for training a 1-layer feedforward neural network:



For simplicity, suppose that \dot{W} is a 1×1 (singleton) matrix with entry w . Thus, the goal of training is to compute $\hat{w} = \operatorname{argmin}_w L$.

- ② Let's suppose we're using the loss function for ordinary linear regression ("ordinary least squares" loss)

$$L = f(w) = \sum_{n=1}^N (y^{(n)} - wx^{(n)})^2$$

MINIBATCH GRADIENT DESCENT

- ③ Now what if you have so many training examples that you can't fit them in memory? (i.e. N is very large). This poses a difficulty for gradient descent:

GRADIENT DESCENT (loss function L , learning rate $\alpha \in \mathbb{R}$):

initialize $w^{(0)}$ to some real number; $t \leftarrow 0$
repeat until happy:
- let update $\sigma^{(t)} \leftarrow -\alpha \cdot \frac{dL}{dw}(w^{(t)}) = -\alpha \cdot \sum_{n=1}^N \frac{dL^{(n)}}{dw}(w^{(t)})$
- let next guess $w^{(t+1)} \leftarrow w^{(t)} + \sigma^{(t)}$
- let $t \leftarrow t+1$

----- this is expensive to compute!

- ④ A common solution is to split the training data into bite-size "minibatches" of a fixed size. Create a partition B_1, B_2, \dots, B_p of the indices $\{1, \dots, N\}$. Then:

MINIBATCH GRADIENT DESCENT (L, α):

initialize $w^{(0)}$ to some real number; $t \leftarrow 0$
repeat until happy:
 for i in $\{1, \dots, p\}$:
 - let update $\sigma^{(t)} \leftarrow -\alpha \cdot \sum_{n \in B_i} \frac{dL^{(n)}}{dw}(w^{(t)})$
 - let next guess $w^{(t+1)} \leftarrow w^{(t)} + \sigma^{(t)}$
 - let $t \leftarrow t+1$

← this is typically called an epoch

MINIBATCH GRADIENT DESCENT

- ⑤ This solves our memory problem, but at what cost?
WHAT DID WE JUST DO? We wanted to compute

$$-\alpha \cdot \frac{dL}{dw}(w^{(t)})$$

but instead we computed:

$$-\alpha \cdot \sum_{n \in B_i} \frac{dL^{(n)}}{dw}(w^{(t)}) \neq -\alpha \cdot \frac{dL}{dw}(w^{(t)})$$

- ⑥ Imagine we had the following data

n	$x^{(n)}$	$y^{(n)}$
1	2	4
2	4	9
3	-1	-2
4	5	15

and we split it into minibatches $B_1 = \{1, 2\}$
 $B_2 = \{3, 4\}$

MINIBATCH GRADIENT DESCENT

- ⑦ When we compute our first update (over minibatch B_1), we are computing the gradient

$$\sum_{n \in B_1} \frac{dL^{(n)}}{dw}(w^{(t)}) = \frac{dL^{(1)}}{dw}(w^{(t)}) + \frac{dL^{(2)}}{dw}(w^{(t)})$$

which is the gradient of loss function

$$L'(w) = L^{(1)}(w) + L^{(2)}(w)$$

- ⑧ When we compute our second update (over minibatch B_2), we are computing the gradient

$$\sum_{n \in B_2} \frac{dL^{(n)}}{dw}(w^{(t)}) = \frac{dL^{(3)}}{dw}(w^{(t)}) + \frac{dL^{(4)}}{dw}(w^{(t)})$$

which is the gradient of loss function

$$L''(w) = L^{(3)}(w) + L^{(4)}(w)$$

- ⑨ Notice that neither of these is our actual loss function:

