

1. Bias-Variance Tradeoff Review

- (a) Show that we can decompose the expected mean squared error into three parts: bias, variance, and irreducible error σ^2 :

$$\text{Error} = \text{Bias}^2 + \text{Variance} + \sigma^2$$

Formally, suppose we have a randomly sampled training set \mathcal{D} (drawn independently of our test data), and we select an estimator denoted $\hat{\theta} = \hat{\theta}(\mathcal{D})$ (for example, via empirical risk minimization). The expected mean squared error for a test input x can be decomposed as below:

$$\mathbb{E}_{Y \sim p(y|x), \mathcal{D}}[(Y - f_{\hat{\theta}(\mathcal{D})}(x))^2] = \text{Bias}(f_{\hat{\theta}(\mathcal{D})}(x))^2 + \text{Var}(f_{\hat{\theta}(\mathcal{D})}(x)) + \sigma^2$$

You may find it helpful to recall the formulaic definitions of Variance and Bias, reproduced for you below:

$$\begin{aligned}\text{Bias}(f_{\hat{\theta}(\mathcal{D})}(x)) &= \mathbb{E}_{Y \sim p(Y|x), \mathcal{D}}[f_{\hat{\theta}(\mathcal{D})}(x) - Y] \\ \text{Var}(f_{\hat{\theta}(\mathcal{D})}(x)) &= \mathbb{E}_{\mathcal{D}}[(f_{\hat{\theta}(\mathcal{D})}(x) - \mathbb{E}[f_{\hat{\theta}(\mathcal{D})}(x)])^2]\end{aligned}$$

Pf: For any variable x , we have $\text{var}(x) = E(x^2) - (E(x))^2$

let $D[x]$ denote the variance of x .

$$\begin{aligned}R.H.S. &= (E[\hat{\theta}(\mathcal{D})(x) - Y])^2 + D[\hat{\theta}(\mathcal{D})(x)] + \sigma^2 \\ &= E[(\hat{\theta}(\mathcal{D})(x) - Y)^2] - D[\hat{\theta}(\mathcal{D})(x) - Y] + D[\hat{\theta}(\mathcal{D})(x)] + \sigma^2 \\ &= E[(\hat{\theta}(\mathcal{D})(x) - Y)^2] - D[\hat{\theta}(\mathcal{D})(x)] - D[Y] + D[\hat{\theta}(\mathcal{D})(x)] + \sigma^2 \\ &= E[(\hat{\theta}(\mathcal{D})(x) - Y)^2] - D[Y] + \sigma^2.\end{aligned}$$

let $\sigma^2 = D[Y]$, we have R.H.S. = L.H.S. Q.E.D.

- (b) Suppose our training dataset consists of $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ where the only randomness is coming from the label vector $Y = X\theta^* + \varepsilon$ where θ^* is the true underlying linear model and each noise variable ε_i is i.i.d. with zero mean and variance 1. We use ordinary least squares to estimate a $\hat{\theta}$ from this data. Calculate the bias and covariance of the $\hat{\theta}$ estimate and use that to compute the bias and variance of the prediction at particular test inputs x . Recall that the OLS solution is given by

$$\hat{\theta} = (X^\top X)^{-1} X^\top Y,$$

where $X \in \mathbb{R}^{n \times d}$ is our (nonrandom) data matrix, $Y \in \mathbb{R}^n$ is the (random) vector of training targets. For simplicity, assume that $X^\top X$ is diagonal.

For a particular input x :

$$\textcircled{1} \quad \text{bias} = \hat{\theta}(x) - \theta^*(x) = x^\top (X^\top X)^{-1} X^\top Y - x^\top \theta^* + \varepsilon$$

$$\textcircled{2} \quad \text{variance} = (\hat{\theta}(x) - E[\theta^*(x)])^2 = (x^\top (X^\top X)^{-1} X^\top Y - x^\top \theta^*)^2$$

Consider the corresponding expectations:

$$\begin{aligned} E[\text{bias}] &= E[x^T(x^T x)^{-1} x^T Y - x^T \theta^* + \epsilon] \\ &= E[x^T (\hat{\theta}^T x)^{-1} x^T Y - \hat{\theta}^T] \\ &= E[x^T (\hat{\theta} - \theta^*)] = 0. \end{aligned}$$

$$\begin{aligned} E[\text{variance}] &= E[(x^T \hat{\theta} - x^T \theta^*)^2] = E[x (\theta^* - \hat{\theta}) (\hat{\theta} - \theta^*) x^T] \\ &= x E[(\theta^* - \hat{\theta}) (\hat{\theta} - \theta^*)] x^T \\ &= x \end{aligned}$$

2. Least Squares and the Min-norm problem from the Perspective of SVD

Consider the equation $Xw = y$, where $X \in \mathbb{R}^{m \times n}$ is a non-square data matrix, w is a weight vector, and y is vector of labels corresponding to the datapoints in each row of X .

Let's say that $X = U\Sigma V^T$ is the (full) SVD of X . Here, U and V are orthonormal square matrices, and Σ is an $m \times n$ matrix with non-zero singular values (σ_i) on the "diagonal".

For this problem, we define Σ^\dagger an $n \times m$ matrix with the reciprocals of the singular values ($\frac{1}{\sigma_i}$) along the "diagonal".

- (a) First, consider the case where $m > n$, i.e. our data matrix X has more rows than columns (tall matrix) and the system is overdetermined. **How do we find the weights w that minimizes the error between Xw and y ?** In other words, we want to solve $\min_w \|Xw - y\|^2$.

$$\text{a)} \quad \text{let } f(w) = \|Xw - y\|^2 = (Xw - y)^T(Xw - y).$$

$$= (w^T X^T - y^T)(Xw - y)$$

$$= w^T X^T X w - y^T X w - w^T X^T y + y^T y.$$

$$\frac{dt}{dw} = 2X^T X w - 2X^T y.$$

$$\text{let } \frac{dt}{dw} = 0, \text{ we have } X^T X w = X^T y.$$

Assume $X^T X$ is nonsingular, we have $\hat{w} = (X^T X)^{-1} X^T y$.

