Lecture 1

Representation Learning: ML -> tell the computer "how to learn" from biara b/c red worlde deata has "screetere". en) Lease square regression neek to hande craftede feature representation. for example: Laptop: cpa, ram, weight ... Problem: Image feature representation is in high-dimensional space (256×256×3) ... Wo + W1 × (pixel 1) + pixel 1 prediction pixel k pixel k regression model : Wax (pixel 2) + 1 WK× (pixel k). Voesnit consider the proximity of pixels (no location information given to the model). U information is contained in the images, but not in individual pixels. (pattern is not captured!) J Solution: a feature extractor to detect pattern. pixels -> feature extractor -> features -> prediction function. building good feature is hard

=> representation learning. Linear Algebra Vector $\vec{k} \in R^{d}$ as a coordinate vector $\vec{k} = \begin{pmatrix} k_{1} \\ k_{2} \end{pmatrix}$ 1 7 $(2\pi) \frac{-7}{\pi} = \begin{pmatrix} 2\\ -3 \end{pmatrix} \in \mathbb{R}^2$ $y' = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$ Writting Vector in coordinate form requires choosing a basis " " Standard Basis : e, ..., e. $\hat{e}^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \hat{e}^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \stackrel{\text{lse couch.}}{\leftarrow 2nk} \quad \text{rock.}$ $50, \vec{x} = (x, \dots, xd)^T$ = $\kappa_1 \hat{e}^{(1)} + \kappa_2 \hat{e}^{(1)} + \cdots + \kappa_d \hat{e}^{(d_1)}$ amount in Ist coort ê'z e_{n}). $-\frac{1}{n^{2}} = (3, -2)^{T} = \kappa_{1} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \kappa_{2} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ = 3 = - 2 In coordinate form: $e^{(i)} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \leftarrow ith coordinate. (Stanbark basis vector).$ where the 1 appears in the i-th place Pot product : in -7 = 112 11 11 0 11 cos 0.

where & is angel between it and i?. $\vec{u}\cdot\vec{v}=0$ $\cdot \vec{u} \cdot \vec{v} = 0 \quad (=) \quad \vec{u} \quad anh \quad \vec{v} \quad orchogonal. \quad -) \quad || \vec{u} \cdot || \cdot |\vec{v} \cdot || \cos \theta = 0$ $\vec{u} \cdot \vec{v} = \vec{u} \cdot \vec{v}$ $= (c_1, \dots, u_h) \begin{pmatrix} V_1 \\ \vdots \\ V_h \end{pmatrix} \quad (assumes stankarh basis).$ = U1. V, + · - - + UA. VA $= \frac{A}{Z} u_i \cdot v_i$ Is there any other bases? Stanback basis : 0+6 other lassis : 3 Orthonormal Bases A set of vectors $\hat{\mathcal{U}}^{(1)}\cdots \hat{\mathcal{U}}^{(d)}$ forms orthonormal basis U for R^d if they are mutually orthogonal: $\hat{\mathcal{U}}^{(i)} \cdot \hat{\mathcal{U}}^{(j)} = 0$ • they are unit vectors : $||\hat{u}^{(i)}|| = |$ $(2\pi) \qquad \begin{array}{c} & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & & \\ & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & &$ $(u^{\prime}) \quad \widehat{\mathcal{U}}_{\alpha}^{\prime\prime\prime} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \qquad \widehat{\mathcal{U}}_{\alpha}^{\prime\prime2} = \frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$

$$\begin{array}{c} (1,1) \cdot (1,1)^{2} \pm (1,1)^{2} \pm (-1,1)^{2} \\ = \frac{1}{2} (0) \\ = 0. \end{array}$$

$$\begin{array}{c} (2^{N} + 2^{N}) \\ \end{array}$$

$$\begin{array}{c} (2^{N} + 2^{N}$$

Alternesively, since
$$\vec{x} = x_i^{-4} \hat{u}^{-1} + x_n^{-4} \hat{u}^{-1}$$

 $\vec{x} \cdot \hat{u}^{-1} = x_n^{-4}$
 $\vec{x} \cdot \hat{u}^{-1} = x_n^{-4}$
 $\vec{x} \cdot \hat{u}^{-1} = x_n^{-4}$
 $z \cdot \hat{u}^{-1} = z_n^{-4}$
 $z \cdot$

solve for n, " & n,". x, " = -1/13 $\pi_2^{U} = \frac{10}{13}$ Function of a vector : f: Rh -> Rh' A transformation of is a function that takes in a vactor and vector of same dimensionality. f: Kd -> Kd vector field. for all · Transformation is a $\int_{1}^{2} \frac{1}{f} \, dr \qquad \int_{1}^{2} \frac{1}{f} \left(\frac{1}{x} \right) = \left(3 \times \frac{1}{x} \right).$ (n) $f(\tilde{k}) = (0, \chi_{1}^{2})^{T}$ $\vec{a} = [\vec{a}] + [\vec{a}] = [\vec{a}]$ $\vec{b} = \vec{L} / \vec{J} = \vec{L} / \vec{b} = \vec{L} / \vec{c}$ $\begin{array}{c} -2 \\ -2 \\ -2 \end{array} = \left(\begin{array}{c} 2 \\ -2 \end{array} \right) = \left(\begin{array}{c} 2 \end{array} \right) = \left(\begin{array}{c} 2 \\ -2 \end{array} \right) = \left(\begin{array}{c} 2 \end{array} \right) = \left$ - -7 - -- -- -ù a transformation is generally: Then $\vec{f}(\vec{n}) = (g_1(\vec{n}), g_2(\vec{n}))^{T}$ Linear Transformation : Transformation is linear if $\left[\overline{f}(a\overline{n} + B\overline{y}) = af(\overline{n}) + B\overline{f}(\overline{y})\right]$ en) To check if a transformation is linear:

