

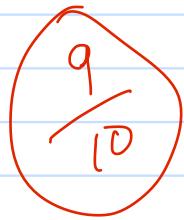
### 1 Backpropagation (2pts)

Consider a MLP consisting of  $N$  input units,  $K$  hidden units and  $N$  output units. The activations are computed as follows:

$$\begin{aligned} z &= \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \\ h &= \sigma(z) \\ y &= \mathbf{x} + \mathbf{W}^{(2)} h + \mathbf{b}^{(2)} \end{aligned}$$

where  $\sigma(\cdot)$  is the sigmoid function, applied element wise. The loss  $\mathcal{L}$  will involve both  $h$  and  $y$ :

$$\begin{aligned} \mathcal{L} &= \mathcal{S} + \mathcal{R} \\ \mathcal{S} &= \frac{1}{2} \|y - s\|_2^2 \\ \mathcal{R} &= r^T h \end{aligned}$$

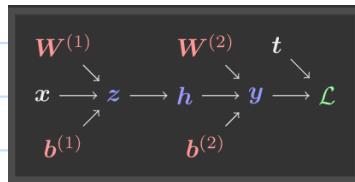


where  $r$  and  $s$  are given vectors.

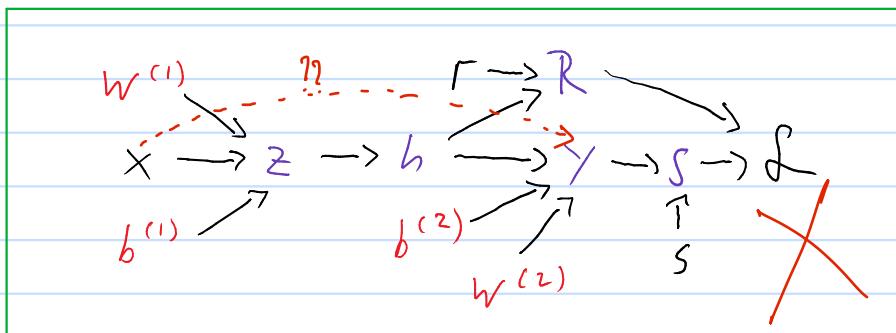
1. (1pt) Draw the computation graph relating  $x, z, h, y, \mathcal{R}, \mathcal{S}$  and  $\mathcal{L}$

2. (1pt) Derive the backpropagation equations for computing  $\frac{\partial \mathcal{L}}{\partial x}$ . You may use  $\sigma'$  to denote the derivation of the sigmoid function, so you do not have to write it out explicitly.

1. From the lecture slides :



So here with different  $\mathcal{L}$ :



2. I'm using the denominator layout convention (as usual in ML), so  $\frac{\partial \mathcal{L}}{\partial v}$  is a column vector, if  $v$  is a column vector and  $\mathcal{L}$  is a scalar.

$$\mathcal{L}' = 1, \quad \mathcal{R}' = \frac{\partial \mathcal{L}}{\partial R} = 1 = S', \quad y' = \mathcal{L}' \left( R \frac{\partial R}{\partial y} + S' \frac{\partial S}{\partial y} \right) = \underbrace{\frac{\partial S}{\partial y}}_{=0} = y - s$$

$$h' = \underbrace{\frac{\partial y}{\partial h} \cdot y'}_{\text{Jacobidian}} + \mathcal{R}' \frac{\partial R}{\partial h} = W^{(2)}(y - s) + r$$

$$z' = h' \odot \frac{\partial h}{\partial z} = (W^{(2)}(y - s) + r) \odot \sigma'(z) \quad \text{is a column vector}$$

elementwise because  $h$  is an element wise function of  $z$

$$x' = \underbrace{\frac{\partial z}{\partial x}}_{\text{Jacobidian}} z' = W^{(1)} z' = W^{(1)} \odot [(W^{(2)}(y - s) + r) \odot \sigma'(z)] = \frac{\partial \mathcal{L}}{\partial x}$$

regular matrix multiplication

You missed one more term.

