

$$y = \begin{pmatrix} y_A \\ y_B \end{pmatrix} \sim N(\mu, \Sigma), \mu = \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}, \Sigma = \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix}$$

$$y \sim N(\mu, \Sigma), y' \sim N(\mu', \Sigma') \Rightarrow y + y' \sim N(\mu + \mu', \Sigma + \Sigma')$$

$$x \sim N(\mu, \sigma^2)$$

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^n |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu))$$

$$p(x) = \frac{1}{\sqrt{2\pi n}} \exp(-\frac{(x - \mu)^2}{2\sigma^2})$$

$$\mathbb{E}(x)^T \mathbb{E}(x)$$

Hongyi Zhang $y_A \sim N(\mu_A, \Sigma_{AA}), y_B \sim N(\mu_B, \Sigma_{BB})$
 $y_A | y_B \sim N(\mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (y_B - \mu_B), \Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA})$

Performance Evaluation

$$\Sigma_{AA} - \Sigma_{AB} \Sigma_{BB}^{-1} \Sigma_{BA}$$

2.1 Performance Measures for Regression

$$2.1.1 \text{ Error Mean absolute error (MAE)} \quad e = \frac{1}{n} \sum_{i=1}^n |f(x_i) - y_i|$$

$$\text{Mean square error (MSE)} \quad e = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$2.1.2 \quad R^2 = 1 - \frac{SS_{RES}}{SS_{TOT}} = \frac{SS_{REG}}{SS_{TOT}} = \frac{1}{n} SS_{REG}$$

$$SS_{RES} = \sum_{i=1}^n (f(x_i) - y_i)^2 \quad SS_{REG} = \sum_{i=1}^n (f(x_i) - \bar{y})^2 \quad SS_{TOT} = \sum_{i=1}^n (y_i - \bar{y})^2$$

2.1.3 Negative log predictive density

$$NLPD = -\log p(y) = -\frac{1}{n} \sum_{i=1}^n \log P(y_i = f(x_i) | x_i)$$

$$\text{Gaussian Process: } NLPD = \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} \log(2\pi\sigma^2) + \frac{(y_i - \mu_i)^2}{2\sigma^2} \right)$$

2.2 Performance Measures for Classification

$$2.2.1 \text{ 0/1 loss: } \frac{1}{n} \sum_{i=1}^n \delta_{h(x_i) \neq y_i} \times 100\% \quad (\text{error})$$

$$1\text{-error: } \frac{1}{n} \sum_{i=1}^n \delta_{h(x_i) \neq y_i} = y_i \times 100\% \quad (\text{accuracy})$$

2.2.2 Confusion Matrix

		+1	-1
prediction	+1	TP	FP
	-1	FN	TN

$$P = TP + FN; \quad N = TN + FP$$

$$FPR = \frac{FP}{N} \quad FNR = 1 - TPR = \frac{FN}{P} \quad TNR = \frac{TN}{N} \quad TPR = \frac{TP}{P} \quad (\text{recall})$$

$$\text{precision} = \frac{TP}{TP+FP} \quad \text{accuracy} = \frac{TP+TN}{P+N} \quad F\text{-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

(only for binary) (change)

2.2.3 Area under the ROC curve (AUC) – threshold parameter t

$X \neq y$: $FPR, Y \neq y$: TPR AUC = probability that a randomly selected positive example is ranked higher than a randomly selected negative example

2.4 Efficiency Runtime (training time / test time) Memory

3. Null hypothesis H_0 : A and B perform equivalently Aim: reject H_0

3.1 Parametric Tests t-Test split 1...n

H_0 : the observed means μ_A and μ_B are the same (same performance)

$$d_i = \bar{a}_{i1}^A - \bar{a}_{i1}^B, \quad d = \mu_A - \mu_B, \quad \hat{\sigma}_d = \sqrt{\frac{\sum (d_i - \bar{d})^2}{n-1}}, \quad t = \frac{\bar{d}}{\hat{\sigma}_d / \sqrt{n}}$$

degree of freedom $df = n-1$; two-tailed test $|t| \geq t_{\alpha/2}$ ($\alpha = 0.05$)

we can reject H_0

3.2 Non-parametric Tests 3.2.1 McNemar's Test (classification)

C_{01} the number of instances misclassified by A correctly classified by B

C_{10} missclassified by B correctly classified by A

$H_0: p(C_{01}) = p(C_{10})$ (same performance)

$$\chi^2_{MC} = \frac{(C_{01} - C_{10} - 1)^2}{C_{01} + C_{10}}$$

Chi-squared distribution with 1 degree of freedom

3.2.2 Sign Test (used on multiple domains)

A_{win} = number of experiments A outperform B. B_{win} = B outperformed A

$$H_0: p(A_{win} > B_{win}) = 0.5 \quad (\text{same performance}) \quad A_{win} \sim B(n, 0.5)$$