- (b) Plug in the SVD $X = U\Sigma V^T$ and simplify. Be careful with dimensions!

b) Plug in $x = U\Sigma V^T$, we have

$$\begin{aligned} \hat{w} &= [(U\Sigma V^T)^T (U\Sigma V^T)]^{-1} (U\Sigma V^T)^T y \\ &= (V\Sigma^T U^T U\Sigma V^T)^{-1} (V\Sigma^T U^T) y. \end{aligned}$$

Because U and V are orthonormal, $U^T U = I_m$, $V^T V = I_n$,

$$\begin{aligned} \hat{w} &= (V\Sigma^T \Sigma V^T)^{-1} (V\Sigma^T U^T) y \\ &= [V \left(\begin{smallmatrix} \frac{1}{\sigma_1^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n^2} \end{smallmatrix} \right) V^T]^{-1} (V\Sigma^T U^T) y \\ &= V \left(\begin{smallmatrix} \frac{1}{\sigma_1^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n^2} \end{smallmatrix} \right) V^T V \Sigma^T U^T y \\ &= V \left(\begin{smallmatrix} \frac{1}{\sigma_1^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n^2} \end{smallmatrix} \right) \Sigma^T U^T y \\ &= V \Sigma^T (\Sigma^T)^{-1} \Sigma^T U^T y \\ &= V \Sigma^T (I_n - 0) U^T y \end{aligned}$$

Denote $\left(\begin{smallmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{smallmatrix} \right)$ as Σ_n , $\Sigma^T = (\Sigma_n^{-1} 0)$

$$\hat{w} = V (\Sigma_n^{-1} 0) (I_n - 0) U^T y$$

$$= V(\Sigma^{-1} 0)U^T y$$

$$= V\Sigma^+ U^T y.$$

- (c) You'll notice that the least-squares solution is in the form $w^* = Ay$. **What happens if we left-multiply X by our matrix A ?** This is why the matrix A of the least-squares solution is called the left-inverse.

$$\text{c)} \quad A = V\Sigma^+ U^T. \quad AX = V\Sigma^+ U^T U \Sigma V^T = V\Sigma^+ \Sigma V^T \\ = V(\Sigma_n^{-1} 0) \begin{pmatrix} \Sigma_n \\ 0 \end{pmatrix} V^T = V I_m V^T = I_n.$$

We get a n -D identity matrix.

- (d) Now, let's consider the case where $m < n$, i.e. the data matrix X has more columns than rows and the system is underdetermined. There exist infinitely many solutions for w , but we seek the minimum-norm solution, ie. we want to solve $\min \|w\|^2$ s.t. $Xw = y$. **What is the minimum norm solution?**

d) Using Lagrange multipliers:

$$L(w, \lambda) = \frac{1}{2} w^T w - \lambda^T (Xw - y).$$

let $\frac{\partial L}{\partial w} = 2w - X^T \lambda = 0$

$$\left| \begin{array}{l} \left\{ \begin{array}{l} X^T \lambda = w \\ Xw = y \end{array} \right. \\ \Rightarrow \left\{ \begin{array}{l} \lambda = (X X^T)^{-1} y \\ w = X^T (X X^T)^{-1} y. \end{array} \right. \end{array} \right.$$

- (e) Plug in the SVD $X = U\Sigma V^T$ and simplify. Be careful with dimensions!

$$\begin{aligned} e) \quad w &= V\Sigma^T U^T (U\Sigma V^T V\Sigma^T U^T)^{-1} y \\ &= V\Sigma^T [U\Sigma\Sigma^T U^T]^{-1} y \\ &= V\Sigma^T [\Sigma\Sigma^T]^{-1} y \\ &= V\Sigma^T (\Sigma\Sigma^T)^{-1} \cdot U^T y \end{aligned} \quad \begin{aligned} &= V\Sigma^T (\Sigma\Sigma^T)^{-1} U^T y \\ &= V\Sigma^+ U^T y. \end{aligned}$$

- (f) You'll notice that the min-norm solution is in the form $w^* = By$. **What happens if we right-multiply X by our matrix B ?** This is why the matrix B of the min-norm solution is called the right-inverse.

$$\text{f)} \quad B = V\Sigma^+ U^T \quad . \quad XB = U\Sigma V^T V\Sigma^+ U^T = U\Sigma(V^T V)\Sigma^+ U^T = U\Sigma\Sigma^+ U^T$$

Denote $\begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{pmatrix}$ as Σ_m .

$$\begin{aligned} XB &= U(\Sigma_m 0) \begin{pmatrix} \Sigma_m^{-1} \\ 0 \end{pmatrix} U^T \\ &= UU^T = I_m \end{aligned}$$

We get a m -D identity matrix.

3. The 5 Interpretations of Ridge Regression

- (a) *Perspective 1: Optimization Problem.* Ridge regression can be understood as the unconstrained optimization problem

$$\underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2, \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is a data matrix, and $\mathbf{y} \in \mathbb{R}^n$ is the target vector of measurement values. What's new compared to the simple OLS problem is the addition of the $\lambda \|\mathbf{w}\|^2$ term, which can be interpreted as a "penalty" on the weights being too big.

Use vector calculus to expand the objective and solve this optimization problem for \mathbf{w} .

a) Sol : Denote $f(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$

$$\begin{aligned} &= (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \\ &= (\mathbf{y}^\top - \mathbf{w}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w} \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\mathbf{w} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} + \lambda \mathbf{w}^\top \mathbf{w}. \end{aligned}$$

$$\frac{d f(\mathbf{w})}{d \mathbf{w}} = 2 \mathbf{X}^\top \mathbf{X}\mathbf{w} - 2 \mathbf{y}^\top \mathbf{y} + 2\lambda \mathbf{w}.$$

let $\frac{d f}{d \mathbf{w}} = 0$. $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{w} = \mathbf{y}^\top \mathbf{y}$.

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{y}^\top$$

- (b) *Perspective 2: "Hack" of shifting the Singular Values.* In the previous part, you should have found the optimal \mathbf{w} is given by

$$n \geq d. \quad \mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

(If you didn't get this, you should check your work for the previous part).