**2 Vector-Jacobian Products (2pts)**

Consider the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  where  $f(x) = vv^T x$ , and  $v \in \mathbb{R}^{n \times 1}$  and  $x \in \mathbb{R}^{n \times 1}$ . Here, we will explore the relative costs of evaluating Jacobians and vector-Jacobian products. We denote the Jacobian of  $f$  with respect to  $x$  as  $J \in \mathbb{R}^{n \times n}$ .

1. (0.5 pt) Compute  $J$  for  $n = 3$  and  $v^T = [1, 2, 3]$  and write down all its values.
2. (0.5 pt) What is the time and memory cost of evaluating the Jacobian of the function  $f$  as a function of  $n$ ?
3. (1.0 pt) Describe how to evaluate  $J^T y$  where  $y \in \mathbb{R}^n$  with a time and memory cost that is linear in  $n$ . Then, compute  $z = J^T y$  where  $v^T = [1, 2, 3]$  and  $y = [1, 1, 1]^T$ .

it's constant?

1.

$$J = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}, & \frac{\partial f_1}{\partial x_2}, & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1}, & \frac{\partial f_2}{\partial x_2}, & \frac{\partial f_2}{\partial x_3} \\ \frac{\partial f_3}{\partial x_1}, & \frac{\partial f_3}{\partial x_2}, & \frac{\partial f_3}{\partial x_3} \end{pmatrix} = v v^T = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{pmatrix}$$

2. Not quite sure what is asked here.  $f$  is a linear function and therefore its Jacobian is a constant (see above), independent of  $x$ . So evaluating it comes at no computational cost at all, the matrix just needs to be stored in memory.

Or did you mean computing  $J$ ?  $\rightarrow$  Yes computing. should be  $O(n^2)$

Evaluation:  $J(x)$  for  $\frac{\partial J}{\partial x} = 0$  has no cost!

3.  $J^T y = v v^T y$ , so instead of computing  $(v^T v)^T y = v v^T y$  it is more efficient to compute the inner product  $v^T y = J$  and then the scalar product  $J \cdot v$ . Both operations have linear time and memory cost, so the entire evaluation still has linear cost as well.

For  $v^T = (1, 2, 3)$ ,  $y^T = (1, 1, 1)$ :

$$z = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} (1, 2, 3) \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} 6 = \begin{pmatrix} 6 \\ 12 \\ 18 \end{pmatrix}$$

Double check:  $J^T y = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 6 \\ 12 \\ 18 \end{pmatrix}$  ✓

### 3 Linear Regression (3pts)

In this problem we will explore some curious properties of gradient descent and its variants for learning over-parameterized models. We will use linear regression for this purpose but some of these properties also hold for deep neural networks.

Given  $n$  pairs of input data with  $d$  features and scalar label  $(x_i, t_i) \in \mathbb{R}^d \times \mathbb{R}$ , we wish to find a linear model  $f(x) = \hat{w}^T x$  with  $\hat{w} \in \mathbb{R}^d$  that minimizes the squared error of prediction on the training samples defined below. For concise notation, denote the data matrix  $X \in \mathbb{R}^{n \times d}$  and the corresponding label vector  $t \in \mathbb{R}^n$ . The training objective is to minimize the following loss:

$$\min_{\hat{w}} \frac{1}{n} \sum_{i=1}^n (\hat{w}^T x_i - t_i)^2 = \min_{\hat{w}} \frac{1}{n} \|X\hat{w} - t\|_2^2 \quad (1)$$

Assume that  $X$  is full rank:  $X^T X$  is invertible when  $n > d$ , and  $XX^T$  is invertible otherwise. Note that when  $d > n$ , the problem is *underdetermined*, i.e., there are less training samples than parameters to be learned. This is analogous to learning an *overparameterized* model, which is common when training of deep neural networks.

1. **Deriving the Gradient (0.5pt)** Write down the gradient of the loss, w.r.t. the learned parameter vector  $\hat{w}$ .