 $f(\vec{n}) = (x_2 - \alpha_1)^T$ $LHS = \widehat{f}(\alpha \widehat{\delta} + \beta \widehat{\gamma}) = \begin{bmatrix} \alpha \widehat{\delta}_2 + \beta \widehat{\gamma}_2 \\ -(\alpha \widehat{\delta}_1 + \beta \widehat{\gamma}_1) \end{bmatrix}$ $RHS = \alpha \overline{f}(\overline{S}) + \beta \overline{f}(\overline{S}) = \alpha \left(\frac{S_2}{S_1} \right) + \beta \left(\frac{S_2}{S_1} \right) = \left(\frac{\alpha S_2 + \beta S_2}{-\alpha S_1 - \beta S_1} \right)$ So LHS = RHS, \vec{f} is linear. · Suppose f is a linear transformation: $= \varkappa_{1} \overline{f}(\hat{e}^{(1)}) + \varkappa_{2} \overline{f}(\hat{e}^{(2)}).$ standards basislinear cransformation to basis vector. · transformation determined by what it does to basis vector. - Note: to not hold for arbitany transformation $\vec{f}(a\vec{z}+\beta\vec{y}) \neq a\vec{f}(\vec{z})+\beta\vec{f}(\vec{y})$ $\int \left(x_{1}, \hat{e}_{1}^{(\prime)} + x_{2} \hat{e}^{(\prime)} \right) \times x_{1} \int \left(\hat{e}^{(\prime)} \right) + x_{2} \int \left(\hat{e}^{(\prime)} \right) + x_{2} \int \left(\hat{e}^{(\prime)} \right)$ $e_{n}) \stackrel{-7}{f}(\bar{x}) = (o, \kappa_{2}^{2})^{T}$ · b/c of linearity, Mong any direction of change only in scale. $\vec{f}(\lambda \hat{n}) = \lambda f(\hat{x}). \qquad \vec{n} \qquad \lambda \hat{n}^{2}$ $\vec{f}(\hat{x}) \qquad f(\lambda \hat{n}) = \lambda f(\hat{x}^{2}).$ Along a given direction in vector field, can be simplified as unit vector direction. Transformation in other basis:

 $\left[\vec{f}(\vec{n})\right]_{u}$ \notin $\left[\vec{x}\right]_{u}$.

 $e_{n} = \begin{pmatrix} -2 \\ -2 \end{pmatrix} = \begin{pmatrix} 2 \\ -2 \end{pmatrix} \begin{pmatrix} 2 \\ -2$ $\overline{\mathcal{U}}^{(2)} = \frac{1}{\sqrt{2}} (1, 1)^{\frac{1}{2}}$ linear transformation $\vec{u}^{(2)} = \frac{1}{\sqrt{2}} (-1, 1)^{T}.$ orthonormal basis $= \int_{1}^{1} \frac{1}{f} (\hat{u}'') + \delta_{2} \frac{1}{f} (\hat{u}'')$ calculation. $\begin{bmatrix} \overline{\mu} \\ \overline{\mu} \end{bmatrix} u = \begin{bmatrix} \overline{\mu} \\ \overline{\mu} \\$ Mutrices $(A\vec{x})_{i} = \sum_{i=1}^{n} A_{ij} x_{j}$ in general: $\begin{pmatrix} 1 & 1 & 1 & 1 \\ a & a & a^{2/3} \\ b & b & b \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1 a^{2/1} + x_2 a^{2/1} + x_3 a^{2/1}$ Then, suppose f liner transformation: $f^{2}(\overline{x}) = f(x, \hat{e}^{(1)} + x_{2}, \hat{e}^{(2)}) = x, f(\hat{e}^{(1)}) + x_{2}f(\hat{e}^{(2)})$ $f(\hat{e}^{(2)}) = f(\hat{e}^{(2)}) = x, f(\hat{e}^{(2)}) + x_{2}f(\hat{e}^{(2)})$ $A = \left(\vec{f}(\hat{e}^{(1)}) - \vec{f}(\hat{e}^{(2)}) \right) \qquad Matrix as a$ $I = \left(\vec{f}(\hat{e}^{(1)}) - \vec{f}(\hat{e}^{(2)}) \right) \qquad form of linear transformation.$ I GPG: A (n×n) matrix chn he interpretely us a compact representation of linear transformation f: K" -> R". $\int_{-2}^{-2} (\bar{x}) = A \bar{x} = A \begin{pmatrix} x \\ y \end{pmatrix}$ $A\vec{x} = \vec{f}(\vec{x}).$ $= x, f(\hat{e}^{(n)}) + x, f(\hat{e}^{(n)})$ 1 e_{n} $\vec{x} = 3\vec{e}^{(1)} - t\vec{e}^{(2)} = \begin{pmatrix} 3 \\ -4 \end{pmatrix}$ $f(\hat{e}'') = -\hat{e}'' + 3\hat{e}^{(2)}$

