

Announcements





Principal Component Analysis

Machine Learning – CSE546

Kevin Jamieson

University of Washington

November 8, 2018

Linear projections

Given $x_1, \dots, x_n \in \mathbb{R}^d$, for $q \ll d$ find a compressed representation with $\lambda_1, \dots, \lambda_n \in \mathbb{R}^q$ such that $x_i \approx \mu + \mathbf{V}_q \lambda_i$ and $\mathbf{V}_q^T \mathbf{V}_q = \mathbf{I}$

$$\min_{\mu, \mathbf{V}_q, \{\lambda_i\}_i} \sum_{i=1}^n \|x_i - \mu - \mathbf{V}_q \lambda_i\|_2^2$$

Linear projections

Given $x_1, \dots, x_n \in \mathbb{R}^d$, for $q \ll d$ find a compressed representation with $\lambda_1, \dots, \lambda_n \in \mathbb{R}^q$ such that $x_i \approx \mu + \mathbf{V}_q \lambda_i$ and $\mathbf{V}_q^T \mathbf{V}_q = \mathbf{I}$

$$\min_{\mu, \mathbf{V}_q, \{\lambda_i\}_i} \sum_{i=1}^n \|x_i - \mu - \mathbf{V}_q \lambda_i\|_2^2$$

Fix \mathbf{V}_q and solve for μ, λ_i :

$$\mu = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\lambda_i = \mathbf{V}_q^T (x_i - \bar{x})$$

Which gives us:

$$\min_{\mathbf{V}_q} \sum_{i=1}^n \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|^2.$$

$\mathbf{V}_q \mathbf{V}_q^T$ is a *projection matrix* that minimizes error in basis of size q

Linear projections

$$\mathbf{z}^T \Sigma \mathbf{z} \geq 0 \quad \forall \mathbf{z}$$

$$\sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|_2^2$$

$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\mathbf{V}_q^T \mathbf{V}_q = I_q$$

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|_2^2 = \min_{\mathbf{V}_q} \underline{\text{Tr}(\Sigma) - \text{Tr}(\mathbf{V}_q^T \Sigma \mathbf{V}_q)}$$

Eigenvalue decomposition of $\Sigma =$ V_2 : maximizes $\text{Tr}(V_2^T \Sigma V_2)$

Linear projections

$$\sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|_2^2$$

$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$
$$\mathbf{V}_q^T \mathbf{V}_q = I_q$$

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|_2^2 = \min_{\mathbf{V}_q} \text{Tr}(\Sigma) - \text{Tr}(\mathbf{V}_q^T \Sigma \mathbf{V}_q)$$

Eigenvalue decomposition of $\Sigma =$

\mathbf{V}_q are the first q eigenvectors of Σ

Minimize reconstruction error and capture the most variance in your data.

Pictures

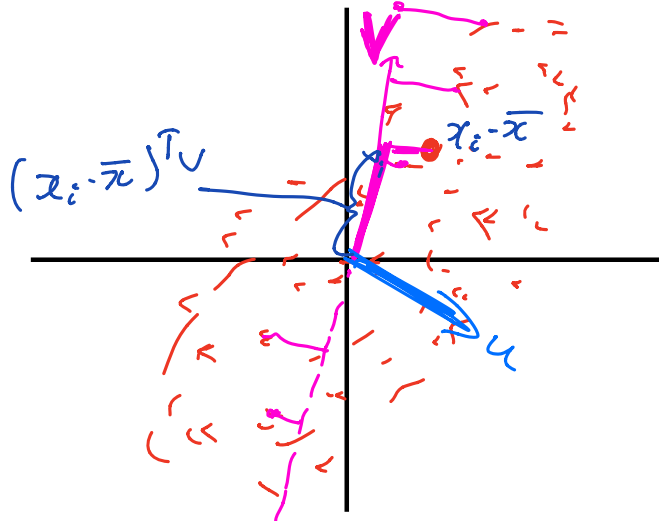
\mathbf{V}_q are the first q eigenvectors of Σ