2. **Underparameterized Model (0.5pt)**  $\Rightarrow$  *overdetermined*

- (a) First consider the underparameterized  $d < n$  case. Write down the solution obtained by gradient descent assuming training converges. Show your work. Is the solution unique?
- (b) Assume that the ground truth labels are generated by a linear target:  $t_i = w^{*T} x_i$ . Show that the solution in part (a) achieves perfect generalization when  $d < n$ , i.e.,  $\forall x \in \mathbb{R}^d$ ,  $(w^{*T} x - \hat{w}^T x)^2 = 0$ .

$$X \in \mathbb{R}^{n \times d}$$

$$\Rightarrow X^T X \in \mathbb{R}^{d \times d},$$

$$X X^T \in \mathbb{R}^{n \times n}$$

1.  $\mathcal{L} = \frac{1}{n} \|Xw - t\|_2^2 = \frac{1}{n} (Xw - t)^T (Xw - t) = \frac{1}{n} (w^T X^T X w - t^T X w + t^T t - w^T X^T t)$   
*constant factor, does not affect minimum*

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial w} = \frac{2}{n} (X^T X w - X^T t)$$

2. Gradient descent:  $w_n = w_{n-1} - \alpha \frac{\partial \mathcal{L}}{\partial w} \Big|_{w=w_{n-1}}$ , with some learning rate  $\alpha$  and initial  $w_0$ .

So for absolute convergence we want:  $w^* := w_n = w_{n \rightarrow \infty}$ , which is the case when we reached the minimum where

$$\frac{\partial \mathcal{L}}{\partial w} \Big|_{w=w^*} = 0 = X^T X w^* - X^T t$$

which leads to the well-known normal equations:

$$X^T X w^* = X^T t \stackrel{(a)}{\Rightarrow} w^* = (X^T X)^{-1} X^T t$$

if  $(X^T X)^{-1}$  exists  $X^T$  is the pseudo inverse

Clearly,  $w^*$  is unique if it exists since it is simply a linear transformation of  $t$ .

(b) for  $t_i = w^{*T} x_i$  we get  $t = X w^*$  and  $w^*$  is the solution:

$$w^* = (X^T X)^{-1} X^T t = (X^T X)^{-1} X^T X w^* = w^* \quad \checkmark$$

$\frac{1}{n} \|Xw^* - t\|_2^2 = \|w^T X - t\|_2^2 = \|w^T X - w^{*T} X\|_2^2 = 0 \quad \forall X$

So clearly:  $(w^{*T} X - \cancel{w^{*T} X})^2 = (w^T X - w^{*T} X)^2 = 0 \quad \forall X$

3. Overparameterized Model: 2D Example (1pt)

- underdetermined
- Now consider the overparameterized  $d > n$  case. We will first illustrate that there exists multiple empirical risk minimizers. For simplicity let  $n = 1$  and  $d = 2$ . Choose  $x_1 = [2, 1]$  and  $t_1 = 2$ , i.e., the one data point and all possible  $\hat{w}$  lie on a 2D plane. Show that there exists infinitely many  $\hat{w}$  satisfying  $\hat{w}^T x_1 = t_1$  on a real line. Write down the equation of the line.
  - We know that multiple empirical risk minimizers exist in overparameterized linear regression and there is only one true solution. Thus, it seems unlikely that gradient descent will generalize if it returns an arbitrary solution. However, we will show that gradient descent tends to find a certain solution with desirable properties. This phenomenon, known as *implicit regularization*, helps explain the success of gradient-based methods to train overparameterized models, like deep neural networks. First consider the 2-dimensional example in the previous part: starting from zero initialization i.e.,  $\hat{w}(0) = 0$ , what is the direction of the gradient? You should write down a unit-norm vector. Does the direction change along the trajectory? Based on this geometric intuition, which solution - along the line of solutions - does gradient descent find? Provide a pictorial sketch or a short description of your reasoning.
  - Give a geometric argument that among all the solutions on the line, the gradient descent solution from the previous part has the smallest Euclidean norm. Hint: Use the Pythagorean Theorem.