$A_{win} > n_{\alpha}$, reject H_0

$$\text{ARD Kernel } K(x, x') = \sigma^2 \exp\left(-\frac{1}{2}(x-x')^T M(x-x') + \sigma^2 \delta_{xx'}\right) \quad M = \text{diag}(\ell)^{-2}$$

Kernel Methods 不需要计算 $\phi(x)$, 只需计算 $\langle \phi(x), \phi(x') \rangle$ After 2. Kernel Trick

$$2.1 \text{ linear model } h(x) = w^T x, \quad (w) = \sum_{i=1}^n (w^T x_i - y_i)^2, \quad g(w) = \frac{\partial L}{\partial w} = \sum_{i=1}^n 2(w^T x_i - y_i)x_i$$

$$w_{\text{test}} = w - sg(w) \rightarrow w_t = \sum_{i=1}^n \alpha_i^{(w)} x_i \rightarrow w_t^T x_i = \sum_{j=1}^n \alpha_j^{(w)} x_j^T x_i \Rightarrow \text{new test}$$

$$y^* = \sum_{i=1}^n \alpha_i^{(w)} x_i^T x^* \xrightarrow{\text{kernel}} y^* = \sum_{i=1}^n \alpha_i^{(w)} \langle \phi(x_i), \phi(x^*) \rangle \quad \phi(x) = \sum_{i=1}^n \alpha_i x_i$$

Definition: Kernel $K(x, x') = \langle \phi(x), \phi(x') \rangle$ Definition: Valid Kernel $\in \mathbb{R}^{n \times n}, \{x_1, \dots, x_m\}$

Dsymmetric $K(x, x') = K(x', x)$ ② positive semi-definite K is a PSD matrix
 $\forall q \in \mathbb{R}^m, q^T A q \geq 0 \quad \forall v \in \mathbb{R}^n, \forall \lambda \geq 0 \quad \lambda I \geq 0 \quad K_{ij} = K(x_i, x_j)$

3.2. Common Kernel Linear Kernel $K(x, x') = x^T x'$ Polynomial Kernel \bar{d} , degree of poly $K(x, x') = (1 + x^T x')^{\bar{d}}$ RBF/Gaussian Kernel $K(x, x') = \exp(-\frac{1}{2\sigma^2} \|x - x'\|^2)$ $\forall i, j, K_{ij} \geq 0$

3.4 Kernel Construction ③ $K(x, x') = x^T x' \quad ④ K(x, x) = C \geq 0 \quad ⑤ K(x, x') = K_1 + K_2$
 $f, g \in \mathbb{R}^n \quad ⑥ K = g(x), g \text{ is polynomial function} \quad ⑦ K = K_1 K_2 \quad ⑧ K = f(x) K_1 f(x') \quad ⑨ K = e^{K_1}$

⑩ $K(x, x') = x^T A x'$, A is PSD ⑪ Kernel Machines

4.1 Kernel Ridge Regression $\min_w \|Xw - y\|^2 + \lambda \|w\|^2 \rightarrow w = (X^T X + \lambda I)^{-1} X y$

$w = X\alpha, \alpha = (X^T X + \lambda I)^{-1} y = (K + \lambda I)^{-1} y$, where $K = X^T X, \alpha \in \mathbb{R}^n$

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

$y = w^T x + \varepsilon \Rightarrow y = w^T \phi(x) + \varepsilon \text{ or } y = K x \alpha, \alpha = (K + \lambda I)^{-1} y, K_x = K(x, X)$

$y = w^T x + \varepsilon$ (Primal) $y = x^T \alpha$ (Dual) Gaussian Process (只对 ε 有噪音)

$p(y|X, w) = N(X^T w, \sigma^2 I), \quad p(y|x, w) = N(w^T x, \sigma^2)$

$p(w|X, y) = \frac{p(y|X, w)p(w)}{p(y|X)} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \quad p(w) = N(0, \sigma_w^2 I)$

$y = f(x) + \varepsilon = w^T \phi(x) + \varepsilon = \sum_{i=1}^n \alpha_i K(x_i, x) + \varepsilon$