Let $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the (full) SVD of the \mathbf{X} . Recall that \mathbf{U} and \mathbf{V} are square orthonormal (norm-preserving) matrices, and Σ is a $n \times d$ matrix with singular values σ_i along the "diagonal". **Plug this into the Ridge Regression solution and simplify. What happens when the singular values when $\sigma_i \ll \lambda$? What about when $\sigma_i \gg \lambda$?**

b) Sol: Plug in $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$. We have:

$$\begin{aligned} \mathbf{w} &= ((\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{U}\Sigma\mathbf{V}^\top + 2\lambda \mathbf{I}_d)^{-1} (\mathbf{U}\Sigma\mathbf{V}^\top)^\top \mathbf{y} \\ &= (\mathbf{V}^\top \mathbf{U}^\top \mathbf{U}\Sigma\mathbf{V}^\top + 2\lambda \mathbf{I}_d)^{-1} \mathbf{V}^\top \mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}^\top \Sigma \mathbf{U}^\top + 2\lambda \mathbf{I}_d)^{-1} \mathbf{V}^\top \mathbf{U}^\top \mathbf{y} \end{aligned}$$

When $\sigma_i \ll \lambda$: $\mathbf{w} \approx [2\lambda \mathbf{I}_d]^{-1} \mathbf{V}^\top \mathbf{U}^\top \mathbf{y} = \frac{1}{2\lambda} \mathbf{I}_d \mathbf{V}^\top \mathbf{U}^\top \mathbf{y} = \frac{1}{2\lambda} \mathbf{V}^\top \mathbf{U}^\top \mathbf{y}$.

Compared with the solution $\hat{\mathbf{w}} = \mathbf{V}^\top \mathbf{U}^\top \mathbf{y}$ in 1(b),

its components are likely to be smaller due to λ .

when $\pi \gg \lambda$: $w \approx V\Sigma^{-1}y$, which is consistent with the soln in 1(b).

generally λ controls how much constraint we want to put on the magnitude of w .

(c) *Perspective 3: Maximum A Posteriori (MAP) estimation.* Ridge Regression can be viewed as finding the MAP estimate when we apply a prior on the (now viewed as random parameters) \mathbf{W} . In particular, we can think of the prior for \mathbf{W} as being $\mathcal{N}(0, I)$ and view the random Y as being generated using $Y = \mathbf{x}^T \mathbf{W} + \sqrt{\lambda}N$ where the noise N is distributed iid (across training samples) as $\mathcal{N}(0, 1)$. At the vector level, we have $\mathbf{Y} = X\mathbf{W} + \sqrt{\lambda}\mathbf{N}$. Note that the X matrix whose rows are the n different training points are not random.

Show that (1) is the MAP estimate for \mathbf{W} given an observation $\mathbf{Y} = \mathbf{y}$.

$$\text{c) Pf: } P(W=w | Y=y) = \frac{P(Y=y | W=w) P(W=w)}{P(Y=y)}$$

Since $P(Y=y)$ is a constant, we only need to consider maximize

$$P(Y=y | W=w) P(W=w) \quad (\text{denoted as } f(w)).$$

$$P(XW + \sqrt{\lambda}N = y | W=w) = P(N = \frac{y - Xw}{\sqrt{\lambda}} | W=w)$$

Note that N and W are independent.

$$P(N = \frac{y - Xw}{\sqrt{\lambda}} | W=w) = P(N = \frac{y - Xw}{\sqrt{\lambda}})$$

$$f(w) = P(N = \frac{y - Xw}{\sqrt{\lambda}}) P(W=w).$$

Because $N \sim \mathcal{N}(0, 1)$, $W \sim \mathcal{N}(0, I)$.

$$f(w) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y - Xw)^T(y - Xw)}{2}} \cdot \frac{1}{\sqrt{\lambda n}} e^{\frac{w^T w}{2}}$$

$$\frac{d}{dw} f(w) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y - Xw)^T(y - Xw)}{2}} \left(\frac{w^T w}{2} + \frac{X^T y + X^T X w}{\lambda} \right)$$

$$\text{let } \frac{d}{dw} f(w) = 0, \quad -\frac{X^T y + X^T X w}{\lambda} + w = 0$$

$$(X^T X + \lambda I_n) w = X^T y.$$

Therefore $w = (X^T X + \lambda I_n)^{-1} X^T y$ is the MAP estimate. Q.E.D.

- (d) *Perspective 4: Fake Data.* Another way to interpret “ridge regression” is as the ordinary least squares for an augmented data set — i.e. adding a bunch of fake data points to our data. Consider the following augmented measurement vector $\hat{\mathbf{y}}$ and data matrix $\hat{\mathbf{X}}$:

$$\hat{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_d \end{bmatrix} \quad \hat{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_d \end{bmatrix},$$

where $\mathbf{0}_d$ is the zero vector in \mathbb{R}^d and $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ is the identity matrix. **Show that the classical OLS optimization problem $\operatorname{argmin}_{\mathbf{w}} \|\hat{\mathbf{y}} - \hat{\mathbf{X}}\mathbf{w}\|_2^2$ has the same minimizer as (1).**

$$\begin{aligned} d). \text{ Note } \|\hat{\mathbf{y}} - \hat{\mathbf{X}}\mathbf{w}\|_2^2 &= \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_d \end{bmatrix} - \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda} \mathbf{I}_d \end{bmatrix} \mathbf{w} \right\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbf{y} - \mathbf{X}\mathbf{w} \\ \sqrt{\lambda} \mathbf{I}_d \mathbf{w} \end{bmatrix} \right\|_2^2 \\ &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \end{aligned}$$

which is of the same form with (1).

- (e) *Perspective 5: Fake Features.* For this last interpretation, let's instead construct an augmented design matrix in the following way:

$$\tilde{\mathbf{X}} = [X \ \sqrt{\lambda} \mathbf{I}_n]$$

i.e. we stack X with $\sqrt{\lambda} \mathbf{I}_n$ horizontally. Now our problem is underdetermined: the new dimension $d + n$ is larger than the number of points n . Therefore, there are infinitely many values $\eta \in \mathbb{R}^{d+n}$ for which $\tilde{\mathbf{X}}\eta = \mathbf{y}$. We are interested in the **min-norm** solution, ie. the solution to

$$\underset{\eta}{\operatorname{argmin}} \|\eta\|_2^2 \text{ s.t. } \tilde{\mathbf{X}}\eta = \mathbf{y}. \quad (2)$$

Show that this is yet another form of ridge regression and that the first d coordinates of η^* form the minimizer of (1).