\mathbf{V}_q are the first q principal components

$$\Sigma := \sum_{i=1}^N \underbrace{(x_i - \bar{x})(x_i - \bar{x})^T}$$

\mathbf{V}_q with $\mathbf{V}_q^T \mathbf{V}_q = I$ maximizes $\text{Tr}(\mathbf{V}_q^T \Sigma \mathbf{V}_q)$

$x_i - \bar{x}$



$q=1 \Rightarrow v_1 \in \mathbb{R}^d$

$$\sum_i v^T (x_i - \bar{x})(x_i - \bar{x})^T v = v^T \Sigma v$$

$$v^T \Sigma v > u^T \Sigma u$$

Pictures

\mathbf{V}_q are the first q eigenvectors of Σ

\mathbf{V}_q are the first q principal components

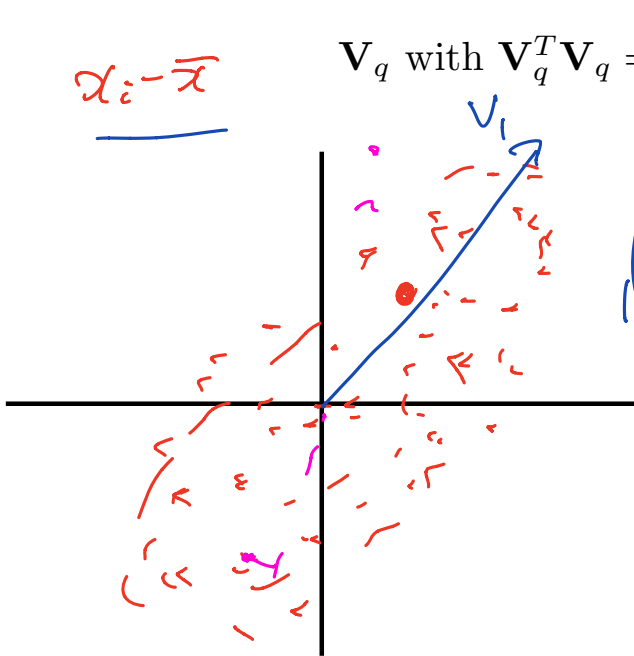
$$\Sigma := \sum_{i=1}^N \underline{(x_i - \bar{x})(x_i - \bar{x})^T}$$

\mathbf{V}_q with $\mathbf{V}_q^T \mathbf{V}_q = I$ maximizes $\text{Tr}(\mathbf{V}_q^T \Sigma \mathbf{V}_q)$

$q=2 \Rightarrow V_2 \in \mathbb{R}^{d \times 2}$

$$= \sum_i V_1^T (x_i - \bar{x})(x_i - \bar{x})^T V_1$$

$$+ \sum_i V_2^T (x_i - \bar{x})(x_i - \bar{x})^T V_2$$



Linear projections

Given $x_i \in \mathbb{R}^d$ and some $q < d$ consider

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|^2.$$

where $\mathbf{V}_q = [v_1, v_2, \dots, v_q]$ is orthonormal:

$$\mathbf{V}_q^T \mathbf{V}_q = I_q$$

\mathbf{V}_q are the first q eigenvectors of Σ

\mathbf{V}_q are the first q principal components

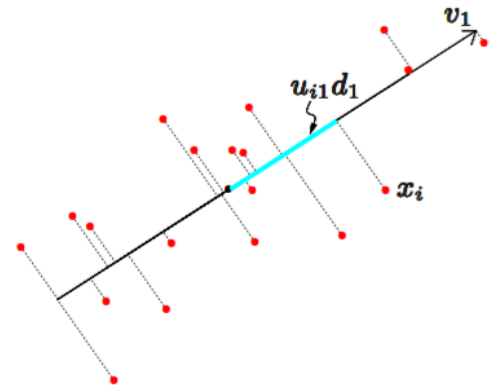
Principal Component Analysis (PCA) projects $(\mathbf{X} - \mathbf{1}\bar{x}^T)$ down onto \mathbf{V}_q

$$(\mathbf{X} - \mathbf{1}\bar{x}^T) \mathbf{V}_q = \mathbf{U}_q \text{diag}(d_1, \dots, d_q)$$