 $f(\hat{e}^{(2)}) = 2\hat{e}^{(1)}$ $s_{0}, \tilde{f}(\tilde{n}) = \varkappa, \hat{f}(\tilde{e}^{(n)}) + \varkappa_{2}\hat{f}(\tilde{e}^{(n)})$ $= 3 \begin{bmatrix} -1 \\ 3 \end{bmatrix} + - 4 \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ Alteratively. $A = \begin{bmatrix} -1 & 2 \\ 3 & 0 \end{bmatrix}$ $f(\vec{x}) = A \kappa = \begin{bmatrix} -1 & 2 \\ 3 & 0 \end{bmatrix} \begin{bmatrix} -3 \\ -4 \end{bmatrix}$ Matrices in Other Bases. $A^{\mu} = \begin{pmatrix} 1 & 1 \\ CF(\hat{\alpha}^{(\prime)}) \overline{1}_{\mu} & \cdots & \overline{LF}(\hat{\alpha}^{(d)}) \overline{1}_{\mu} \\ u & u \end{pmatrix}$ $\vec{f}(\vec{n}) = \vec{f}(x, \vec{u}'' + x, \vec{u}'') = x, \vec{f}(\vec{u}'') + x, \vec{f}(\vec{u}'') + x, \vec{f}(\vec{u}'')$ $= 7 \left[f(\bar{n}) \right]_{\mathcal{U}} = \left[\varkappa_{,}^{\mu} \bar{f}(\hat{u}'') + \varkappa_{,}^{\mu} \bar{f}(\hat{u}'') \right]_{\mathcal{U}} = \varkappa_{,}^{\mu} \left[\bar{f}(\hat{u}'') \right]_{\mathcal{U}} + \varkappa_{,}^{\mu} \left[\bar{f}(\hat{u}'') \right]_{\mathcal{U}}$ by U= orchonomal = Au [ze,]u basis. Consider of which mirror a vector over 45° line. $\begin{array}{c} x_{1} \\ x_{2} \\ x_{1} \\ y_{2} \\ y_{2} \\ x_{2} \end{array} \xrightarrow{\left(\begin{array}{c} x_{2} \\ x_{2} \end{array} \right)} = \left(\begin{array}{c} x_{2} \\ x_{1} \end{array} \right)}{\left(\begin{array}{c} x_{1} \\ x_{2} \end{array} \right)} = \left(\begin{array}{c} x_{2} \\ x_{1} \end{array} \right)}{\left(\begin{array}{c} x_{2} \\ x_{2} \end{array} \right)} = \left(\begin{array}{c} x_{2} \\ x_{2} \end{array} \right)}$ $A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_n \end{bmatrix}$ χ, $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ $e_{n})$ $\hat{u}^{(1)} = \frac{1}{\sqrt{2}} (1, 1)^{T}$

 $\hat{\mu}^{(n)} = \frac{1}{\sqrt{2}} \left(-1, 1 \right)^{T}$

=) $\vec{f}(\vec{u}'') = \vec{u}''$ transformation of orchogonal basis $\hat{\mu}^{(2)}_{F_{\chi}}$ F_{χ} $\vec{f}(\hat{u}^{(2)}) = -\hat{u}^{(2)}, \qquad \text{if } u \text{ is orthornormal},$ > f(a") $= \begin{bmatrix} \hat{u}^{(1)} \cdot \hat{u}^{(1)} & - \hat{u}^{(2)} \cdot \hat{u}^{(1)} \\ \hat{u}^{(4)} \cdot \hat{u}^{(2)} & - \hat{u}^{(2)} \cdot \hat{u}^{(4)} \end{bmatrix}$ $= \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$ linear transformation of orthornormal basis CS q matrix representation. Eigenvector & Eisenvalue $A \vec{J} = \lambda \vec{J}$ $\gamma \vec{V}$ $eigenvetor \quad eignuclue \lambda$ of A linear transformation. f => eigenvector of f with eigenvalue) is a nonzero vector i s.t. f(ご)= x ご. i den: When f applies to one of its eigenvector, f scales it. (not changing any direction) $f(\vec{v}) = \lambda \vec{v}$ Finding eigenvector => O Graphically, find vector that are in the same direction after transformation f.

(2) Note: not all matrices haven even one eigenvector = a set of n eignvector that are mutually orthogonal. Thm: Lee A be an non "symmetric" matrix. (elerc can be many spe) (Spectral => thre exist n eignvators of A which are all matually orthogonal. "symmetric" (=> AT = A O symmetric linear transformation have axes of symmetry. 3 The axes of symmetry are orthogonal to one another. 3 The action of f along Ax 1s of symmetry scaleth input. (Size of scaling might hifferent $f\left(\frac{-2^{\prime\prime}}{u}\right) = \lambda, \overline{u}^{\prime\prime}$ Note, if A is diagonal, its eigenvector are $f\left(\frac{-2}{u}\right) = \lambda_2 \frac{-2}{u}$ the standard basis vector Total Symmetry : infinity many eigenvector every vector is an eigenvector of A! Eigenvector as basis vector: If A is symmetry matrix, pick a of its eigenvector $\hat{\mathcal{U}}^{(1)}$... $\hat{\mathcal{U}}^{(d)}$ form an orthonormal basis. (eigen hecomposition). Any veccor The can be written in this basis : $=) [\vec{n}] u = \begin{bmatrix} v_1 \\ \vdots \\ \vdots \end{bmatrix}$ $\vec{x} = b, \hat{u}'' + \cdots + b_{a}\hat{u}^{(a)}$ $\mathcal{U} = \{\hat{u}'', \cdots, \hat{u}'^{b}\}$ (eigen-basis) as eigen-basis expression. $f(\vec{n}) = A \vec{n} = A(b, \hat{u}'' + \cdots + b_{d} \hat{u}'^{d})$

 $= b_1 (A\hat{\mu}^{(\prime)}) + \cdots + b_4 (A\hat{\mu}^{(d_2)})$

$$= b_{1} (\lambda \lambda^{-1}) + \dots + b_{k} (\lambda \lambda^{-1})$$

$$= (b_{1} \lambda) \lambda^{-1} + \dots + (b_{k} \lambda \lambda) \lambda^{-1} (b_{k})$$
So, if A is generating, pipe base U is as natural basis.

$$f(\lambda^{2}) = \lambda \lambda^{2} + \lambda b \lambda^{-1} + \lambda b \lambda \lambda^{-1} + \dots + \lambda b b \lambda^{-1} \lambda^{-1}$$
sole if out and and inter the of CE-Jac by eigenvalue λ :
Eigenverse as optimizer
er) find while vacer Z s.t. $\|A \overline{X}\|$ largues.

$$\frac{1}{|A|^{2}} b^{-1} + \lambda b \lambda^{-1} \lambda^{-1} + \lambda b \lambda^{-1} \lambda^{-1} + \lambda b b \lambda^{-1} + \lambda b \lambda^{-1} + \lambda^{-1} + \lambda b \lambda^{-1$$