o). Using Lagrange multipliers:

$$L(\eta, \lambda) = \frac{1}{2}\eta^T\eta + \lambda^T(\mathbf{y} - \tilde{\mathbf{X}}\eta)$$

$$\text{let } \frac{\partial L}{\partial \eta} = \eta - \tilde{\mathbf{X}}^T\lambda = 0 \quad \dots \textcircled{1}$$

$$\mathbf{y} - \tilde{\mathbf{X}}\eta = 0 \quad \dots \textcircled{2}$$

From \textcircled{1}.\textcircled{2}, we can deduce $\begin{cases} \lambda = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\mathbf{y}, \\ \eta = \tilde{\mathbf{X}}^T(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\mathbf{y} \end{cases}$

$$\begin{aligned} \eta &= [\begin{smallmatrix} \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \mathbf{I}_n \end{smallmatrix}] ([\mathbf{x} \ \mathbf{x}^T \mathbf{I}_n] [\begin{smallmatrix} \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \mathbf{I}_n \end{smallmatrix}])^{-1} \mathbf{y} \\ &= [\begin{smallmatrix} \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \mathbf{I}_n \end{smallmatrix}] (\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} \mathbf{y} \end{aligned}$$

The first d dimension of η is $\mathbf{x}^T(\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$,

The soln in (1) is $(\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I}_d)^{-1} \mathbf{x}^T \mathbf{y}$.

We need to further prove $\mathbf{x}^T(\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I}_d)^{-1} \mathbf{x}^T$

which is $(\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I}_d) \mathbf{x}^T (\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} = \mathbf{x}^T$

$$\begin{aligned} \text{L.H.S.} &= (\mathbf{x}^T \mathbf{x} \mathbf{x}^T + \lambda \mathbf{x}^T) (\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} = \mathbf{x}^T (\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n) (\mathbf{x} \mathbf{x}^T + \lambda \mathbf{I}_n)^{-1} \\ &= \mathbf{x}^T. \end{aligned}$$

Q.E.D

- (g) We know that the solution to ridge regression (1) is given by $\hat{\mathbf{w}}_r = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$. What happens when $\lambda \rightarrow \infty$? It is for this reason that sometimes ridge regularization is referred to as "shrinkage."

When $\lambda \rightarrow \infty$, $\hat{\mathbf{w}}_r \rightarrow (\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \lambda^{-1} \mathbf{X}^T \mathbf{y}$.

Thus $\|\hat{\mathbf{w}}_r\| = \frac{1}{\lambda} \|\mathbf{X}^T \mathbf{y}\| \rightarrow 0$.

- (h) What happens to the solution of ridge regression when you take the limit $\lambda \rightarrow 0$? Consider both the cases when X is wide (underdetermined system) and X is tall (overdetermined system).

when $\lambda \rightarrow 0$, $\hat{w} \rightarrow (X^T X)^{-1} X^T y$.

The solution becomes the OLS solution

no matter X is tall or wide.

- (f) We know that the Moore-Penrose pseudo-inverse for an underdetermined system (wide matrix) is given by $A^\dagger = A^T (AA^T)^{-1}$, which corresponds to the min-norm solution for $A\eta = z$. That is, the optimization problem

$$\operatorname{argmin} \|\eta\|^2 \text{ s.t. } A\eta = z$$

is solved by $\eta = A^\dagger z$. Let \hat{w} be the minimizer of (1).

Use the pseudo-inverse to show that solving to the optimization problem in (2) yields

$$\hat{w} = X^T (XX^T + \lambda I)^{-1} y$$

Then, show that this is equivalent to the standard formula for Ridge Regression

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y$$

+> $\text{Pf}:$ ^① Use pseudo-inverse:

$$\begin{aligned}\eta &= \hat{X}^T (\hat{X} \hat{X}^T)^{-1} y \\ &= \begin{bmatrix} X^T \\ \lambda I_n \end{bmatrix} \left([X \quad \lambda I_n] \begin{bmatrix} X^T \\ \lambda I_n \end{bmatrix} \right)^{-1} y \\ &= \begin{bmatrix} X^T \\ \lambda I_n \end{bmatrix} (XX^T + \lambda I_n)^{-1} y \\ &= \begin{bmatrix} X^T (XX^T + \lambda I_n)^{-1} y \\ \lambda I_n (XX^T + \lambda I_n)^{-1} y \end{bmatrix}\end{aligned}$$

Therefore $\hat{w} = X^T (XX^T + \lambda I_n)^{-1} y$.

^② To prove $(X^T X + \lambda I_n)^{-1} X^T y = X^T (XX^T + \lambda I_n)^{-1} y$,

We can first prove $(X^T X + \lambda I_n)^{-1} X^T = X^T (XX^T + \lambda I_n)^{-1}$

Since $(X^T)^{-1} (X^T X + \lambda I_n)^{-1} X^T = (X^T X^T + \lambda X^T)^{-1} X^T$

$$= [X^T (XX^T + \lambda I_n)]^{-1} X^T$$

$$= (XX^T + \lambda I_n)^{-1} (X^T)^{-1} X^T = (XX^T + \lambda I_n)^{-1}.$$

Hence

$$(X^T X + \lambda I_n)^{-1} X^T = X^T (XX^T + \lambda I_n)^{-1}$$

4. General Case Tikhonov Regularization

Consider the optimization problem:

$$\min_{\mathbf{x}} \|W_1(A\mathbf{x} - \mathbf{b})\|_2^2 + \|W_2(\mathbf{x} - \mathbf{c})\|_2^2$$

Where W_1 , A , and W_2 are matrices and \mathbf{x} , \mathbf{b} and \mathbf{c} are vectors. W_1 can be viewed as a generic weighting of the residuals and W_2 along with c can be viewed as a generic weighting of the parameters.

- (a) **Solve this optimization problem manually** by expanding it out as matrix-vector products, setting the gradient to 0, and solving for \mathbf{x} .

$$f(\mathbf{x}) = (A\mathbf{x} - \mathbf{b})^T W_1^T W_1 (A\mathbf{x} - \mathbf{b}) + (\mathbf{x} - \mathbf{c})^T W_2^T W_2 (\mathbf{x} - \mathbf{c}).$$

$$= (\mathbf{x}^T A^T - \mathbf{b}^T) W_1^T W_1 (A\mathbf{x} - \mathbf{b}) + (\mathbf{x}^T - \mathbf{c}^T) W_2^T W_2 (\mathbf{x} - \mathbf{c}).$$

$$\frac{df}{d\mathbf{x}} = 2A^T W_1^T W_1 A \mathbf{x} - 2A^T W_1^T W_1 \mathbf{b} + 2W_2^T W_2 \mathbf{x} - 2W_2^T W_2 \mathbf{c}.$$

$$\text{let } \frac{df}{d\mathbf{x}} = 0. \quad (A^T W_1^T W_1 A + W_2^T W_2) \mathbf{x} = A^T W_1^T W_1 \mathbf{b} + W_2^T W_2 \mathbf{c}$$