$p(y^*|x^*, X, y) = N(y^*|y_{\text{pred}}, \Sigma_{\text{pred}}) \quad y_{\text{pred}} = K_*^T (K + \lambda I)^{-1} y$

$\Sigma_{\text{pred}} = K_{**} - K_*^T (K + \lambda I)^{-1} K_* \quad K = K(x, X) \in \mathbb{R}^{n \times n}, K_* = K(x^*, X) \in \mathbb{R}^n, K_{**} = K(x^*, X)$

2. Gaussian Process joint distribution $f \sim GP(\mu, K)$ μ : mean $K(x, x')$: covariance

3. Gaussian Process Regression (GPR) $f(x) = f(x) \sim GP(\mu(x), K(x, x'))$

3.1 Noise free observation $D = \{(x_i, f_i)\}, f \in \mathbb{R}^n$ $p(f|X) = N(\mu(X), K(X, X))$

test point (x_t, f_t) , $p(f_t|x_t) = N(\mu(x_t), K(x_t, x_t))$

$P\left[\begin{bmatrix} f(x) \\ f(x_t) \end{bmatrix}\right] = N\left(\begin{bmatrix} \mu(x) \\ \mu(x_t) \end{bmatrix}, \begin{bmatrix} K(x, x) & K(x, x_t) \\ K(x_t, x) & K(x_t, x_t) \end{bmatrix}\right)$

$P(f_t|x_t, X, f) = N\left(\mu(x_t) + K(x_t, X) K(X, X)^{-1} (f - \mu(X)), K(x_t, x_t) - K(x_t, X) K(X, X)^{-1} K(X, x_t)\right)$

$E(f_t) = \bar{f}_t = \mu(x_t) + K(x_t, X) K(X, X)^{-1} (f - \mu(X))$

$\text{Cov}(f_t) = V(f_t) = K_{**} - K_*^T K^{-1} K_*$

3.2 Noise Observation $D = \{(x_i, y_i)\}, y = f(x) + \varepsilon, \varepsilon \sim N(0, \sigma^2)$

$P\left[\begin{bmatrix} f(x) \\ f(x_t) \end{bmatrix}\right] = N\left(\begin{bmatrix} \mu(x) \\ \mu(x_t) \end{bmatrix}, \begin{bmatrix} K(x, x) + \sigma^2 I & K(x, x_t) \\ K(x_t, x) & K(x_t, x_t) + \sigma^2 I \end{bmatrix}\right)$

$\bar{f}_t = \mu(x_t) + K(x_t, X) (K(X, X) + \sigma^2 I)^{-1} (y - \mu(X))$

$\text{Cov}(f_t) = K_{**} - K_*^T (K + \sigma^2 I)^{-1} K_*$

3.3 Noisy prediction y_t $\text{cov}(y_t) = \text{cov}(f_t) + \sigma^2$ 4. Algorithm

Train (X, y) , K (covariance) σ^2 (noise) X_t (test) Rescale y_t : zero-mean GP

$\tilde{X} = \text{inverse Cholesky}(K + \sigma^2 I)$ $\tilde{L} = y - \tilde{X} \tilde{\alpha} = V \Rightarrow \tilde{\alpha} = L^T V$

$\tilde{f}_t = \tilde{X} \tilde{\alpha}, V = L^T L$

$V(f_t) = K(x_t, x_t) - V^T V$ 4.2 Training GP (hyperparameter) max marginal likelihood

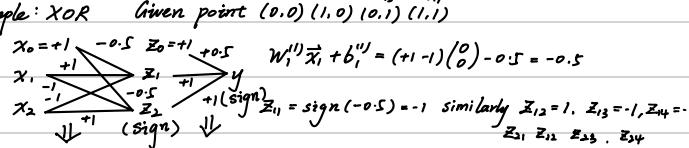
$\sigma^2 = \text{argmax}_{\sigma^2} \text{GP}(y, X, \sigma^2)$

$$\text{Perception: } h(\vec{x}) = \text{sign}(\vec{w}^T \vec{x}) \quad L_2\text{-norm } \|\vec{x}\|_2^2 = \vec{x}^T \vec{x} = \sum x_i^2$$

$$l_1 + \lambda \|w\|_1 = \lambda \sum |w_j| \quad \text{Eigenvector } Av = \lambda v \Rightarrow (A - \lambda I)v = 0 \quad \det(A - \lambda I) = 0 \Rightarrow \lambda_1, \dots, \lambda_n$$