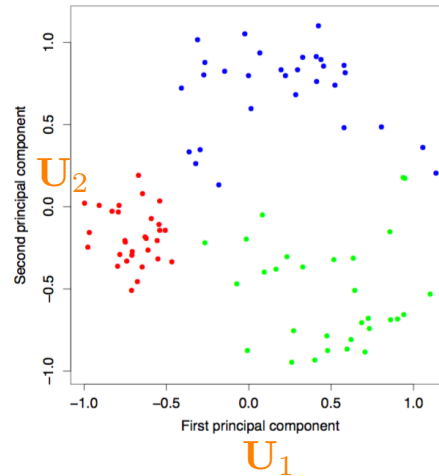
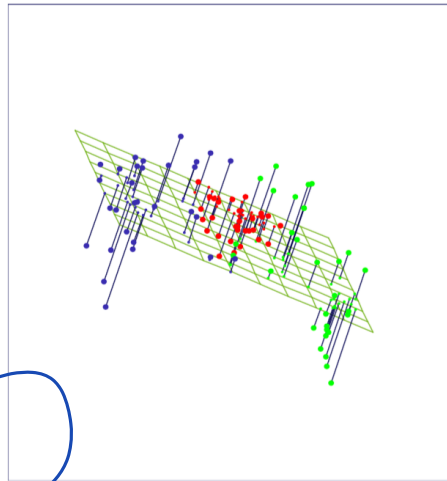
$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\mathbf{U}_q^T \mathbf{U}_q = I_q$$

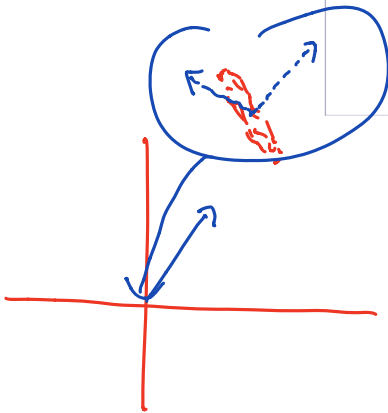


Dimensionality reduction

V_q are the first q eigenvectors of Σ and SVD



$$X - \mathbf{1}\bar{x}^T$$



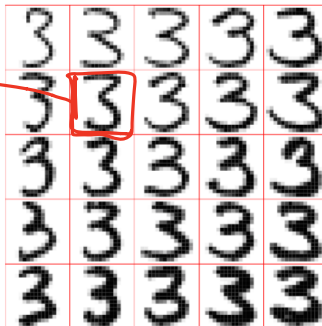
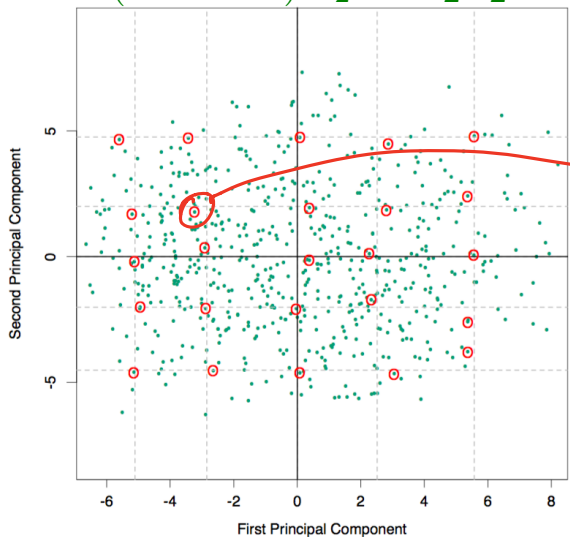
Dimensionality reduction

V_q are the first q eigenvectors of Σ and SVD $\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\mathbf{S}\mathbf{V}^T$