Assume $A^T W_1^T W_1 A + W_2^T W_2$ is invertible,

$$\mathbf{x} = (A^T W_1^T W_1 A + W_2^T W_2)^{-1} (A^T W_1^T W_1 \mathbf{b} + W_2^T W_2 \mathbf{c})$$

- (b) Construct an appropriate matrix C and vector d that allows you to rewrite this problem as

$$\min_{\mathbf{x}} \|Cx - d\|^2$$

and use the OLS solution ($\mathbf{x}^* = (C^T C)^{-1} C^T d$) to solve. Confirm your answer is in agreement with the previous part.

$$f(\mathbf{x}) = \mathbf{x}^T (A^T W_1^T W_1 A + W_2^T W_2) \mathbf{x}$$

$$- (\mathbf{b}^T W_1^T W_1 A + \mathbf{c}^T W_2^T W_2) \mathbf{x} - \mathbf{x}^T (A^T W_1^T W_1 \mathbf{b} + W_2^T W_2 \mathbf{c}) + \mathbf{b}^T \mathbf{b}$$

$$(Cx - d)^T (Cx - d) = \mathbf{x}^T C^T C \mathbf{x} - d^T C \mathbf{x} - \mathbf{x}^T C^T d + d^T d$$

$$\begin{cases} A^T W_1^T W_1 A + W_2^T W_2 = C^T C \\ b^T W_1^T W_1 A + c^T W_2^T W_2 = d^T C. \end{cases}$$

[I have no idea how to solve them].

$$\begin{aligned} \text{Therefore } \mathbf{x}^* &= (C^T C)^{-1} C^T d = (A^T W_1^T W_1 A + W_2^T W_2)^{-1} (b^T W_1^T W_1 A + c^T W_2^T W_2) \\ &= (A^T W_1^T W_1 A + W_2^T W_2)^{-1} (A^T W_1^T W_1 \mathbf{b} + W_2^T W_2 \mathbf{c}) \end{aligned}$$

- (c) Choose a W_1 , W_2 , and c such that this reduces to the simple case of ridge regression that you've seen in the previous problem, $\mathbf{x}^* = (A^T A + \lambda I)^{-1} A^T b$.

Let $\begin{cases} w_1 = I \\ w_2 = \lambda I \\ c = 0 \end{cases}$, we have $f(x) = \|Ax - b\|_2 + \lambda \|x\|_2$, which is the same as 3(a).

5. Coding Fully Connected Networks

In this coding assignment, you will be building a fully-connected neural network from scratch using NumPy. You will have the choice between two options:

Use Google Colab (Recommended). Open [this url](#) and follow the instructions in the notebook.

Use a local Conda environment. Clone <https://github.com/gonglinyuan/cs182hw1> and refer to `README.md` for further instructions.

For this question, please submit a .zip file your completed work to the Gradescope assignment titled “Homework 1 (Code)”. Please answer the following question in your submission of the written assignment:

- (a) Did you notice anything about the comparative difficulty of training the three-layer net vs training the five layer net?

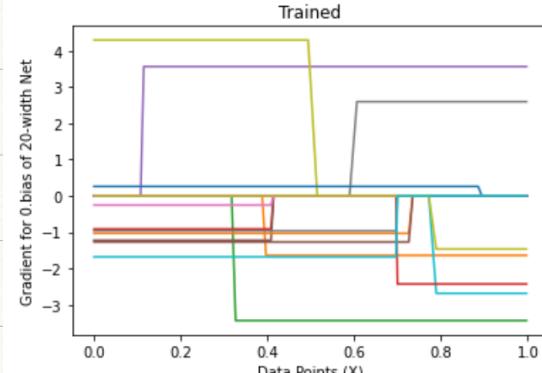
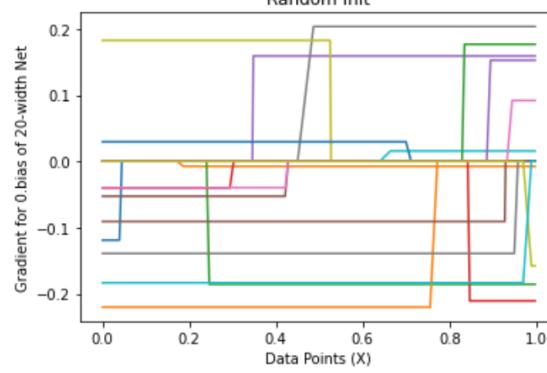
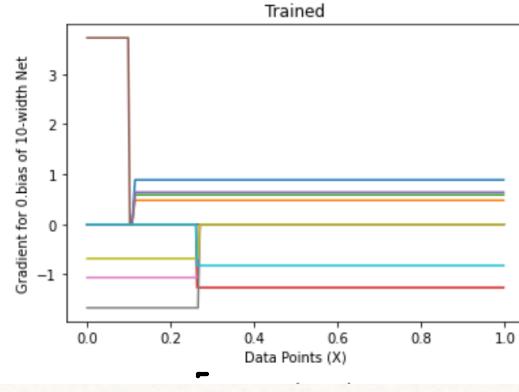
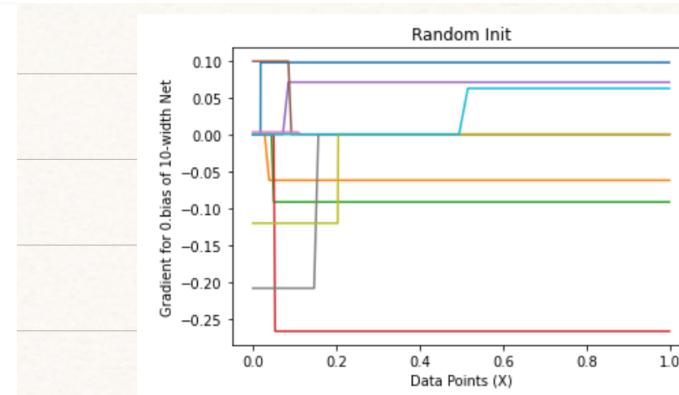
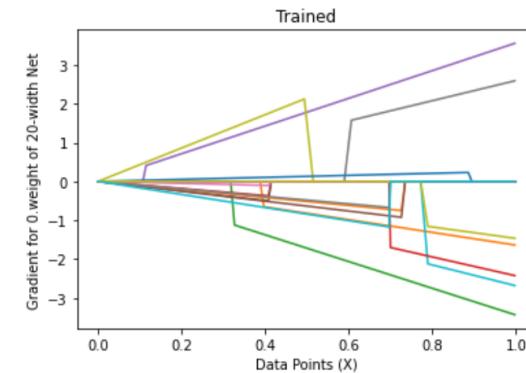
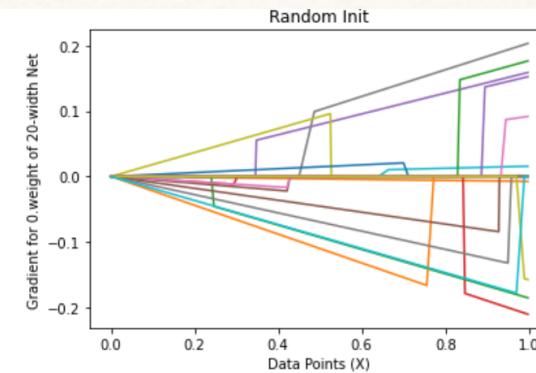
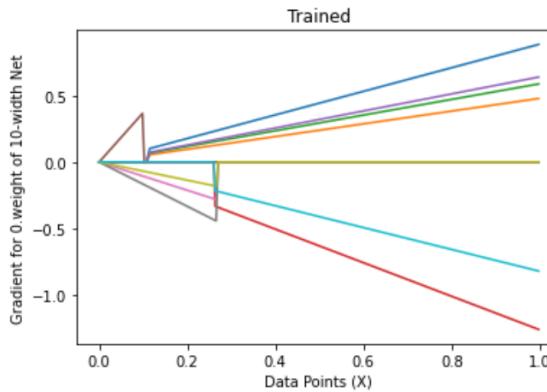
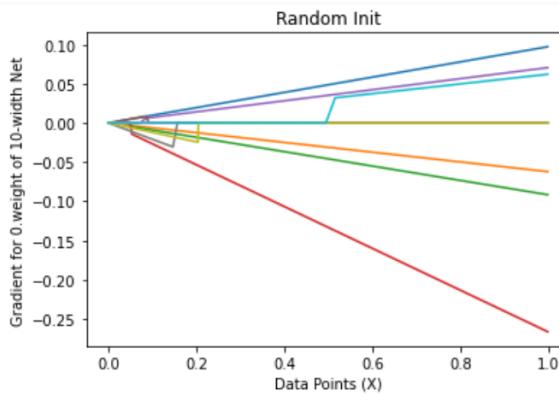
- ① It seems that five layer net is more sensitive to weight-scale ,
the loss of which may explode even if there is only slight
change to the parameters
- ② It takes longer time to train a five layer net.