3. (a) The equation of the line is  $w^T x_1 = t_1 = w_x x_{1x} + w_y x_{1y}$

$$\text{Or rearranged: } w_y = \frac{t_1}{x_{1y}} - \frac{x_{1x}}{x_{1y}} w_x = 2 - 2 w_x$$

for  $x_1^T = (2, 1)$ ,  $t_1 = 2$

So one of the components of  $w$  is a free parameter.

Choosing any arbitrary value for either  $w_x$  or  $w_y$  determines the value of the other to satisfy  $w^T x_1 = t_1$ .

any  $w = \frac{t_1 x_1}{\|x_1\|_2^2} + \lambda \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} x_1 = \frac{2}{5} \begin{pmatrix} 2 \\ 1 \end{pmatrix} + \lambda \begin{pmatrix} 1 \\ -2 \end{pmatrix}, \lambda \in \mathbb{R} \text{ is a solution!}$

(b) Here:  $X = x_1^T$ , since  $n=1$ . So the gradient for  $w_0 = 0$  is:

$$\frac{\partial f}{\partial w} \Big|_{w=0} = \frac{1}{n} (2 X^T x - 2 X^T t) = -2 x_1^T t$$

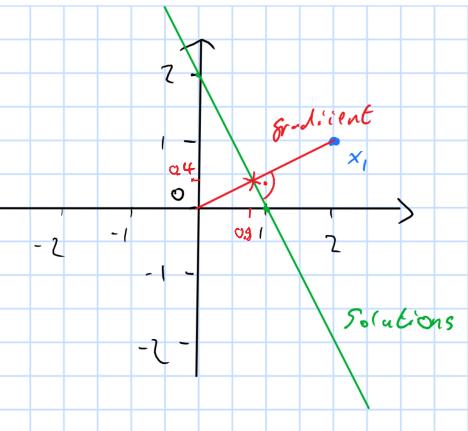
$$\text{So the direction of } x_1. \text{ Normalized: } \frac{x_1}{\|x_1\|_2} = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

The direction of the gradient does not change along the trajectory:

$$\lambda \in \mathbb{R}: \frac{\partial f}{\partial w} \Big|_{w=\lambda x_1} = 2 x_1^T \lambda x - 2 x_1^T t = 2 \lambda x_1^T x_1 - 2 x_1^T t$$

$= 2(\lambda \|x_1\|_2^2 - 1)x_1$ , still points in  $x_1$  direction!

some other constant



The solution obtained by gradient descent will be  $\frac{t_1}{\|x_1\|_2^2} x_1 = \frac{2}{5} x_1 = (0.8, 0.4)^T$ . The closest point on the solution line to the starting point  $w_0 = 0$ .

(c) Gradient descent will always find the solution closest to the starting position  $w_0$ . Which is the intersection

of  $\frac{b_i x_i}{\|x_i\|_2^2} + \lambda \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} x_i$  with  $w_0 + \beta x_i$ ,  $\lambda, \beta \in \mathbb{R}$   
 $\in \text{range}(x_i)$   $\in \text{null}(x_i)$

$$\text{so: } w^* = \frac{b_i x_i}{\|x_i\|_2^2} + \lambda^* \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} x_i = w_0 + \beta^* x_i$$

It is easy to see, that to minimize the Euclidean norm  $\|w^*\|_2^2$ ,  $\lambda^* = 0$ .

So any  $w_0 = \frac{b_i x_i}{\|x_i\|_2^2} - \beta^* x_i$  will find

the solution with minimal Euclidean norm:  $\tilde{w} = \frac{b_i x_i}{\|x_i\|_2^2} = \min_w \|w^*\|_2^2$

Just like in part (b) with  $w_0 = 0$  ( $\beta^* = \frac{b_i}{\|x_i\|_2^2}$ ).

Simple argument in words: Starting at  $w_0$  SD will always find the solution closest to  $w$  because the gradient always points from  $w$  to the hyperplane of solutions (perpendicular to the hyperplane), so  $w_0 = 0$  will always lead to the solution with minimal Euclidean norm.