To maximize Ti ATi over unic vector pick to be top eigenvector pf: $\vec{x} = A\vec{n} = (b, \hat{u}''' + b_2 \hat{u}'') \cdot (b, \lambda, \hat{u}'' + b_2 \lambda_2 \hat{u}'^2).$ $= b_{1}^{2} \lambda_{1}^{1} + o + o + b_{2}^{1} \lambda_{2}^{1}$ $= b_1^{2} \lambda_1^{2} + b_2^{2} \lambda_2^{L}$ ··· (same proof as above). e_{x} Max 4x, $^{2} + 2x_{2}^{2} + 3x$, x_{2} subject to x, $^{2} + x_{3}^{2} = 1$. || || 2 || = | $s_{0}, \vec{n}, A\vec{n} = \begin{pmatrix} x_{1} \\ x_{2} \end{pmatrix} A \begin{pmatrix} x_{1} \\ x_{2} \end{pmatrix} = a_{x_{1}}^{1} + (b+c)x_{1}x_{2} + b_{x_{2}}^{2}$ $T \qquad T \qquad T \qquad T$ $= 4x_{1}^{1} + 3x_{1}x_{2} + 2x_{2}, A symmetry$ $=7 a = 4 \qquad A = \begin{bmatrix} a & b \\ c & 4 \end{bmatrix}$ d = 2To maximize $\bar{\lambda}^2 \cdot A \bar{\lambda}^2 \quad s. \epsilon. \quad ||\bar{\lambda}^2|| = 1, \quad A = \begin{pmatrix} 4 & 1.5 \\ 1.5 & 2 \end{pmatrix}$ finh it * , top eigenvector of A Eigenvector as Equilibria : iken f(x) rotates x toward the "top" eigenvector. \vec{V} is an equilibrian $\vec{f}(\vec{u}) = \lambda \vec{u}$ Methole of computing top eigen vector/value. (power method). · inicialize 2 · repeat until convergence: -> (1+1) = A = "1) / 11 A = "11

Change of Basis Matrices $\vec{X} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = a_1 \hat{e}^{(1)} + a_2 \hat{e}^{(2)} \quad 1$ Represent change of basis as matrix. $\vec{X} = \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} = a_1 \hat{e}^{(1)} + a_2 \hat{e}^{(2)} \quad 1$ $A = \begin{bmatrix} f(\hat{e}^{(1)}) & f(\hat{e}^{(2)}) \end{bmatrix} = \begin{pmatrix} \hat{e}^{(1)} & \hat{u}^{(2)} & \hat{e}^{(2)} & \hat{u}^{(1)} \\ \hat{e}^{(2)} & \hat{u}^{(2)} & \hat{e}^{(2)} & \hat{u}^{(2)} \\ \hat{e}^{(2)} & \hat{u}^{(2)} & \hat{e}^{(2)} & \hat{u}^{(2)} \end{pmatrix}$ $Suppose \qquad \hat{u}^{(1)} \quad \hat{u}^{(2)} = \begin{pmatrix} \hat{e}^{(1)} & \hat{u}^{(2)} & \hat{e}^{(2)} & \hat{u}^{(2)} \\ \hat{e}^{(2)} & \hat{u}^{(2)} & \hat{e}^{(2)} & \hat{u}^{(2)} \end{pmatrix}$ $\tilde{x}^2 = b_1 \hat{x}^2 + b_1 \hat{x}^{(2)}$ > So change of basis as maximal; $\vec{f}(\vec{x}) = [\vec{x}]_{u} = \mathcal{U}\vec{x}.$ e_{x}) $\hat{u}^{\prime \nu} = (\frac{\sqrt{3}}{2}, \frac{1}{2})^{7}$ $\hat{\mathcal{U}}^{(2)} = \left(-\frac{1}{2}, \frac{\sqrt{3}}{2}\right)^{T}$ If U is change of basis $\vec{\lambda}^{2} = (1/2, 1)^{T}$ matrix : $S_{\sigma_{1}}[\bar{x}]_{u} = (l\bar{x}) = \begin{pmatrix} 3/2 & 1/2 \\ -1/2 & 5/2 \end{pmatrix} \begin{pmatrix} 1/2 \\ 1 \end{pmatrix}$ $\overline{x}^2 = \mathcal{U}^{T}(\overline{x}^2)\mathcal{U}$ $= \begin{bmatrix} -\sqrt{3}/4 + \sqrt{2} \\ -\sqrt{4}/4 + \sqrt{3}/2 \end{bmatrix}$ Diagonalization Recall matrix as linear transformation f $A = \left(\begin{array}{cc} -\frac{1}{f} & -\frac{1}{e} \\ f & -\frac{1}{e} \end{array}\right) \cdots f & f & f \\ \frac{1}{e} \end{array}$

If we use a difference basis $U = \frac{1}{2} \widehat{u}^{(1)} \cdots \widehat{u}^{(d_1)} \frac{1}{2}$ $A_{\mu} = \left(\mathcal{L} \, \tilde{f}(\tilde{u}^{\mu}) \right)_{\mu} \cdots \mathcal{L} \, \tilde{f}(\tilde{u}^{\mu})_{\mu} \right)_{\mu}$ $[\bar{x}]_{u} = U\bar{x}$ => $\vec{y} = A \vec{x}^2$ => Suppose A is matrix, find basis U $\begin{bmatrix} \vec{y} \end{bmatrix}_{\mathcal{U}} = A_{\mathcal{U}} \begin{bmatrix} \vec{x} \end{bmatrix}_{\mathcal{U}}.$ S. C. Au is diagonal. U If A is symmetric matrix, pick à of its eigenvector û", ... û da, to form orthogonal basis. $Au = \left(\begin{array}{c} [f(\overline{u}^{3})']_{u} & \cdots \end{array} \right) = \left(\begin{array}{c} \lambda_{i} & 0 \\ \vdots \\ 0 & \lambda_{k} \end{array} \right)$ $diagonal \stackrel{l}{=} \left(\begin{array}{c} 0 & \lambda_{k} \\ 0 & \lambda_{k} \end{array} \right)$ why? matrix of linear transformation f in basis of orthonormal
 its eigenvector is a diagonal matrix.
 entries are eigenvalue. $Compute f(\tilde{n}^2)$ $\vec{y} = A\vec{x}$ OK C. change bass to eigenbasis U: $U\bar{y}^{2} = UA\bar{x}^{2}$ Un $[J]_{u} = UAU^{T}[\bar{x}]_{u}$ 2 apply f in eighbabis vie Au. $Au[\vec{x}]_u = UAu^{T}[\vec{x}]_u$ $[\vec{y}] = An[\vec{x}]u.$ Au=UAUi 3. go back to standard basis. b/c $u' = u^{7}$ =7 $A = U^T A u U$ $\mathcal{U}^{\intercal}[\overline{\mathcal{F}}]_{\mathcal{U}}$