Handwritten 3's, 16x16 pixel image so that $x_i \in \mathbb{R}^{256}$

$$\begin{aligned}\hat{f}(\lambda) &= \bar{x} + \lambda_1 v_1 + \lambda_2 v_2 \\ &= \text{3} + \lambda_1 \cdot \text{3} + \lambda_2 \cdot \text{3}.\end{aligned}$$

$$(\mathbf{X} - \mathbf{1}\bar{x}^T)\mathbf{V}_2 = \mathbf{U}_2\mathbf{S}_2 \in \mathbb{R}^{n \times 2}$$



$\text{diag}(\mathbf{S})$

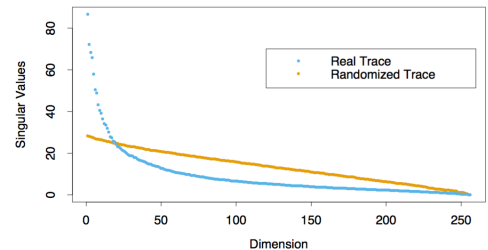


FIGURE 14.24. The 256 singular values for the digitized threes, compared to those for a randomized version of the data (each column of \mathbf{X} was scrambled).

Singular Value Decomposition (SVD)

$$e_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ with } 1 \text{ at } i\text{th pos}$$

Theorem (SVD): Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $\text{rank } r \leq \min\{m, n\}$. Then $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ where $\mathbf{S} \in \mathbb{R}^{r \times r}$ is diagonal with positive entries, $\mathbf{U}^T\mathbf{U} = \mathbf{I}$, $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

$$\mathbf{V} = [v_1, \dots, v_n]$$

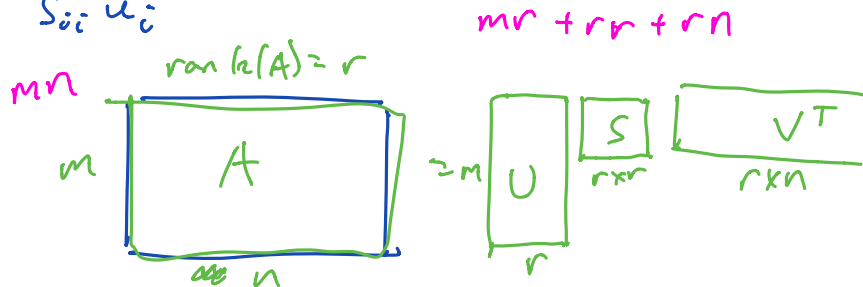
$$[\mathbf{V}^T\mathbf{V}]_{ij} = v_i^T v_j$$

$$\mathbf{A}^T \mathbf{A} v_i = (\mathbf{U}\mathbf{S}\mathbf{V}^T)^T (\mathbf{U}\mathbf{S}\mathbf{V}^T) v_i = \mathbf{V} \mathbf{S} \mathbf{U}^T \mathbf{U} \mathbf{S} \mathbf{V}^T v_i = \mathbf{V} \mathbf{S}^2 \mathbf{V}^T v_i = \mathbf{V} \mathbf{S}^2 e_i$$

$$\mathbf{U} = [u_1, \dots, u_n]$$

$$\mathbf{A} \mathbf{A}^T u_i = \mathbf{U} \mathbf{S} \mathbf{V}^T (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T u_i = \mathbf{U} \mathbf{S}^2 \mathbf{U}^T u_i = \mathbf{S}^2 u_i$$

$$= s_{ii}^2 v_i$$



Singular Value Decomposition (SVD)

Theorem (SVD): Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with rank $r \leq \min\{m, n\}$. Then $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ where $\mathbf{S} \in \mathbb{R}^{r \times r}$ is diagonal with positive entries, $\mathbf{U}^T\mathbf{U} = \mathbf{I}$, $\mathbf{V}^T\mathbf{V} = \mathbf{I}$.

$$\mathbf{A}^T\mathbf{A}v_i = \mathbf{S}_{i,i}^2v_i$$

$$\mathbf{A}\mathbf{A}^T u_i = \mathbf{S}_{i,i}^2u_i$$

\mathbf{V} are the first r eigenvectors of $\mathbf{A}^T\mathbf{A}$ with eigenvalues $\text{diag}(\mathbf{S})$