In fact, any  $w_0 = \lambda x_i$  will converge to that solution (see part 4).

4. Overparameterized Model: General Case (1pt)

- (a) Now we will generalize the previous geometric insight developed to general  $d > n$ . Show that gradient descent from zero initialization i.e.,  $\hat{w}(0) = 0$  finds a unique minimizer if it converges. Write down the solution and show your work. Hint: Recall the result on the direction of the gradient in the previous 2D example, which suggests that the gradient vector is always spanned by the rows of  $X$ . What does this tell you about the gradient descent solution?
- (b) Given the gradient descent solution from the previous part  $\hat{w}$  and another zero-loss solution  $\hat{w}_1$ , evaluate  $(\hat{w} - \hat{w}_1)^T \hat{w}$ . Use this quantity to show that among all the empirical risk minimizers for  $d > n$ , the gradient descent solution has the smallest Euclidean norm.

In the case of  $d > n$  and  $X$  of full rank,  $Xw = t$  is underdetermined and has infinitely many solutions.

Any solution must have the form:  $w = \{\vec{w} + w_{\perp} \mid w_{\perp} \in \text{Null}(X^T)\}$  where  $\vec{w} = X^T(X^T X)^{-1}t$  is the only solution in  $\text{range}(X^T)$  because  $XX^T$  has full rank!

$$(\text{Easy to verify: } X\vec{w} = \underbrace{X X^T(X^T X)^{-1}}_1 t = t)$$

It is easy to see this is the least norm solution, because any other solution  $w$  must satisfy:

$$Xw = X(\vec{w} + (w - \vec{w})) = X\vec{w} + X(w - \vec{w}) = t \Rightarrow X(w - \vec{w}) = 0$$

$$\text{So } w - \vec{w} \in \text{null}(X) \text{ and } \|w\|_2 = \|\vec{w} + (w - \vec{w})\|_2 \geq \|\vec{w}\|_2$$

because  $\vec{w} \perp (w - \vec{w})$

In words:  $\vec{w} = X^T(X^T X)^{-1}t$  is the only solution of  $Xw = t$  in the range of  $X^T$ . Adding any vector from  $\text{null}(X)$  to this unique solution also gives a solution, and increases the norm.  
 ( $\text{range}(X^T)$  is the span of the rows of  $X$ )

(a) Now, to show that GD converges to  $\vec{w}$  for  $w_0 = 0$ :

$$\text{Again: } \mathcal{L} = \frac{1}{n} \|Xw - t\|_2^2 \stackrel{1.}{\Rightarrow} \frac{\partial \mathcal{L}}{\partial w} = \frac{2}{n} (X^T X w - X^T t) \\ = \frac{2}{n} X^T (Xw - t)$$

$\text{range}(X^T) \hookleftarrow \text{some vector}$

$$\text{and also: } X^T X w \in \text{range}(X^T) \nparallel w. \quad \boxed{\text{So } \frac{\partial \mathcal{L}}{\partial w} \in \text{range}(X^T).}$$

So for any  $w_0 \in \text{range}(X^T)$ , any gradient descent update will stay in  $\text{range}(X^T)$  and the only solution in  $\text{range}(X^T)$  is  $\vec{w}$  as shown above. Therefore, if GD converges, it converges to  $\vec{w}$  when starting in  $\text{range}(X^T)$ , for example by choosing  $w_0 = 0$ .

(b) I already showed this. But again:

$$X \vec{w}_1 = X(\vec{w} + (\vec{w}_1 - \vec{w})) = X\vec{w} + X(\vec{w}_1 - \vec{w}) = t \Rightarrow X(\vec{w} - \vec{w}) = 0$$

which means that for any solution  $\vec{w}_1 \neq \vec{w}$ ,  $(\vec{w} - \vec{w}_1)$  is orthogonal to  $\vec{w}$ :

$$(\vec{w} - \vec{w}_1)^T \vec{w} = 0 \Rightarrow \vec{w}^T \vec{w} - \vec{w}_1^T \vec{w} = \vec{w}_1^T \vec{w} = \vec{w}^T \vec{w}$$