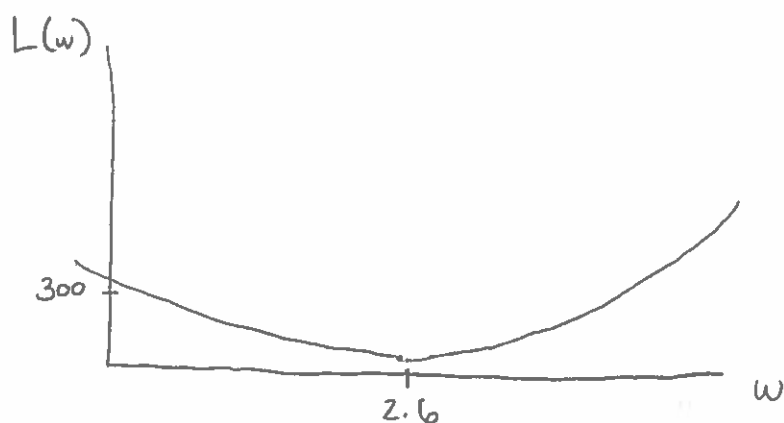
$$L(w) = \sum_{n=1}^4 L^{(n)}(w) = (4-2w)^2 + (9-4w)^2 + (-2+w)^2 + (15-5w)^2$$

$$L'(w) = L^{(1)}(w) + L^{(2)}(w) = (4-2w)^2 + (9-4w)^2$$

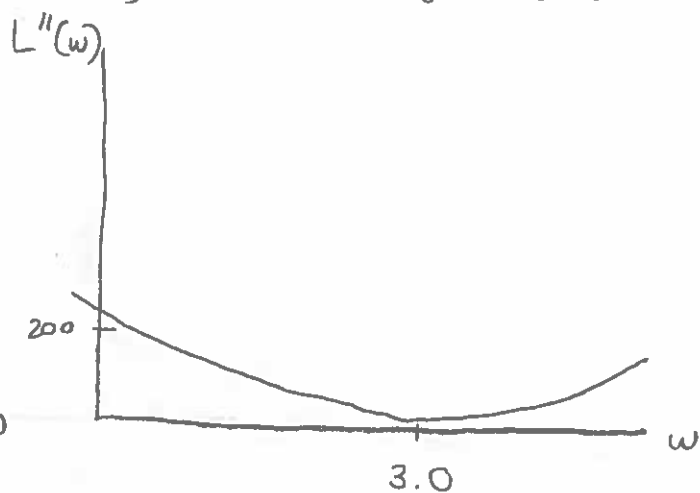
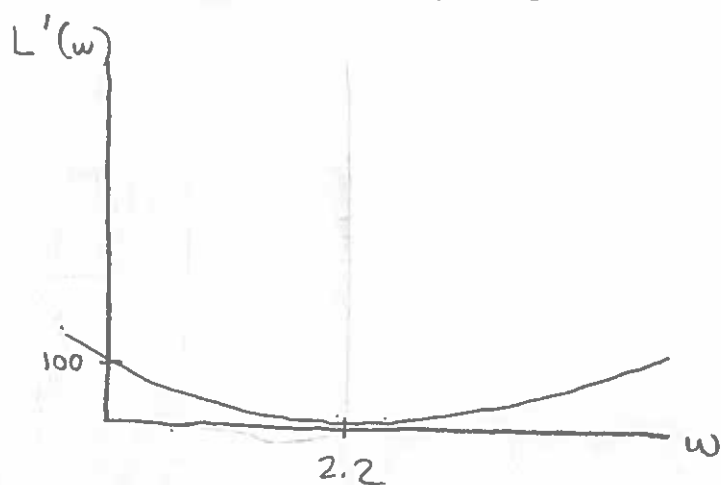
$$L''(w) = L^{(3)}(w) + L^{(4)}(w) = (-2+w)^2 + (15-5w)^2$$

