

# Announcements



- Google form feedback <https://tinyurl.com/yb2tprkl>



# The previous and future weeks

Machine Learning – CSE546

Kevin Jamieson

University of Washington

November 9, 2017

# So far...



Supervised learning:  $x_i \in \mathbb{R}^d$   $y_i \in \mathbb{R}$  for  $i = 1, \dots, n$ . Learn  $f : x \rightarrow y$

Loss functions:  $LA = \sum_{i=1}^n \ell(f(x_i), y_i)$

$$\ell(y, f(x)) = (f(x) - y)^2 \quad \max(0, 1 - yf(x))$$
$$= y \log(f(x)) + (1 - y) \log(1 - f(x)) \Leftrightarrow \log(1 + \exp(-yf(x)))$$

$\frac{\partial}{\partial yf(x)}$

Methods:

Linear } SUM  
- Lasso  
- Ridge

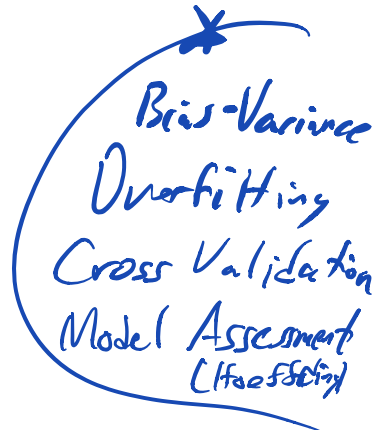
Trees, Boosting, Trees

Nearest Neighbor  
Kernel Machines

MLE, MAP

Neural Networks

Bagging



# Method comparison

**TABLE 10.1.** *Some characteristics of different learning methods. Key: ▲ = good, ◆ = fair, and ▼ = poor.*

Characteristic	Neural Nets	SVM	Trees	Boosting Trees	k-NN, Kernels
Natural handling of data of “mixed” type	▼	▼	▲	▲	▼
Handling of missing values	▼	▼	▲	▲	▲
Robustness to outliers in input space	▼	▼	▲	▼	▲
Insensitive to monotone transformations of inputs	▼	▼	▲	▼	▼
Computational scalability (large $N$ )	▼	▼	▲	▲	▼
Ability to deal with irrelevant inputs	▼	▼	▲	▲	▼
Ability to extract linear combinations of features	▲	▲	▼	▼	◆
Interpretability	▼	▼	◆	▲	▼
Predictive power	▲	▲	▼	◆	▲

# To come



- Unsupervised learning
  - SVD
  - Clustering
  - Density estimation
- Machine learning street fighting tools
  - Tips, tricks, data pre-processing, output post-processing
  - Domain specific data (images, sequences)
- Reinforcement learning
- Learning theory



# Principle Component Analysis

Machine Learning – CSE546

Kevin Jamieson

University of Washington

November 9, 2017

# Linear projections

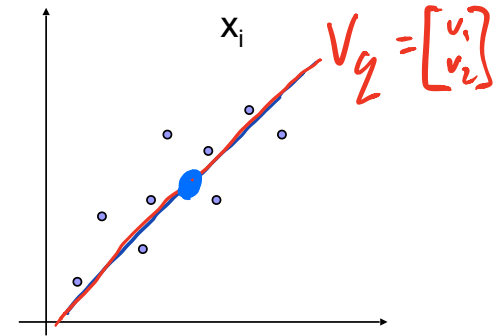
Given  $x_i \in \mathbb{R}^d$  and some  $q < d$  consider

$$\min_{\mu, \{\lambda_i\}, \mathbf{V}_q} \sum_{i=1}^N \|\underline{x}_i - \underline{\mu} - \mathbf{V}_q \lambda_i\|^2.$$

where  $\lambda_i \in \mathbb{R}^q$  and  $\mathbf{V}_q = [v_1, v_2, \dots, v_q]$  is orthonormal:

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

$\mathbf{V}_q \in \mathbb{R}^{d \times q}$



# Linear projections

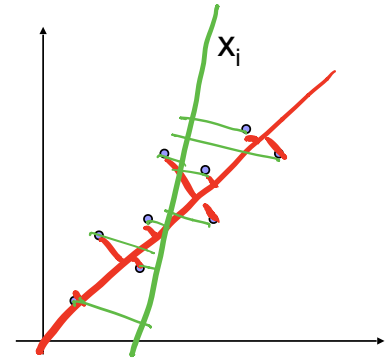
Given  $x_i \in \mathbb{R}^d$  and some  $q < d$  consider

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

where  $\lambda_i \in \mathbb{R}^q$  and  $\mathbf{V}_q = [v_1, v_2, \dots, v_q]$  is orthonormal:

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

Natural choices for  $\mu, \lambda_i$  ?