 $\vec{f}^{2} = \mathcal{U}^{T} \Big[An \Big[\mathcal{U} \vec{z}^{2} \Big] \Big] = A \vec{z}^{2}$ 300 => U'Aull = A Thm A be nxn symmetric matrix. Forthogonal matrix U "U diagonalize A". hiagonal matiix N S.t. $\left(A = \mathcal{U}^T \wedge \mathcal{U}\right)$ where rows of U are eigenvacior of A & entries of A are eigenvalue. PCA High dimensional data [" Rebuce dimensionality by minimizing loss of information () R^a -> R' (z² -> Z.) en) two features $\begin{pmatrix} \varkappa_{1} \end{pmatrix}$ phone will the and weight. $\begin{pmatrix} \varkappa_{2} \end{pmatrix}$ ψ ZER' = combination of X, and X, (mixture of feature). Z = -2 -2 (linear combination). mixture coofficience, ne assume 11 ul 11 = 1 iden : x² in this of u, scale to some magnitude. x² ∧ x u Since u. x = 11 ull . 11 x 11 105 0 7 11 2 11 COSB = (|x || cus 0.

So , is define a direction

$$\vec{Z}^{ij} = \vec{X}^{ij} \cdot \vec{d}$$
 represe position/magnitude =t \vec{X} along \vec{d} .
Just plit,
along bits,
along bits,

Simple extensione matrix C is back matrix where i,j early is
defined so be 60.
Gij =
$$\frac{1}{2}$$
, $\frac{1}{2}$,

idea top eigenvector of covariance matrix points in direction of max voriance. (proof:) let $\vec{\mu}$ be unit vector. $||\vec{\mu}|| = |$ Z⁽ⁱ⁾ = $\vec{k}^{(i)} \cdot \vec{\mu}$, new features for $\vec{k}^{(i)} \in \mathbb{R}^{d}$. $V_{ur}(z) = \frac{1}{n} \frac{2}{z} (z'' - \mu_z)^2 = \frac{1}{n} \frac{2}{z} (\vec{z}'' - \mu_z)^2.$ $[\mathcal{U}_{2}: \text{ them of } \frac{1}{2} \frac{1}$ $= \vec{u} \left(\frac{1}{2} \sum_{i=1}^{n} \vec{v}^{(i)} \right)$ Variance of data in direction is: = u. Mx The data points $V_{ar}(z) = \frac{1}{2} \sum_{i=1}^{2} \left(\frac{1}{2} \sum_{i=1}^{i} \left(\frac{1}{2} \sum_{i=1}^{i} \frac{1}{2} - \frac{1}{2} \sum_{i=1}^{2} \right)^{2}$ $=\frac{1}{2}\left(\vec{x}^{\prime\prime},\vec{u}-\vec{u}\cdot\mu_{x}\right)^{2}$ $= \frac{1}{2} \sum_{i=1}^{\infty} \left(\frac{1}{u} \left(\frac{1}{x} - \frac{1}{u} \right) \right)^{2} = g(\frac{1}{u})$ whose are given by data. we can henote Var(z) as g(ii) Goal : finh unit vector il maximizes function 9. max g(ū) s.t. ||ū||=1 Observation $D: g(\vec{u})$ can be written as $g(\vec{u}) = \vec{u}^T (\vec{u})$ covariance matrix of data points 2" i=1 ... n. $\left(\left(=\frac{1}{2}Z^{T}Z, Z=X-\left(-\frac{u^{T}}{2}\right)\right)\right)$ Observation D: our problem becomes $max \quad \mu^{2} C u^{2} \qquad \Longrightarrow \qquad The solution is:$ $max \quad \mu^{2} C u^{2} \qquad \Longrightarrow \qquad \int \mu^{2} * = C's \quad top (largero) \quad eism vecusion$ $5.\tau. ||\vec{u}||_{1} = 1$ $\int g(\vec{u}^*) = \lambda_{max}$ of C.

PCA Algorithm: Given data points z⁽¹⁾ ... z⁽ⁿ⁾ ER^d. $\mathcal{M}_{\mathsf{X}} = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{x^{i}}$ O Compute covariance matrix C @ (ompute the top eigenvector i of C $\zeta = \frac{1}{2} Z^T Z$ 3 For i E = 1,..., n', create new feature. $Z = \chi - \begin{bmatrix} -\mu_x' \\ \vdots \\ -\mu_y' \end{bmatrix}$ $z'' = \vec{u} \cdot \vec{z}'$ $= \begin{bmatrix} & \chi^{(y)} & - & \mu_{x} \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\$ PCA for k principal components (/2 2 1) 2 features: $\vec{z} = (\vec{z}_1, \vec{z}_2)$. $Z_{1} = \overline{U}^{(1)} - \overline{x}^{(2)} = \mu U_{1}^{(1)} x_{1} + \cdots + \mu U_{k}^{(2)} \kappa d_{k}$ $Z_{2} = \tilde{u}^{(2)} \cdot \tilde{\kappa} = \mu^{(2)}_{1} \chi_{1} + \cdots + \mu^{(2)}_{d_{1}} \chi_{d_{1}}$ How to find? O choose $u^{(2)}$ to be orthogonal to $u^{(2)}$ => No "redundance" (~ 2" & ~ no overlap). 00000 du 000 du 000 X2 weight (2) Since for symmetric matrix, if i and i are eigenvectors with