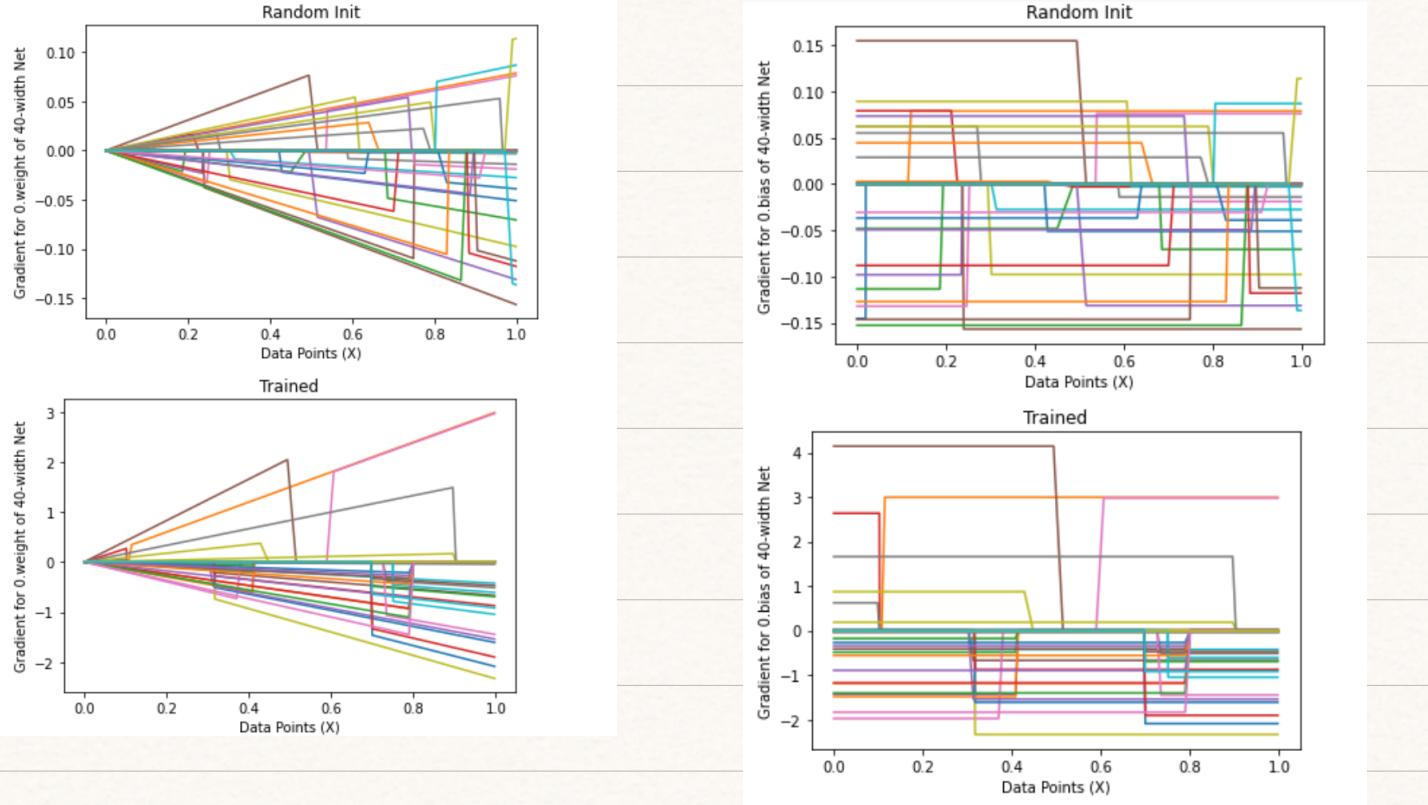
6. Visualizing features from local linearization of neural nets

This problem expects you to modify the Jupyter Notebook you were given in the first discussion section for the course to allow the visualization of the effective “features” that correspond to the local linearization of the network in the neighborhood of the parameters.