# Linear projections

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

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

where  $\lambda_i \in \mathbb{R}^q$  and  $\mathbf{V}_q = [v_1, v_2, \dots, v_q]$  is orthonormal:

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

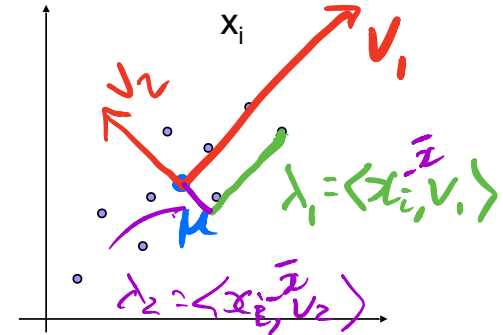
Natural choices for  $\mu, \lambda_i$  ?

$$\hat{\mu} = \bar{x}, \quad \hat{\lambda}_i = \mathbf{V}_q^T (x_i - \bar{x}).$$

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

Which gives us:

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



$$\hat{x}_i = \bar{x} + d_i v_i$$

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

# Linear projections

$\text{Tr}(V^T A V) = \text{Tr}(A V V^T)$   
 $\text{Tr}((A+B)C) = \text{Tr}(AC) + \text{Tr}(BC)$   
 $= I - V_2 V_2^T$

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

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

$$V_q^T V_q = I_q$$

$$\Rightarrow \|(I - V_2 V_2^T)(x_i - \bar{x})\|_2^2$$

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

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

$$= \sum_i \text{Tr}((I - V_2 V_2^T) (x_i - \bar{x})(x_i - \bar{x})^T)$$

$$= \sum_i \text{Tr}((x_i - \bar{x})(x_i - \bar{x})^T) - \text{Tr}(V_2 V_2^T (x_i - \bar{x})(x_i - \bar{x})^T)$$

$$= \text{Tr}(\Sigma) - \sum_{i=1}^n \text{Tr}(V_2^T (x_i - \bar{x})(x_i - \bar{x})^T V_2)$$

$$= \text{Tr}(\Sigma) - \text{Tr}(V_2^T \Sigma V_2)$$

# Linear projections

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

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

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

$\lambda, v$  are eigenvalue/vector pair if

$$\Sigma v = \lambda v$$

$$\tilde{\Lambda} = \begin{bmatrix} \lambda_1 & & 0 \\ 0 & \lambda_2 & \\ & & \ddots \\ 0 & & & 0 \end{bmatrix}$$

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

$$\Sigma = \tilde{\mathbf{V}} \tilde{\Lambda} \tilde{\mathbf{V}}^T$$

$$\max_{\mathbf{V}_q} \text{Tr}(\mathbf{V}_q^T \tilde{\mathbf{V}} \tilde{\Lambda} \tilde{\mathbf{V}}^T \mathbf{V}_q)$$

$$\Rightarrow \mathbf{V}_q = \tilde{\mathbf{V}}_q$$

# 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  $\Sigma$

*with the largest  $q$  eigenvalues*

# 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:

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

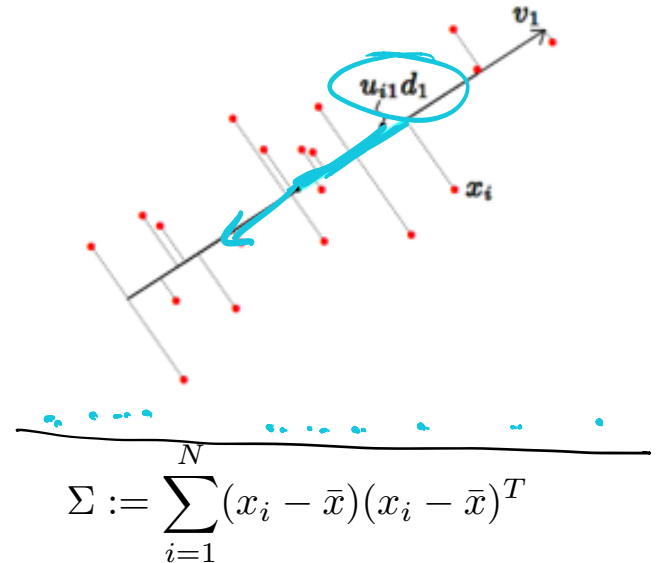
$\mathbf{V}_q$  are the first  $q$  eigenvectors of  $\Sigma$