distince eigenvalues, they are automatically orchogonal. So, choose U to be an eispave-tor of C. U called the second principle component. U⁽²⁾ is the second eigenvector of C. (out of all vectors orthogonal to the prinicple component, points in direction of max variance), O (ompute covariance matrix C, top $k \leq d$ eigenvectors \mathcal{U} , \mathcal{U} . @ For any vector \$ ER, its new representation in R is: $\vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \text{in matrix form:} \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}, \dots, \vec{z}_{k})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T}, \text{ where } \qquad \vec{z}^{2} = (\vec{z}, \dots, \vec{z})^{T$ $Z_{1} = \frac{1}{2} \cdot \vec{u}^{(2)}$ $\overline{\mu}^{2} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{n} E^{ij} E^{j}$ $\frac{1}{Z_{k}} = \frac{1}{X} \cdot \vec{u}^{(k)}$ $\mathcal{L} = \frac{1}{n} \mathcal{W}^{\mathsf{T}} \mathcal{W} ,$ $W_{n\times h} = \chi_{n\times h} - \begin{bmatrix} -\overline{\mu} & \overline{\mu} & \overline{\mu} \\ \vdots & \vdots \\ -\overline{\mu} & \overline{\mu} & \overline{\mu} \end{bmatrix}$ lee $\mathcal{U}_{dxk} = \begin{bmatrix} -> & \cdot \\ \mathcal{U} & \cdot \\ \end{bmatrix} = \begin{bmatrix} -> & \cdot \\ \mathcal{U} & \cdot \\ \\ \cdot & \cdot \\ \end{bmatrix}$ Mis ('s top i eigenvector Thus, $\vec{z}' = \begin{bmatrix} -\vec{z} & (D)^T \\ -\vec{z} & (A)^T \end{bmatrix} \begin{bmatrix} \vec{z} \\ \vec{x} \end{bmatrix} = \begin{bmatrix} 0^T \\ \vec{z} \end{bmatrix}$ $= \begin{pmatrix} z' \\ \vdots \\ z_k \end{pmatrix} \in \mathcal{K}^k$

$$\begin{array}{c} \mathbb{P}(A \quad \text{reservation} \quad \text{rote:} \\ \mathbb{R}^{k} \rightarrow \mathbb{R}^{k} ? \\ \mathbb{S}_{uppose} \quad new \quad \text{referencession} \quad of \quad \mathbb{X} \quad is \quad \mathbb{R}_{i} : \\ \mathbb{Q}_{i} \in \mathbb{R}^{k} \cdot \mathbb{Z}^{(k)} \\ \mathbb{Q}_{i} (\text{recomposed}) \\ \mathbb{Q}_{i}$$

Ever PCA, we are assembly doing ...
choose etc bans {
$$a^{(i)} \cdots a^{(k)}$$
} to be special bank :
the exp k approveder of C, where C is the containance matrix of X.
(Call = $h w(w)$)
bhoffie : $h = \frac{1}{2}$ has now variance
 $h = 2$: reconservation ever is minimize $\| X - V \vec{z} \|^{\lambda}$.
 $\# 3$: $2 : 2 : 3$ are uncorrelated!
(An factor).
View #1 Maximizing Variance of the new dark.
(choosing ness: "investing" new factors which are not rebundant).
"Tooch Variance": sum of variances of relevan Z (up to principle componence b).
 $Z = \begin{bmatrix} -\frac{3}{2}^{(i)} & -\frac{3}$

PCA minimizes the reconstruction error. "the "best" projection of points onto a linear subspace of dimensionalize k.

total reconstruction error = $\sum_{i=1}^{n} ||\vec{x}^{(i)} - ||\vec{z}^{(i)}||^2$.

Claim: choosing U to be top eigenvectors of C minimizes reconstruction error, among all choices of orthonormal U. View #3 Perorrelation 12CA learns new representation by rotating data into a basis where features are uncorrelated. $\sqrt{}$ (natural basis vectors are the principle directions). proof: $2: \begin{bmatrix} 2\\ \vdots\\ 2k \end{bmatrix}, \\ 2nxk = \begin{bmatrix} 2^{ny} \\ \vdots\\ 2^{nxk} \end{bmatrix}$ Eisenvectors of Covariance matrix $\bar{\mu}_{z} = \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2}$ Relation between $Z \mathcal{L} \times : \vec{Z} = \mathcal{V}^{T} \vec{z}, \quad \mathcal{V} = \begin{bmatrix} \vec{u}^{n} \cdots \vec{u}^{(k)} \end{bmatrix}$ $= \overline{\mu_{z}} = \frac{1}{2} \overline{z_{i=1}^{(i)}} = \frac{1}{2} \overline{\mu_{x}}$ So, to show Z; & Z; uncorrelated, V i = j => show (or(Zi,Zi)=0, Vi≠j => (ov of Z should be a diagonal matrix. $lov_2 = \frac{1}{n} \widetilde{Z}^{\dagger} \widetilde{Z}$, $\widetilde{Z} = Z_{n \times k} - \begin{bmatrix} -M_2^{\dagger} T \\ \vdots \\ -M_2^{\dagger} \end{bmatrix}$. $\left(\widetilde{z}_{j}=\left(\overline{z}'^{j}-\overline{\mu}_{z}\right)^{T}=\left(\overline{v}^{T}\overline{z}'^{j}-\overline{v}^{T}\overline{\mu}_{z}\right)^{T}$ $= \left(\frac{2}{\lambda} \left(j \right) - \frac{1}{\mu_{n}} \right)^{T} U$ (Conterch original later $=) \quad \widetilde{Z} = \begin{bmatrix} & -(\widetilde{z}^{2}^{(1)} - \widetilde{\mu}^{2}_{x}^{(1)}) - \\ & & U = WU \\ & & U = WU \\ & & & U = WU \end{bmatrix}$ $\int C_{onv_z} = \frac{1}{n} (wv)^T Wv$ = 1 UTWTWV = UT (+ WTW) U = UT CU, where C is Conv. of re.