We provide you with some starter code on [Google Colab](#). For this question, **please do not submit your code to Gradescope**. Instead, just include your plots and comments regarding the questions in the subparts.

- (a) **Visualize the features corresponding to $\frac{\partial}{\partial w_i^{(1)}} y(x)$ and $\frac{\partial}{\partial b_i^{(1)}} y(x)$ where $w_i^{(1)}$ are the first hidden layer's weights and the $b_i^{(1)}$ are the first hidden layer's biases.** These derivatives should be evaluated at least both the random initialization and the final trained network. When visualizing these features, plot them as a function of the scalar input x , the same way that the notebook plots the constituent “elbow” features that are the outputs of the penultimate layer.





As illustrated, the gradients of $w^{(i)}$ and $b^{(i)}$ are randomly distributed, and after trained, they tend to cluster around certain values.

- (b) During training, we can imagine that we have a generalized linear model with a feature matrix corresponding to the linearized features corresponding to each learnable parameter. We know from our analysis of gradient descent, that the singular values and singular vectors corresponding to this feature matrix are important.

Use the SVD of this feature matrix to plot both the singular values and visualize the “principle features” that correspond to the d -dimensional singular vectors multiplied by all the features corresponding to the parameters.

(HINT: Remember that the feature matrix whose SVD you are taking has n rows where each row corresponds to one training point and d columns where each column corresponds to each of the learnable features. Meanwhile, you are going to be plotting/visualizing the “principle features” as functions of x even at places where you don’t have training points.)

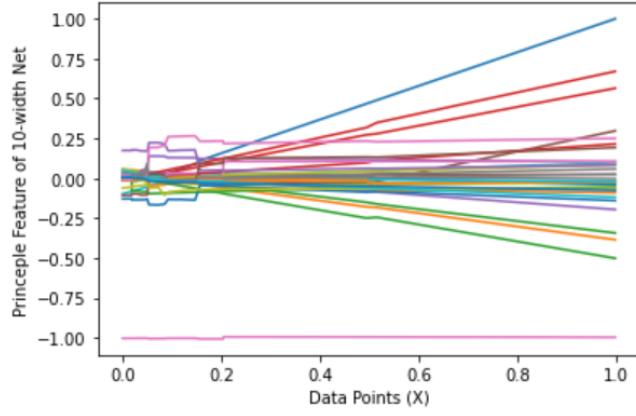
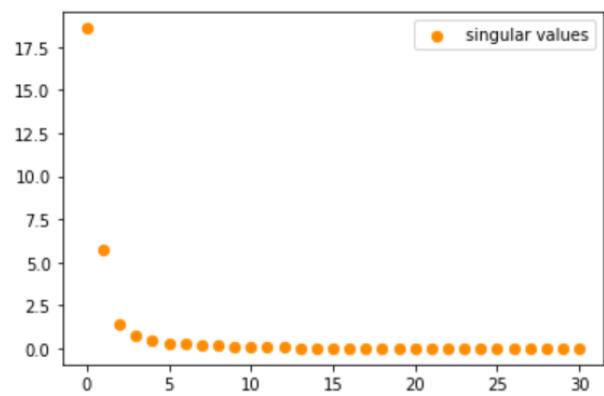
Consider a $m \times n$ feature matrix X (row: training pts, column: feature vecs).

Perform SVD on X : $X = U\Sigma V^T$. ($U \in \mathbb{R}^{m \times d}$, $\Sigma \in \mathbb{R}^{d \times d}$, $V \in \mathbb{R}^{d \times n}$, $d = \min(m, n)$)

Right multiply X by V , we get a $m \times d$ matrix, which is the matrix containing principle features.

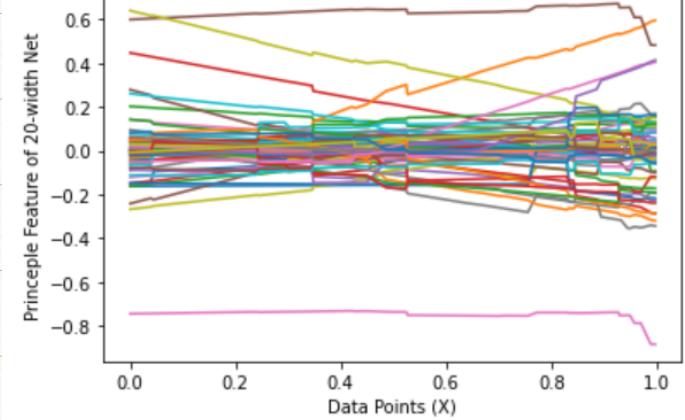
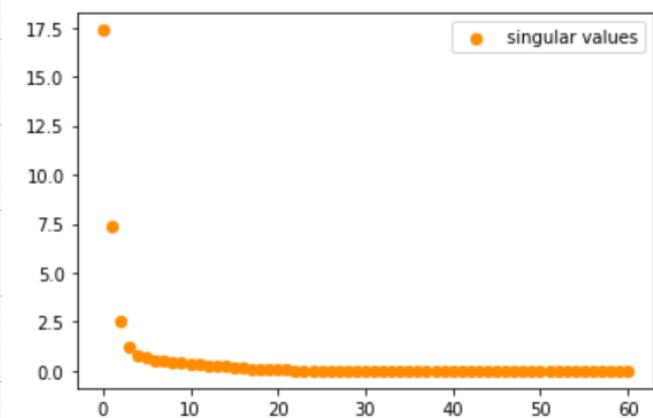
Note: At places where we don't have training pts, we need to corresponding rows with 0.