$\mathbf{V}_q$  are the first  $q$  principle components

Principle component Analysis (PCA) projects  $(\mathbf{X} - \mathbf{1}\bar{x})$  down onto  $\mathbf{V}_q$

$$(\mathbf{X} - \mathbf{1}\bar{x}^T) \mathbf{V}_q = \mathbf{U}_q \mathbf{d} \mathbf{I} \mathbf{a} \mathbf{g}(d_1, \dots, d_q)$$

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



# Linear projections

$$\bar{\mathbf{X}} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_n \end{bmatrix}$$

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

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

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  $\Sigma$

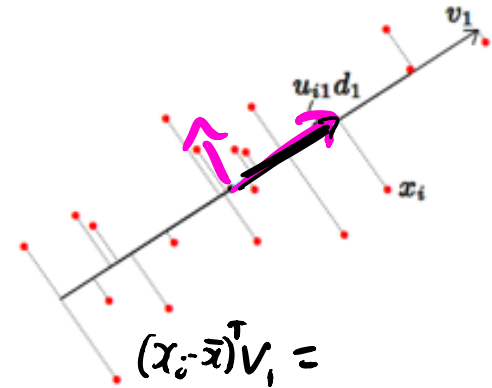
$\mathbf{V}_q$  are the first  $q$  principle components

Principle component Analysis (PCA) projects  $(\mathbf{X} - \mathbf{1}\bar{x})$  down onto  $\mathbf{V}_q$  (if  $d < n$ )

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

Singular Value Decomposition defined as

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



$$(x_i - \bar{x})^T \mathbf{V}_1 =$$

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

$$\begin{aligned} \mathbf{U} &\in \mathbb{R}^{n \times d} \\ \mathbf{S} &\in \mathbb{R}^{d \times d} \\ \mathbf{V} &= \mathbb{R}^{d \times d} \end{aligned}$$

# 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  $\Sigma$

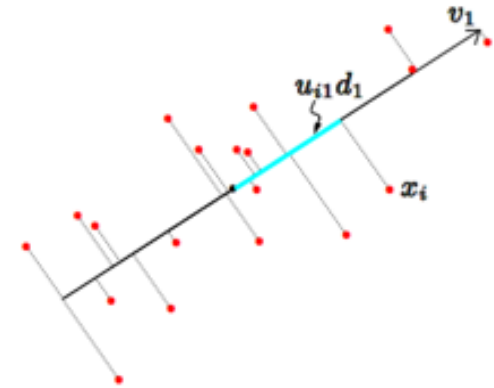
$\mathbf{V}_q$  are the first  $q$  principle components

Principle component Analysis (PCA) projects  $(\mathbf{X} - \mathbf{1}\bar{x})$  down onto  $\mathbf{V}_q$

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

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

$$\mathbf{U}_q^T \mathbf{U}_q = I_q \quad \Sigma = \sum_{i=1}^q v_i v_i^T d_i^2$$



Singular Value Decomposition defined as

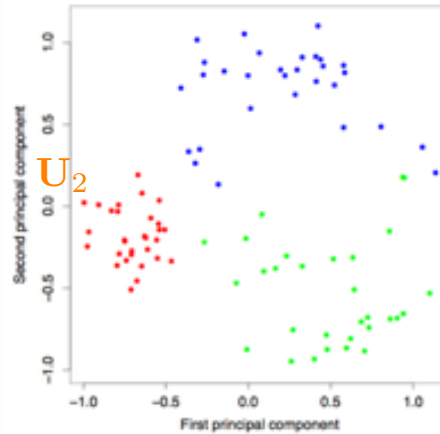
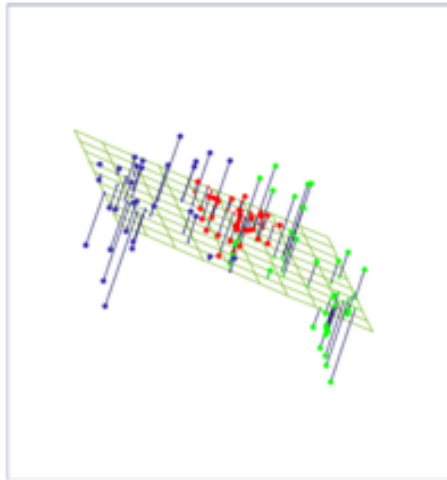
$$\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

How do the *eigenvalues* of  $\Sigma$  relate to the *singular values* of  $\mathbf{X} - \mathbf{1}\bar{x}^T$ ?