Neural Network 2. Multi-layer Perceptron

Example: XOR Given point $(0,0), (1,0), (0,1), (1,1)$



$$W^{(1)} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}, b^{(1)} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}, W^{(2)} = \begin{pmatrix} 1 \end{pmatrix}, b^{(2)} = 0.5$$

$W_j^{(l)}$ j-th row for unit j

3. Parameter Learning Notation: $W^{(l)}$ weight parameter matrix for layer l
 $\{l \in \{1, \dots, L\}$

$b^{(l)}$ bias terms for layer l ; $a^{(l)} = W^{(l)}z^{(l-1)} + b^{(l)}$, $z^{(l)} = g_l(a^{(l)})$, m_l number of units

3.1 Forward Propagation $z^{(0)} = x \in R^d$, $W^{(l)} \in R^{m_l \times m_{l-1}}$, $z^{(l)} = g_l(W^{(l)}z^{(l-1)} + b^{(l)})$

3.2. Compute Loss (Regression) $L(z, y) = \frac{1}{2}(z - y)^2$ squared loss. AD update rules.

$$\text{layer } l: L(z, y) = \frac{1}{2}(g_l(W^{(l)}z^{(l-1)} + b^{(l)}) - y)^2; \frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial a_{ij}} \frac{\partial a_{ij}}{\partial w_{ij}}$$

$$\alpha - \text{AD learning rate} = (z^{(l-1)} - y) \cdot g'_l(a^{(l)}) \cdot z_j^{(l-1)} = \delta_j^{(l)} \cdot z_j^{(l-1)}$$

$$W^{(l)} \leftarrow W^{(l)} - \alpha \Delta W^{(l)} \Rightarrow W^{(l)} \leftarrow W^{(l)} - \alpha \delta^{(l)} z^{(l-1)T} \in R^{m_l \times m_{l-1}}$$

$$3.3. \text{ Back Propagation } \frac{\partial L}{\partial w_{ijk}} = \frac{\partial L}{\partial g_l} \frac{\partial g_l}{\partial a_{ij}} \frac{\partial a_{ij}}{\partial z_i} \frac{\partial z_i}{\partial a_{ik}} \frac{\partial a_{ik}}{\partial w_{ijk}}$$

$$= (\underbrace{z^{(l-1)} - y}_{\delta^{(l)}}) \cdot g'_l(a^{(l)}) \cdot W^{(l)} \begin{bmatrix} g'(a_1^{(l)}) & \dots & 0 \\ 0 & \dots & g'(a_{m_{l-1}}^{(l)}) \\ 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ z_k^{(l-1)} \\ \vdots \\ 0 \end{bmatrix} \leftarrow j\text{-th row}$$

$$= \delta^{(l)} \begin{bmatrix} W_{11}^{(l)} & \dots & W_{1j}^{(l)} & \dots & W_{1m_{l-1}}^{(l)} \\ \vdots & & \vdots & & \vdots \\ 0 & \dots & g'(a_{1j}^{(l)}) & \dots & 0 \end{bmatrix} \begin{bmatrix} g'(a_{1j}^{(l)}) z_k^{(l-1)} \\ \vdots \\ 0 \end{bmatrix} = \delta^{(l)} W_{1j}^{(l)} g'(a_{1j}^{(l)}) z_k^{(l-1)}$$

$$\frac{\partial L}{\partial b_i^{(l)}} = \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial a_i} \frac{\partial a_i}{\partial b_i} = (z^{(l-1)} - y) \cdot g'(a_i^{(l)}) = \delta_i^{(l)} \quad b_i^{(l)} \leftarrow b_i^{(l)} - \alpha \delta_i^{(l)}$$

3.4 Training Algorithm

$$\text{Input: } D = \{(x_i, y_i)\}_{i=1, \dots, n}, \text{NN: } \begin{cases} l \\ m_l \end{cases} \quad \text{Initialization: } \vec{b}^{(l)}, W^{(l)}$$

repeat for each training example $(\vec{x}, y) \in D$ do $z^{(0)} = \vec{x}$

Forward Prop for $l=1, \dots, L$ do $a^{(l)} = W^{(l)}z^{(l-1)} + b^{(l)}$, $z^{(l)} = g_l(a^{(l)})$

$$\delta^{(l)} = (z^{(l)} - y) \cdot g'_l(a^{(l)}) \quad (\text{compute loss in output layer})$$