Width 10

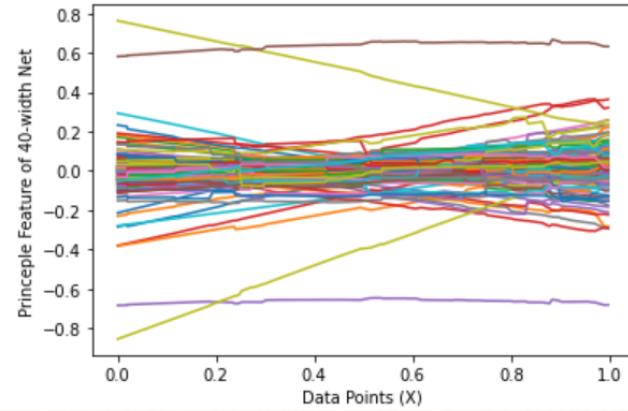
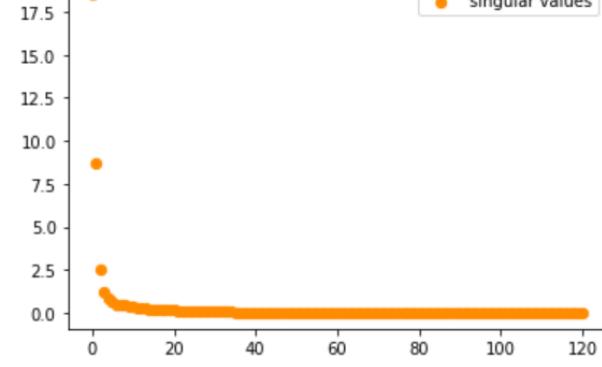


Width 20

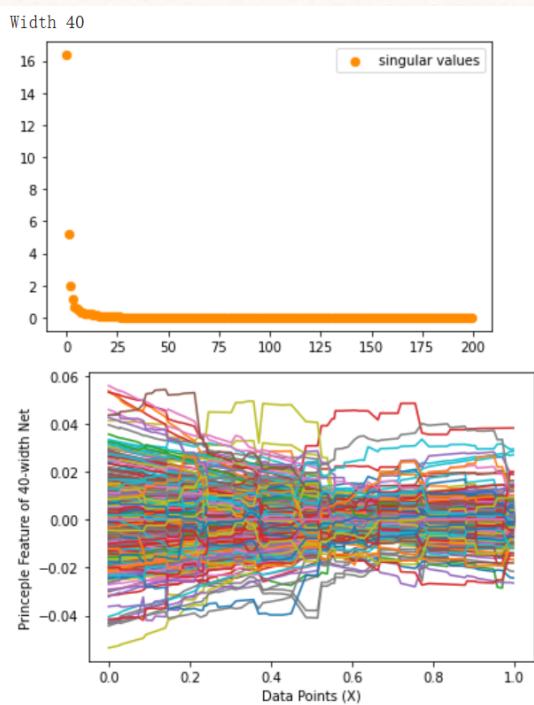
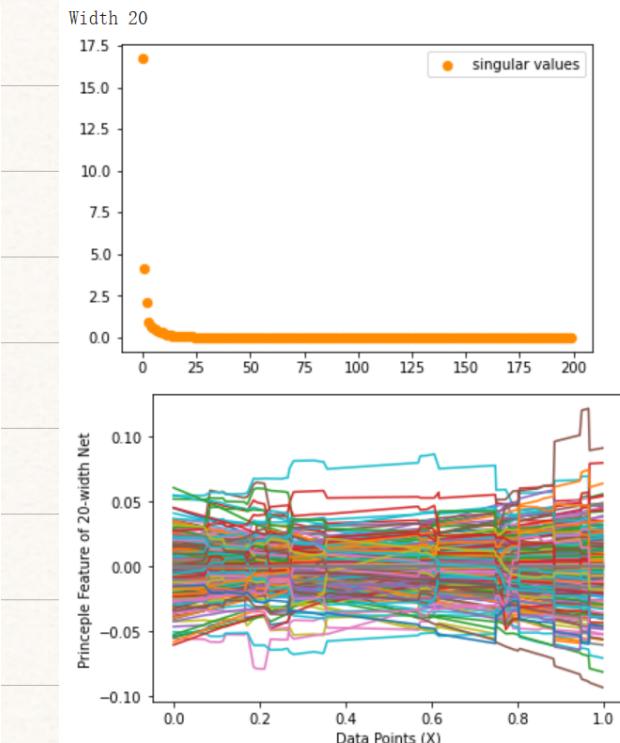
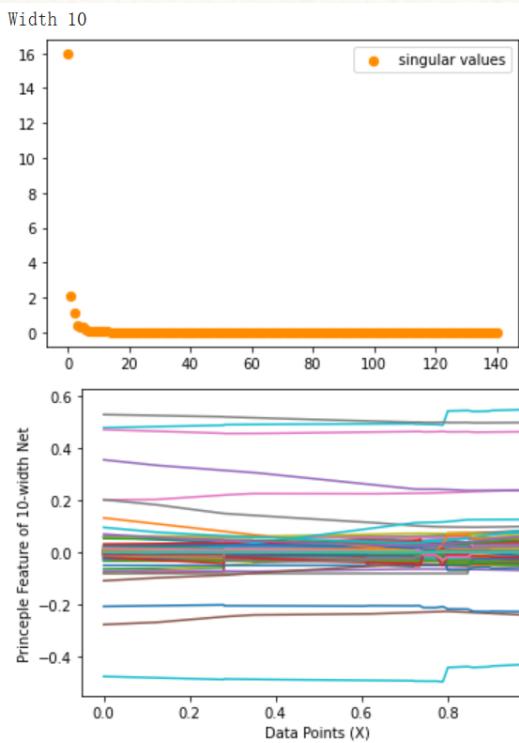
Width 20



Width 40



- (c) Augment the jupyter notebook to add a second hidden layer of the same size as the first hidden layer, fully connected to the first hidden layer. Allow the visualization of the features corresponding to the parameters in both hidden layers, as well as the “principle features” and the singular values.



7. Homework Process and Study Group

Citing sources and collaborators are an important part of life, including being a student!

We also want to understand what resources you find helpful and how much time homework is taking, so we can change things in the future if possible.

- (a) **What sources (if any) did you use as you worked through the homework?**
- (b) **If you worked with someone on this homework, who did you work with?**
List names and student ID's. (In case of homework party, you can also just describe the group.)
- (c) **Roughly how many total hours did you work on this homework? Write it down here where you'll need to remember it for the self-grade form.**

a) <https://zhuanlan.zhihu.com/p/29846048>

<https://numpy.org/doc/stable/reference/routines.linalg.html>

b)	Names	Xiang Fei
	SID	3038733024.

c) 20 h.