$$\begin{aligned} \text{So: } (\vec{w} - \vec{w}_1)^T (\vec{w} - \vec{w}_1) &= \vec{w}^T \vec{w} - 2\vec{w}_1^T \vec{w} + \vec{w}_1^T \vec{w}_1 \\ (\simeq \|\vec{w} - \vec{w}_1\|_2^2) &= \vec{w}_1^T \vec{w}_1 - \vec{w}^T \vec{w} = \|\vec{w}_1\|_2^2 - \|\vec{w}\|_2^2 \\ \Rightarrow \|\vec{w}_1\|_2^2 &= \|\vec{w}\|_2^2 + \|\vec{w} - \vec{w}_1\|_2^2 \geq \|\vec{w}\|_2^2 \end{aligned}$$

Equality only for  $\|\vec{w} - \vec{w}_1\|_2 = 0$  so  $\vec{w} = \vec{w}_1$ .

which means  $\vec{w}$  is indeed the least norm solution.

□

5. Benefit of Overparameterization: Visualize and compare underparameterized with overparameterized polynomial regression (see attached jupyter notebook). Include your code snippets for the `fit_poly` function in the write-up. Does overparameterization (higher degree polynomial) always lead to overfitting, i.e., larger test error?

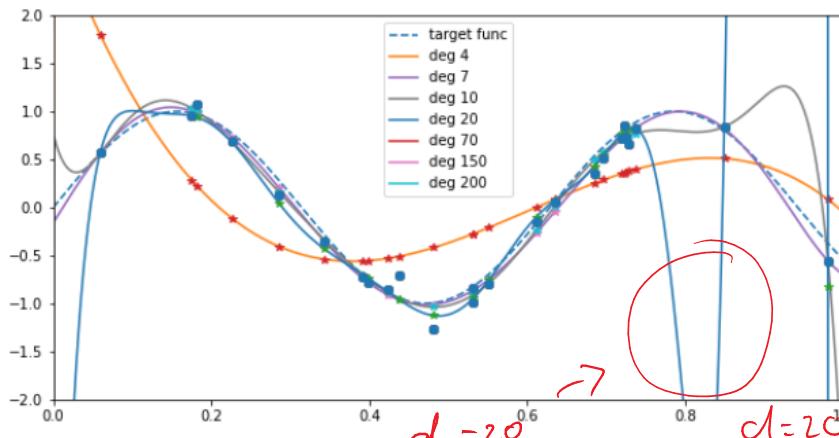
```
In [3]: def fit_poly(X, d, t):
    X_expand = poly_expand(X, d=d, poly_type=poly_type)
    if d > n:
        W = X_expand.T @ np.linalg.inv(X_expand @ X_expand.T) @ t
    else:
        W = np.linalg.inv(X_expand.T @ X_expand) @ X_expand.T @ t
    return W
```

d) Val-loss

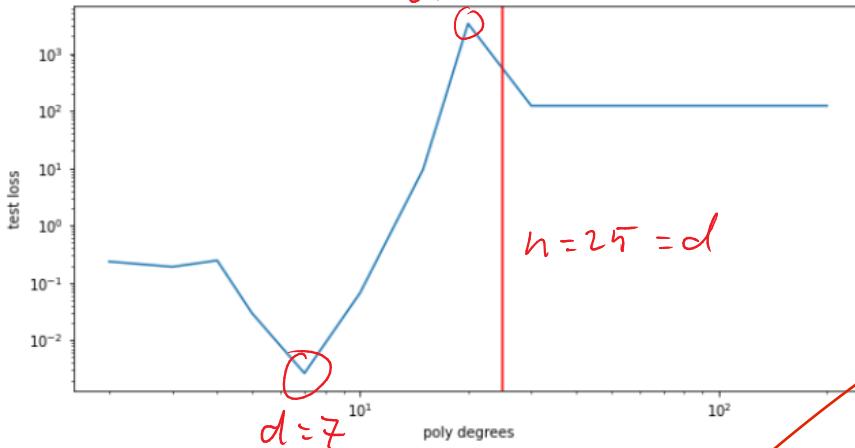
```
2 0.23473638175555447
3 0.19020505096352697
4 0.24537346900176515
5 0.029631242075101965
7 0.002650135121507636
10 0.06641133900488616
15 9.582366859521946
20 3312.0732386744157
30 122.37598629957265
50 122.37598629957265
70 122.37598629957265
100 122.37598629957265
150 122.37598629957265
200 122.37598629957265
```