MINIBATCH GRADIENT DESCENT

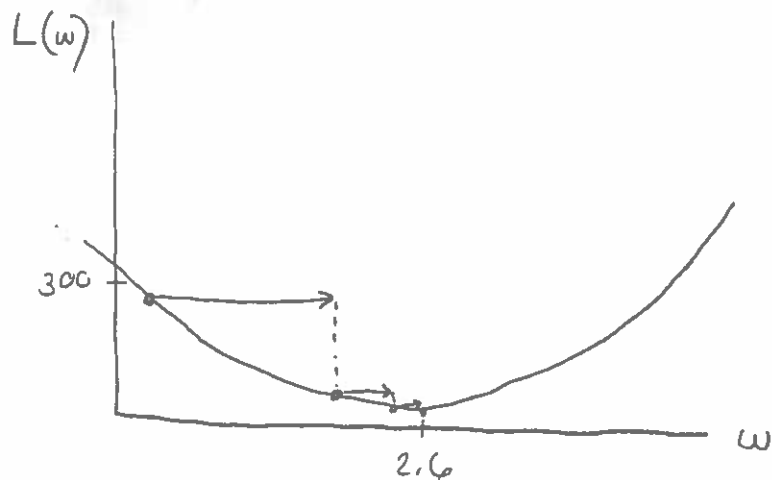
⑩ While $L(w)$ looks as follows:



The other two losses are similar, but not the same:

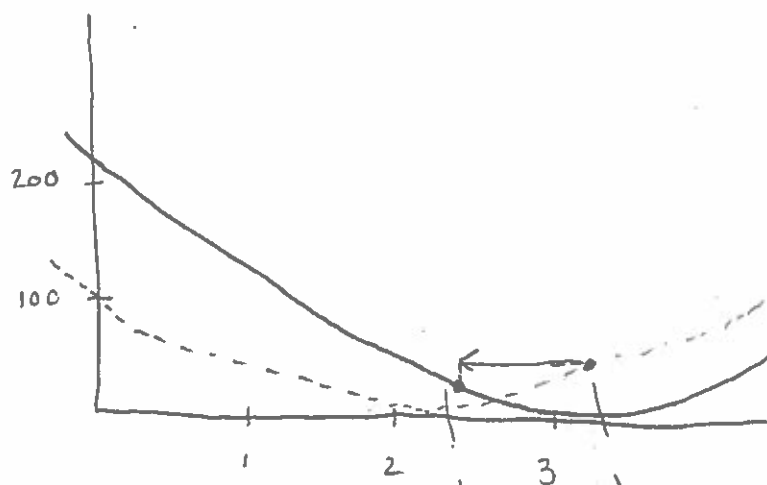


⑪ Standard gradient descent repeatedly finds the gradient, then takes a step in that direction (with the magnitude of the step depending on the steepness of the gradient):



MINIBATCH GRADIENT DESCENT

- ⑫ With minibatch descent, the ground beneath our feet changes with every step!



we start on one curve...

and land on the next!

- ⑬ Why should we expect this to approximate gradient descent on the true loss function?

Well, at the very least we know that if we assume that our data is drawn independently from an identical distribution, then as our batch size increases:

$$\sum_{n \in B_i} L(w) \longrightarrow \sum_{n=1}^N L(w)$$

MINIBATCH GRADIENT DESCENT

- ⑭ So with a sufficiently large batch size, each minibatch loss function should be approximately equal to the true loss function. Thus minibatch gradient descent will behave similarly to gradient descent.
-

- ⑮ What batch size should you use?

Consider the update computation:

$$\sigma^{(t)} \leftarrow -\alpha \cdot \sum_{n \in B_t} \frac{dL^{(n)}}{dw} (w^{(t)})$$

large batch size (points to the summation)
gradient descent (points to the derivative)

This summation can be expressed as a matrix/tensor operation and thus can be better processed with batch sizes between 32 and 256 (too small underutilizes the parallel computation power of GPUs; too large overflows the capacity).

- ⑯ Sometimes people observe that small batch sizes actually help reduce overfitting (i.e. the trained models generalize better because noisy pieces of data do not impact every single update computation).

However, with extremely small batch sizes, the training may be unstable (since the loss function can fluctuate wildly), so you need to worry about the learning rate more.