\mathbf{U} are the first r eigenvectors of $\mathbf{A}\mathbf{A}^T$ with eigenvalues $\text{diag}(\mathbf{S})$

Linear projections

\mathbf{V}_q are the first q eigenvectors of Σ

\mathbf{V}_q are the first q principal components

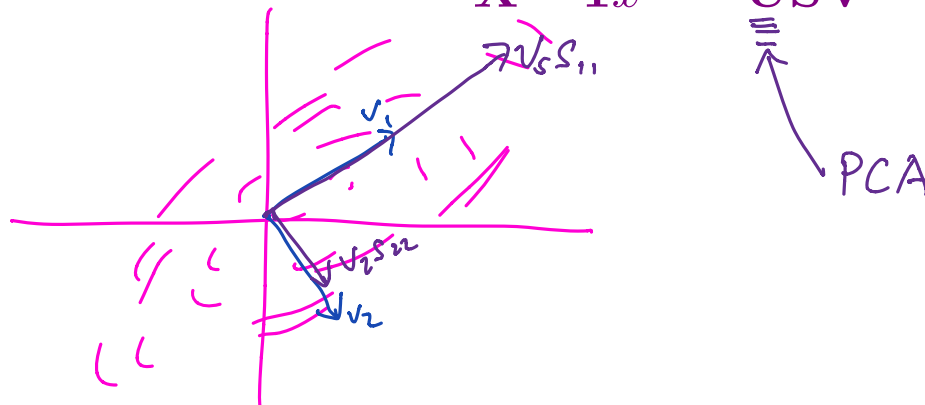
$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

Principal Component Analysis (PCA) projects $(\mathbf{X} - \mathbf{1}\bar{x}^T)$ down onto \mathbf{V}_q


$$(\mathbf{X} - \mathbf{1}\bar{x}^T)\mathbf{V}_q = \mathbf{U}_q \text{diag}(d_1, \dots, d_q) \quad \mathbf{U}_q^T \mathbf{U}_q = I_q$$

Singular Value Decomposition defined as

$$\overset{n \times d}{\mathbf{X} - \mathbf{1}\bar{x}^T} = \overset{n \times 2}{\mathbf{U}} \overset{2 \times 2}{\Sigma} \overset{2 \times d}{\mathbf{V}^T}$$



Pictures, intuition!

$$\begin{aligned}
 \mathbf{J}X &= X - \frac{1}{n} \mathbf{1} \mathbf{1}^T X / n \\
 &= \sum_{i=1}^n x_i^T / n \\
 &= X - \mathbf{1} \bar{x}
 \end{aligned}$$


- Fill in the missing plots: $\mathbf{U}, \mathbf{S}, \mathbf{V} = \text{svd}(\mathbf{J}X)$

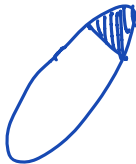
$$\mathbf{J} = \mathbf{I} - \mathbf{1} \mathbf{1}^T / n$$

$$\mathbf{J}X = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

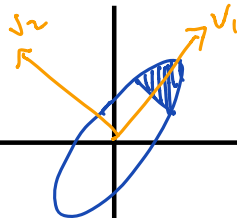
$$\mathbf{V} = [v_1, v_2]$$

$$= \mathbf{U} \mathbf{S} \mathbf{V}^T \mathbf{V} \mathbf{S}^{-1} = \mathbf{U}$$

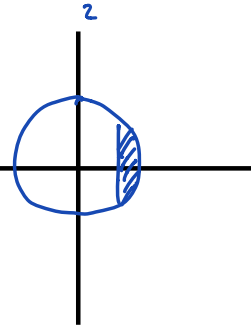
\mathbf{X}



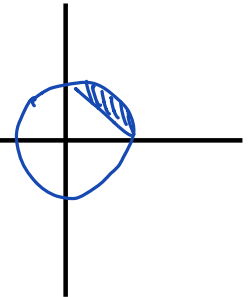
$\mathbf{J}X$



$\mathbf{J}X \mathbf{V} \mathbf{S}^{-1}$



$\mathbf{J}X \mathbf{V} \mathbf{S}^{-1} \mathbf{V}^T$



$$(\mathbf{J}X \mathbf{V})_{i,1} = (x_i - \bar{x})^T v_1$$

Kernel PCA

\mathbf{V}_q are the first q eigenvectors of Σ and SVD $\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\mathbf{S}\mathbf{V}^T$