$$\text{Back Prop } \delta^{(L-1)} = g'_{L-1}(a^{(L-1)}) \cdot W^{(L)T} \delta^{(L)} \quad W^{(L)T} \leftarrow W^{(L)T} - \alpha \delta^{(L)} z^{(L-1)T}$$

$$\text{for } l=L-1, \dots, 1 \text{ do } \delta^{(l-1)} = g'_{l-1}(a^{(l-1)}) \cdot W^{(l)T} \delta^{(l)}$$

$$\Delta W^{(l)} = \delta^{(l)} z^{(l-1)T} \quad W^{(l)} \leftarrow W^{(l)} - \alpha \Delta W^{(l)}$$

$$\Delta b^{(l)} = \delta^{(l)} \quad b^{(l)} \leftarrow b^{(l)} - \alpha \Delta b^{(l)}$$

until error $= \frac{1}{n} \sum_i (z_i^{(L)} - y_i) < \epsilon$

5. Model Choice 5.1 Output Units (activation function used in output unit)

linear classification $a \in R^k$
regression-unit $h(a) = a$, binary-softmax $h(a) = \frac{e^a}{\sum e^a}$, multi-class-softmax $h(a) = \frac{e^a}{\sum e^a}$

5.2 loss function regression: squared loss $L(\vec{x}, y) = \frac{1}{2}(z - y)^2$ binary classification

cross-entropy loss $L(\vec{x}, y) = -[\ln z_0^{(0)} + y \ln z_1^{(0)}]$, z_0, z_1 : sigmoid

multi-class cross-entropy loss $L(\vec{x}, y) = -\sum c_i y_i \ln z_i^{(0)}$

5.3 Activation Function $\text{sigmoid}(a) = \frac{1}{1+e^{-a}}$ $\text{sigm}(a) = \text{sigm}(a) / (1 - \text{sigm}(a))$

$\text{sign}(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$ $\tanh(a) = 1 - \tanh^2(a)$

$\text{relu}(a) = \max(a, 0)$, $\text{relu}'(a) = \begin{cases} 0 & a < 0 \\ 1 & a > 0 \end{cases}$

6. Autoencoder (unsupervised) $\vec{x} \in R^m \xrightarrow{g} \vec{z} \in R^d \xrightarrow{f} \vec{x}$ over $m > d$ deep autoencoder

undercomplete autoencoder, where $m < d$. loss function $L(\vec{x}, f(g(\vec{x}))$

identity activation $g(\vec{x}) = W\vec{x} + b$ and f get non-linear generalization of PCA

Deep Learning challenges: overfitting, difficult optimization

2. Regularization 2.1 Weight Decay $W_i^{(l)}$: i-th row in $W^{(l)}$

$$a_i = W^{(l)T} \vec{x} + b_i^{(l)} = W_i^{(l)T} \vec{x} + b_i^{(l)} \quad L_2 \text{ regularization: } \frac{\lambda}{2} \|\vec{w}\|_2^2 = \frac{\lambda}{2} \vec{w}^T \vec{w}$$

$$\frac{\partial \vec{w}}{\partial \vec{w}} = \lambda \vec{w} \quad \text{gradient update } W_i^{(l)} = w = w - \alpha (\text{grad} + \lambda w) = (1 - \alpha \lambda) w - \alpha \text{grad}$$

constraint

2.2 Constrained Optimization $\Omega(\vec{w}) \leq k$ new objective $L(x, y) + \lambda(\Omega(w) - k)$

2.3. Early Stopping (use validation set) 2.4 Transfer learning as regularizer

2.5 Ensemble Learning and Dropout 2.6 Batch Normalization (γ speeds up training)

2.8 Data Augmentation 2.9 Adversarial Training 3. Optimization

(1) SGD \rightarrow mini-batch (2) adaptive learning rate α (3) momentum method + (4) conjugate GD (avoid compute Hessian) 4. Parameter Initialization 4.1 Random Init

4.2 Greedy Pretraining Supervised Pretraining (omit $W^{(l)}$ for later layer)

layer Unsupervised Pretraining (x consider one per layer, like autoencoder, reconstruction)

Dimensionality Reduction 3. PCA (project $x \in R^d$ to first r principal components) $(r < d)$

find $U \in R^{d \times r}$ s.t. $\{v_i\}_{i=1}^n$, $v_i = U^T x_i$ have maximum spread $\bar{x} = \frac{1}{n} \sum x_i \in R^d$

center the data $C = \frac{1}{n-1} \sum (x_i - \bar{x})(x_i - \bar{x})^T$ Step 1: $x_i \leftarrow x_i - \bar{x}$ Step 2: $U = \arg \max_{U^T U = I} (U^T C U) \Rightarrow$