*best!* ! ! !  
*worst!* ! ! !  
*stable* }  
 for  $d \gg n$

for  $n=25$



$d=20$        $d=20 < n=25$



Stable result for  
larger  $d$  because  
of implicit  
regularization!

Overparameterization does not always make the result worse here, because of the implicit  $L_2$ -regularization.  
 However, generally, one should always lean towards a smaller degree for polynomial interpolation, as clearly evident here.

### 3 Linear Regression and Optimization (3pts)

In this question we will study the solutions to linear regression from the variants of gradient descent that we saw in class.

- Stochastic Gradient Descent (1pt)** Recall that SGD estimates the gradient from a single randomly sampled training example and updates the parameters until convergence. In the overparameterized  $d > n$  case, assume that SGD also converges to a solution  $\hat{w}$  such that  $X\hat{w} = t$ . Show that the SGD solution is identical to the minimum norm solution  $w^*$  obtained by gradient descent, i.e.,  $\hat{w} = w^*$ . Hint: Is  $x_i$  contained in the span of  $X$ ? Do the update steps SGD ever leave the span of  $X$ ?
- Mini-Batch SGD (1pt)** Recall that SGD estimates the gradient from mini-batches of data  $B \in \mathbb{R}^{b \times d}$ , where  $1 < b < n$  and  $B$  is taken from the rows of  $X$ . In this case, does mini-batch SGD obtain the minimum norm solution on convergence?

1. Again: We have already shown in Problem 3 that

$w^* = X^T(X^T X)^{-1}t$  is the only solution in  $\text{range}(X^T)$  and that it is the least norm solution.

So here we choose some  $\tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_b^T \end{bmatrix} \in \mathbb{R}^{b \times d}$

where the  $\tilde{x}_i^T$  are also rows of  $X$ . So:  $\text{range}(\tilde{X}^T) \subseteq \text{range}(X^T)$ .  
And the gradient of the new loss function

Will the trajectory of  $w$  if  
i.e.  $w(0) \notin \text{range}(X^T)$ ,  $w(0) \neq 0$ ?

$$\tilde{\mathcal{L}} = \frac{1}{b} \| \tilde{X} \tilde{w} - t \|_2^2,$$

-0.5

$$\frac{\partial \tilde{\mathcal{L}}}{\partial \tilde{w}} = \frac{2}{b} (\tilde{X}^T \tilde{X} \tilde{w} - \tilde{X}^T t) \in \text{range}(\tilde{X}^T) \subseteq \text{range}(X^T)$$

is still constraint to  $\text{range}(X^T)$ . So if SD converges, the only solution it can converge to is still  $w^*$ , the least norm solution.

2. Same as in 1. just that we approximate  $\nabla \mathcal{L}$ :

$$\frac{\partial \mathcal{L}}{\partial w} \approx \frac{1}{m} \sum_{i=1}^m \frac{\partial \tilde{\mathcal{L}}}{\partial \tilde{w}} \Big|_{\tilde{X} = \text{random choice of } b \text{ rows from } X}$$

So the approx gradient is a superposition of vectors (gradients) that are all elements of  $\text{range}(X^T)$ . Therefore, the approximated gradient is itself still element of  $\text{range}(X^T)$  and can again only find the single solution in that set, which is the least norm solution  $w^*$ .

Please be more mathematically precise.