$$(\mathbf{X} - \mathbf{1}\bar{x}^T)\mathbf{V}_q = \mathbf{U}_q\mathbf{S}_q \in \mathbb{R}^{n \times q}$$

$$\underline{\mathbf{J}\mathbf{X} = \mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\mathbf{S}\mathbf{V}^T} \quad \mathbf{J} = \mathbf{I} - \mathbf{1}\mathbf{1}^T/n$$

$$(\mathbf{J}\mathbf{X})(\mathbf{J}\mathbf{X})^T = \mathbf{U}\mathbf{S}\mathbf{V}^T\mathbf{V}\mathbf{S}\mathbf{U}^T = \mathbf{U}\mathbf{S}^2\mathbf{U}^T \in \mathbb{R}^{n \times n}$$

$$(\mathbf{J}\mathbf{X})^T(\mathbf{J}\mathbf{X}) = \mathbf{V}\mathbf{S}\mathbf{U}^T\mathbf{U}\mathbf{S}\mathbf{V}^T = \mathbf{V}\mathbf{S}^2\mathbf{V}^T \in \mathbb{R}^{d \times d}$$

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

PCA Algorithm

PCA

input

A matrix of m examples $X \in \mathbb{R}^{m,d}$

number of components n

if ($m > d$)

$$A = X^T X$$

Let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be the eigenvectors of A with largest eigenvalues

else

$$B = X X^T$$

Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be the eigenvectors of B with largest eigenvalues

for $i = 1, \dots, n$ set $\mathbf{u}_i = \frac{1}{\|X^T \mathbf{v}_i\|} X^T \mathbf{v}_i$

output: $\mathbf{u}_1, \dots, \mathbf{u}_n$



Cool tricks with SVD

Machine Learning – CSE546

Kevin Jamieson

University of Washington

November 8, 2016

Ridge Regression revisited

$$\hat{w}_{ridge} = \arg \min_w \|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda \|w\|_2^2$$

$$\hat{w}_{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y} \quad (\text{Assume data centered})$$

Singular vector decomposition (SVD): $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \underline{\mathbf{U} \mathbf{S} \mathbf{V}^T}$

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

Ridge Regression revisited

$$\hat{w}_{ridge} = \arg \min_w \|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda \|w\|_2^2$$

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

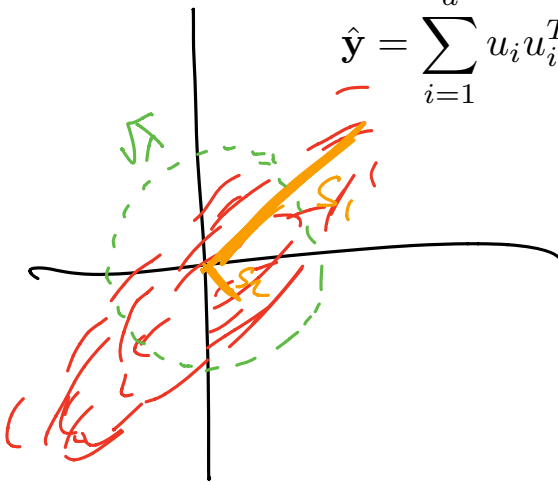
Singular vector decomposition (SVD): $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$

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

$$\hat{\mathbf{y}} = \sum_{i=1}^d u_i u_i^T \frac{s_i^2}{s_i^2 + \lambda} y_i$$

$$\mathbf{U} = [u_1, \dots, u_d]$$

$$\mathbf{S} = \text{diag}(s_1, \dots, s_d)$$



$$\frac{s_i^2}{s_i^2 + \lambda} \approx \begin{cases} 0 & s_i^2 \ll \lambda \\ 1 & s_i^2 \gg \lambda \end{cases}$$