I den: data expressed with a dimension, but it really confined to k-dimensional region. Original hara = h k < d. hidden structure = k this is called a manifold (a lover bimensional space). · d is the ambient dimension · k is the intrinsic dimension. Manifold Learning -> recover the low-dimensional monifold Data in high-dimensional linen manifold (PCA) Non-linear manifold (laplacion eigenmap). Euclideen Distance (dist. in data space) Greadesic distance. (hist. in manifold space) Euclidean vs. Geodesic Distances in this example. Euclidean distance: the "straight-line" distance d(x,y), d^{*}(x,y) x & y are close in R (de is small), but • Geodesic distance: the distance along the manifold $\sqrt{9}(\vec{x}, \vec{y})$ 11 the geodesic distance d⁹ is large : ž in R^k, 3x, 3y should be far away. Euclidean Geodesic if data 2 linear manifold, geoliesic 2 euclidean. Problem nap have in Kt to Kk, K<h s.z. if $\vec{n} \notin \vec{j}$ close in geodesic distance to \mathbb{R}^k , K<d $d_{cF}^{9} < d_{FP}^{9}$ in \mathbb{R}^{d} $\Rightarrow d_{c'F'}^{e} < d_{F'b'}^{e}$, in \mathbb{R}^{K} also close in Euclidean disconce. e close in geodesic close in Euclidean in Rd (d=2) A, B, C, D, E ER² A', B', c', D', $E' \in \mathbb{R}^1 \qquad \mathbb{R}^k$ (k=1)

Solution. O Given data in Kh, XERh. @ Built a similarity graph from points 3 Vinensionalius reduction: find spectral embedding in graph laplacian. Build a similaring graph: Why graphs? Step 1: graph is 1D structure estimating build a similarity graph the manifold from points cire. graph can approximate geodesic distance well!)) -> approximate geodesic distance well. mistakes ! distance but far in geodesic distance (Shortest path on the graph is large) Three Approaches 1) Eplison neighbors graph. Creace a graph with one notice i per point Z'is All edge becker i and j if $||\vec{z}^{(i)} - \vec{z}^{(i)}|| \leq \varepsilon$. => Unneighteh graph. adjacency matrix Work = [wij] wij = 1 if has edge else O. clarge E) Approach #1: Epsilon Neighbors Graph What will the graph look like when ϵ is small? What · E too small, graph about when it is large? noe gook 81 < 82 < 83 < 84 under connected 82 23 · · . . · E too large, graph OVER Connectation noe good. (small E)

2 K - neighbors Graphs · Create a graph with one notice i per point Z" each note i and its K close neighbors. Add edge between ο Οιαμιι if K= 3 x v v => unweighter graph. k=3 ", ER Wn= [wij] e.g. Wnin will be symmetrics asymptric. $W_{AB} = 1$, $W_{BA} = 0$ "if k=3, B is A's neighbor, but A is not B's neighbor! oh (directed graph (not desirable) => WAB = WBA = 1 K4 Approach #2: k-Neighbors Graph Is it possible for a k-neighbors graph to be disconected? k1 < k2 < k3 < K4 K 3 Fully Connected Graph . Create a graph with one note i per point 2013 · Abb edge between every pair of nodes. Assign weight of h(Zⁱⁱ), Zⁱⁱ). =) weighteh graph. in particlar, two close points has large neight. Common similarity function: $h(\vec{x}, \vec{y}) = e^{-||\vec{x} - \vec{y}||^2/6^2}$ $h(\vec{x},\vec{y})$ value of 6 influence the weight. far dise close dist.



Less is low if similar points are class
(prive 2.8.5 and circler if
$$(f_1, f_2)^2$$
 small?
 $(cose(\vec{f}) = \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} w_{ij}(f_1 - f_2)^k$.
 $(cose(\vec{f}) = \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} w_{ij}(f_1 - f_2)^k$.
where w_{ij} is mariple between i and j is
similarity matrix.
 $if w_{ij} \approx 0$, $(f_1, -f_2)^k$ can be large.
 $if w_{ij} \approx 1$, $(f_1 - f_2)^k$ should be small.
goal: min cose(\vec{f}) problem: cose is clamays minimizab when $\vec{f} = 0$.
 \vec{f}
 $(cose(\vec{f}) = \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} w_{ij}(f_1 - f_3)^k$, given and similarity matrix.
 $cose(\vec{f}) = \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} w_{ij}(f_1 - f_3)^k$, given and similarity matrix.
 $where D is the keyree macrix.$
 $where D is the keyree macrix.$
 $where D is the keyree macrix.
 $where D is the keyree macrix.$
 $where D is the keyree macrix.
 $where (\vec{f}) = \sum_{j=1}^{\infty} \sum_{j=1}^{\infty} w_{ij}(f_1 - f_3)^k = 2\vec{f}^T L \vec{f}$
 $proof:$
 $(f_1 - f_3)^k = f_1^k - 2f_1f_2 + f_3^k$.
 $\sum w_{ij}(f_1 - f_3)^k = Z_i w_{ij}f_1^k + T_j w_{ij}f_2^k - 2 \sum_{ij}^{\infty} w_{ij}f_{ij}$$$