$\arg \max_{U^T U = I} U^T C U \Rightarrow \max_{U^T U = I} U^T C U$ s.t. $U^T U = I \Rightarrow \alpha(U, \lambda) = U^T C U - \lambda(U^T U - I) \Rightarrow$

$\frac{\partial \alpha}{\partial U} = 2C_U - 2\lambda U = 0 \Rightarrow U = \lambda U \Rightarrow \max_{U^T U = I} U^T C U = \max \lambda \Rightarrow$ pick the r eigenvectors

with the largest eigenvalues (principal components) $U = [u_1, \dots, u_r] \in R^{d \times r}$

Data Projection: $\vec{x}_i = U^T \vec{x}_i$ (C symmetric U orthogonal $U^T U = I$, $U^T U = I$)

reconstruction original x U (orthogonal) error

3.3 Reconstruction Error $\vec{x}_i = U \vec{z}_i = \sum_{j=1}^r z_{ij} u_j$ $\vec{x} = U \vec{z} = \sum_{i=1}^n \vec{z}_i \vec{x}_i$ $\vec{x} = U \vec{z} = \sum_{i=1}^n \vec{x}_i \vec{z}_i$ $\|\vec{x} - \vec{x}\|^2 = \sum_{i=1}^n \sum_{j=1}^d (z_{ij} \vec{x}_i)^2 = \sum_{i=1}^n \sum_{j=1}^d \vec{z}_i^T U \vec{x}_i \vec{x}_i^T U \vec{z}_i = (n-1) \sum_{i=1}^n \vec{z}_i^T U \vec{U} \vec{z}_i = (n-1) \sum_{i=1}^n \vec{z}_i^T \lambda_i \vec{u}_i \vec{u}_i^T \vec{z}_i = (n-1) \sum_{i=1}^n \lambda_i \|\vec{z}_i\|^2$

(compute U)

4.1. SVD Input x_1, \dots, x_n, r ; $\bar{x} = \frac{1}{n} \sum x_i \in R^d$, $X = [x_1, \dots, x_n] \in R^{d \times n}$, $X = X - \bar{x}$ (center)

$X = UDV^T$

$[U, D, V] = \text{svd}(X)$ $U_{i,j}$ sorted left SVs of X , where $d_{i,i} \dots = \text{diag}(U)$

$U = [u_1, \dots, u_r] \in R^{d \times r}$, $Z = U^T X$

4. MDS (PCA equivalent) $G = X^T X \in R^{n \times n}$, $V = [v_1, \dots, v_n]$ eigenvectors of G , $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_n$, $\lambda_k = \text{diag}(G)_{kk}$ or $Z = P V r^T$

4.3. $X = X_{\text{raw}} - \bar{x}$ (data centering) $\tilde{X} = \bar{X} - \bar{x}$, $\tilde{\lambda}_k = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ (data standardization)

$\tilde{X} = D^{-1} U^T X$, D 奇偶值矩阵 (特征值方阵), U (from PCA) (data whitening)

5 Non-linear Dimensionality Reduction.

5.1 MDS. Kernelize MDS by replacing the $G = X^T X$ by Kernel Matrix $K = K(X, X)$

5.2 Kernel PCA $C = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T \Rightarrow C u_i = \lambda_i u_i$ — eigenvector

$\Rightarrow \frac{1}{n} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T u_i = \lambda_i u_i \Rightarrow \frac{1}{n} \sum_{i=1}^n K(x_i, x_i) \sum_{j=1}^n a_{ij} \phi(x_i) \phi(x_j) = \lambda_i \sum_{j=1}^n a_{ij} K(x_i, x_j)$

$(K(x_i, x_j) = \phi(x_i)^T \phi(x_j)) \Rightarrow K u_i = \lambda_i n K u_i \Rightarrow K u_i = \lambda_i n a_i u_i$ — eigenvector

$[z_i]_j = u_i^T \phi(x_j) = \sum_{j=1}^n a_{ij} \phi(x_i)^T \phi(x_j) = a_{ij}^T K \cdot e_i$ $\forall i = 1, \dots, r$

6. Autoencoder. identity activation - PCA.

non-linear f, g , non linear surface. deep \uparrow highly non-linear