# $A$ is singular if $\exists x \neq 0: Ax = 0$ Dimensionality reduction

$V_q$  are the first  $q$  eigenvectors of  $\Sigma$  and SVD  $\underline{X - 1\bar{x}^T} = \underline{USV^T}$

$$(AB)^T = B^T A^T$$



singular values  
of  $X - 1\bar{x}^T$

eigenvalues of  $\Sigma$   
are singular values  
squared.

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

$U_1$

$$\Sigma = \overbrace{(X - 1\bar{x}^T)^T (X - 1\bar{x}^T)} = \underbrace{V^T}_{\Rightarrow} \underbrace{S^T U^T U S}_{I} V^T = V S^2 V^T$$

eigenvalues of  $\Sigma$



# Dimensionality reduction

$V_q$  are the first  $q$  eigenvectors of  $\Sigma$  and SVD  $X - 1\bar{x}^T = USV^T$

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

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

$$(X - 1\bar{x}^T)V_2 = U_2S_2 \in \mathbb{R}^{n \times 2}$$

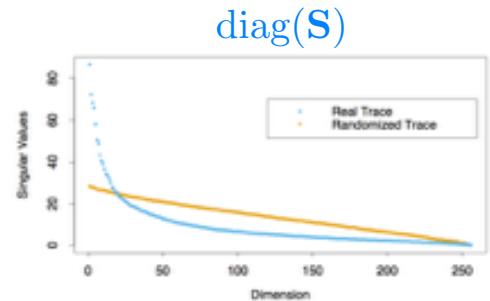
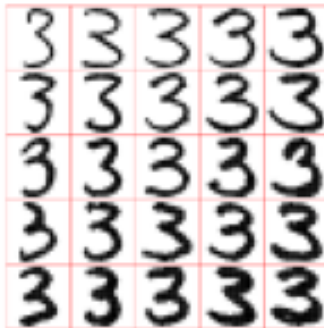
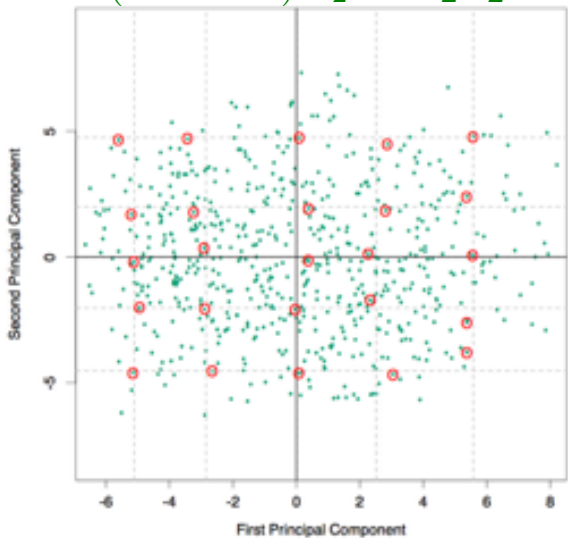


FIGURE 14.24. The 256 singular values for the digitized threes, compared to those for a randomized version of the data (each column of  $X$  was scrambled).

# Kernel PCA

$V_q$  are the first  $q$  eigenvectors of  $\Sigma$  and SVD  $\underline{X - 1\bar{x}^T = USV^T}$

$$(X - 1\bar{x}^T)V_q = \underline{U_q S_q} \in \mathbb{R}^{n \times q}$$

$$\bar{x}^T = \mathbf{1}^T X / n$$

$$\underline{JX = X - 1\bar{x}^T = USV^T} \quad \underline{J = I - \mathbf{1}\mathbf{1}^T / n}$$

$$\begin{aligned} (JX)(JX)^T &= \underline{J} \underline{X X^T} \underline{J} = (X - \mathbf{1}\bar{x}^T)(X - \mathbf{1}\bar{x}^T)^T \\ &= USV^T V S U^T \\ &= U S^2 U^T \end{aligned}$$