3. **Adaptive Methods (1pt)** We will now consider the behavior of adaptive gradient descent methods, in particular the AdaGrad<sup>a</sup> method. Let  $w_i$  denote the  $i$ -th parameter. A scalar learning rate  $\eta$  is used. At time  $t$  for parameter  $i$ , the update step for AdaGrad is as follows:

$$w_{t+1,i} = w_{t,i} - \frac{\eta}{\sqrt{G_{t,i}} + \epsilon} \nabla_{w_i} \mathcal{L}(w)_{i,i}$$

$$G_{t,i} = G_{t-1,i} + (\nabla_{w_i} \mathcal{L}(w))_{i,i}^2$$

The term  $\epsilon$  is a fixed small scalar used for numerical stability. Intuitively, Adagrad can be thought of as adapting the learning rate in each dimension to efficiently move through badly formed curvatures.

Consider the overparameterized linear regression model ( $d > n$ ). Assume that AdaGrad converges to a solution. Provide a proof or a counterexample for whether AdaGrad always obtains the minimum norm solution. Hint: Compute the 2D case i.e.,  $x_1^\top = [2, 1]$ ,  $w^\top = [0, 0]$ , and  $t = [2]$ .

Does this result hold true for other adaptive methods (RMSProp, Adam) in general?

Yes! ✓

<sup>a</sup><https://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>

My intuition is: No, because adjusting the components individually (not the entire gradient as a whole vector) allows the construction of vectors outside of range( $X^\top$ ). Which also holds true for other adaptive methods that scale individual gradient components!

Counterexample:  $X^\top w = t \Leftrightarrow x_1^\top w = t \Leftrightarrow (2, 1)^\top w = 2$

$$\Rightarrow w^* = X^\top (X^\top X)^{-1} t = x_1^\top (x_1^\top x_1)^{-1} t$$

$$= \begin{pmatrix} 2 \\ 1 \end{pmatrix} \cdot \frac{1}{5} \cdot 2 = \frac{2}{5} x_1$$

is the least norm solution.

$$\frac{\partial \mathcal{L}}{\partial w} = 2 \left( x_1^\top x_1^\top w - x_1^\top t \right) = \begin{pmatrix} 8 & 4 \\ 4 & 2 \end{pmatrix} w - 4 x_1$$

```
In [1]: import numpy as np
In [2]: x1 = np.array([2, 1])
t = 2
In [3]: def dL(w):
    return np.array([[8, 4], [4, 2]]) @ w - 4 * x1
In [4]: def ada_step(w, eta, G):
    grad = dL(w)
    w[0] -= eta / (np.sqrt(G[0])) * grad[0]
    w[1] -= eta / (np.sqrt(G[1])) * grad[1]
    G += grad ** 2
    return w, G
In [5]: eta = 0.01
G = np.ones(2)
w_ada = np.zeros(2)
w_regular = np.zeros(2)

for i in range(1000000):
    w_ada, G = ada_step(w_ada, eta, G)
    w_regular -= eta * dL(w_regular) # regular GD
In [9]: print("Theoretical Solution:", 2/5*x1)
print("Gradient at theoretical Solution:", dL(2/5*x1))
print("GD Solution:", w_regular)
print("Gradient at GD Solution:", dL(w_regular))
print("Prediction of GD Solution:", x1 @ w_regular)
print("AdaGrad Solution:", w_ada)
print("Gradient at AdaGrad Solution:", dL(w_ada))
print("Prediction of Adagrad Solution:", x1 @ w_ada)
```

✗

Theoretical Solution: [0.8 0.4]  
 Gradient at theoretical Solution: [0. 0.]  
 GD Solution: [0.8 0.4]  
 Gradient at GD Solution: [-5.32907052e-15 -2.66453526e-15] ≈ 0  
 Prediction of GD Solution: 1.999999999999987 ≈ 0  
 AdaGrad Solution: [0.68022007 0.63955986]  
 Gradient at AdaGrad Solution: [-6.37712105e-13 -3.18856053e-13] ≈ 0  
 Prediction of Adagrad Solution: 1.999999999998406 ≈ 0

✓

Clearly, AdaGrad converges to a different solution!