# Kernel PCA

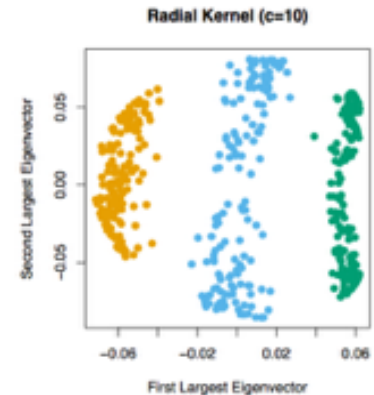
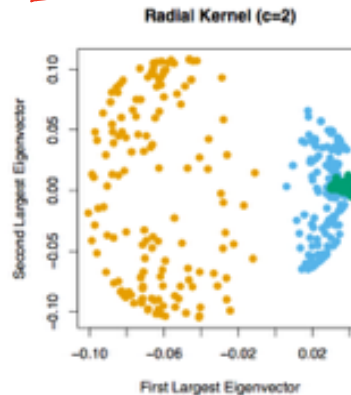
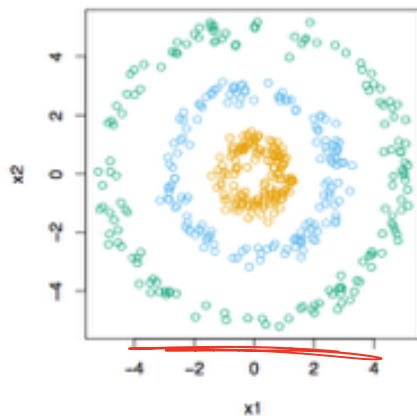
$V_q$  are the first  $q$  eigenvectors of  $\Sigma$  and SVD  $\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\mathbf{S}\mathbf{V}^T$

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

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

$$K_{ij} = e^{-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}}$$

$$(\mathbf{J}\mathbf{X})(\mathbf{J}\mathbf{X})^T = \mathbf{U}\mathbf{S}^2\mathbf{U}^T = \mathbf{J}\mathbf{K}\mathbf{J}$$



# 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:**  $\mathbf{u}_1, \dots, \mathbf{u}_n$

# Ridge Regression revisited

$(AB)^T = B^T A^T$      $V^T V = I$      $V^T V = I$   
 $V^{-1} = V^T$

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

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

Assume  $\mathbf{X}$  is centered

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

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

$$= \mathbf{U} \mathbf{S} \mathbf{V}^T (\mathbf{V} \mathbf{S}^2 \mathbf{V}^T + \lambda I)^{-1} \mathbf{V} \mathbf{S} \mathbf{U}^T \mathbf{y}$$

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

$$= \mathbf{U} \mathbf{S} \mathbf{V}^T (\mathbf{V} (\mathbf{S}^2 + \lambda I) \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{S} \mathbf{U}^T$$

$$= \mathbf{U} \mathbf{S} \mathbf{V}^T (\mathbf{V}^{-1} (\mathbf{S}^2 + \lambda I)^{-1} \mathbf{V}^{-1}) \mathbf{V} \mathbf{S} \mathbf{U}^T$$

$$= \mathbf{U} \mathbf{S} \mathbf{V}^T \mathbf{V} (\mathbf{S}^2 + \lambda I)^{-1} \mathbf{V}^T \mathbf{V} \mathbf{S} \mathbf{U}^T$$

$$= \mathbf{U} \mathbf{S} (\mathbf{S}^2 + \lambda I)^{-1} \mathbf{S} \mathbf{U}^T = \sum_{i=1}^d u_i u_i^T \frac{s_i^2}{s_i^2 + \lambda}$$

# Ridge Regression revisited

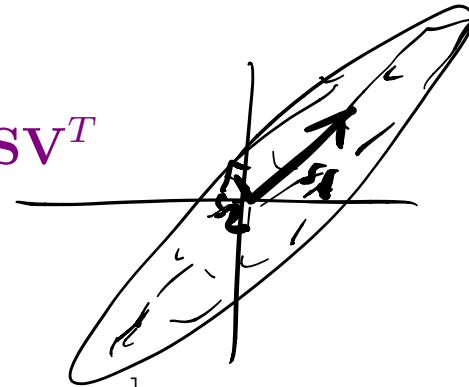
$$\hat{w}_{ridge} = \arg \min_w \|\mathbf{X}w - \mathbf{y}\|_2^2 + \lambda \|w\|_2^2$$

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

Singular vector decomposition (SVD):  $\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\mathbf{S}\mathbf{V}^T$

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

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



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

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