 $\sum_{ij} W_{ij} f_i^2 = \sum_i f_i^2 \sum_j W_{ij} = \sum_i d_i f_i^2 \quad (d_i = \sum_j W_{ij}, which is the degree).$ $\sum_{i,j} w_{i,j} f_{i,j}^{*} = \overline{\Sigma}_{i,j} f_{i,j} f_{j,j}^{*} \iff \overline{\Sigma}_{i,j} d_{i,j} f_{i,j}^{*}.$ 5_{0} $2\overline{2}_{i}h_{i}f_{i}^{2} - 2\overline{2}_{ij}w_{ij}f_{i}f_{j} = 2(\overline{2}_{i}h_{i}f_{i}^{2} - \overline{2}_{ij}w_{ij}f_{i}f_{j}).$ Laplacian matrix L = D - W, $= 2\vec{f}^T L \vec{f}$. Note that : degree (i) = $\sum_{j=1}^{n} w_{ij} = sum of weights associated with node ;$ Degree matrix $D = \begin{bmatrix} a_{11} \\ a_{22} \end{bmatrix}$ $d_{ii} = degree(i) = \sum_{i=1}^{n} w_{ij}$ Our optimization becomes: $\begin{array}{c}
\min_{f} \cos \left(\left(\frac{1}{f} \right) \right) = \frac{1}{f} \left(\frac{1}{f} \right) \\
\frac{1}{f} & \sum L = D - w, \quad \left(L \quad sgmmetrg \right).
\end{array}$ Subject to $\|\tilde{f}\| = |\mathcal{L}| - (1 - 1)^T$ By ideas of PCA and also resolve the constraints $\vec{F}^{\perp}(1\cdots 1)^{T}$. Solution: choose \vec{f} to the eigenvector of L with smallest eigenvalue >0(orthogonal to $(1, \dots 1)^7$ by spectral Thm), 2 k > 1

Find bottom & eigenvectors with eigenvalues 70.

Thus, with k eigenvector $\vec{f}^{(i)} \cdots \vec{f}^{(k)}$, each node is mapped to a point in 12k. so, consider note i: Practical Issue: new coordinate $f_i^{(\prime)}, f_i^{(2)} \cdots f_i^{(k)}$ 1 Lnorm = 12 - 1/2 L 12 - 1/2. $f^{(1)} = f^{(2)} \cdots f^{(k)}$ where $15^{-1/2}$ is the diagonal matrix whose ith diagonal entry is $1/\sqrt{h_{ii}}$. izh _> o o ... o Normalize laplacian to find enbeling for note i: z", Eple eisenverors . (Laplacian Eigenmaps). Supervised Learninz Example: given others movie rating -> predice your rating of the movie. $\overline{\mathcal{H}} = (\mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3, \mathcal{X}_4, \mathcal{X}_5)$ Regression -> number $H(\vec{x}) = prediction.$ fpreduction function. classification -> class labels. Empirical Risk Minimization (ERM). O choose a hypothesis class Deboose a loss function

3 minimize expected loss (empirical risk).

H: a see of $H(\vec{z})$ Hypothesis class H. prediction $H(\vec{n})$ La set of possible prediction functions). function. \mathbf{b} · H := linear function The more complex the hypothesis class, the greater the danger of overfitting · It := decision tree · H := neural ness Occamis Razar: assume It is simple l Assume lineur preblicion function: $H(\vec{x}) = W_0 + W_1 \kappa_1 + W_2 \kappa_2 + \cdots$ (parameterization). $\vec{w} = (w_0, \dots)^T.$ · If there are a features, there are at parameters. $|f(\vec{x}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d.$ $= w_{3} + \sum_{i=1}^{h} w_{i} \times i$ Augmentele feature vector \longrightarrow Aug (\vec{x}) . $\vec{x} = \begin{pmatrix} x_i \\ \vdots \\ x \\ k \end{pmatrix}$, Aug $(\vec{x}) = \begin{pmatrix} x_i \\ \vdots \\ x_{k} \end{pmatrix}$ $\in \mathbb{R}^{d+1}$ $\stackrel{=}{\longrightarrow}$ $H(\vec{x}) = Aug(\vec{x}) \cdot \vec{w}$. · The surface of a prediction function is a hyperplane Loss function quantifies how wrong a single prediction is ... L (H(zⁱⁿ), Yi) T grounde truch. prediction of ith draga. Cxumple loss function: () Absoluce loss: | H(z'') - Y; | ∂ 5 quare las: (H(zⁱ)) - Ji)². $Optim; ze s.t. H(\tilde{z}^{(i)}) \approx Y;$

A gook H is gook on average over all data. Expected loss (empirical risk): average loss across al data $R(H) = \frac{1}{2} \sum_{i=1}^{n} L(H(\vec{x}^{i'}), Y_i)$ $\begin{cases} n^{ii}, y_i \end{cases}$ en) Experied square loss (MSE) $K_{sq}(H) = \frac{1}{2} \sum_{i=1}^{n} \left(H(\overline{x}^{i}) - y_{i} \right)^{2}$ Minimizing Expande Loss Find It minimizing $R_{sq}(H) = \frac{1}{2} \sum_{i=1}^{n} (H(x^{ii}) - Y_i)^{2}.$ $\int \int \frac{1}{\omega^2} \pi min Rsq = \min_{i=1}^{\infty} \frac{1}{\omega^2} \left(\overline{w} \cdot Aug(\overline{z}^{ij}) - y_i \right)^2$ loseh form solucion: (MSE) features 21, ~ 264. $\overline{w}^{*} = (x^{T} x)^{'} x^{T} \overline{y}^{'}$ $X \text{ is hesign matrix } X = \begin{pmatrix} Aug(\overline{x}^{'''}) \\ \vdots \\ Aug(\overline{x}^{'''}) \end{pmatrix} \begin{pmatrix} I x_{1}^{'''} x_{2}^{'''} \cdots x_{d}^{'''} \\ I x_{2}^{'''} x_{2}^{'''} \cdots x_{d}^{'''} \end{pmatrix} points$ $Aug(\overline{x}^{'n'}) \end{pmatrix} \begin{pmatrix} I x_{1}^{'''} x_{2}^{'''} \cdots x_{d}^{'''} \\ \vdots \vdots \vdots \\ I x_{1}^{'''} x_{2}^{'''} \cdots x_{d}^{''''} \end{pmatrix}$ Classification US. Regression linear classifier from regressor: Reyrission —> number (lassification -> class lubel Convert ourput to 2-1,13 using $sign(z) = \begin{cases} 1 & 2 > 0 \\ -1 & 2 < 0 \\ 0 & otherwise. \end{cases}$

prediction: sign (H(Z)). Decision boundary is place to dassify the data $(w_0 + w, \varkappa, z) = 0$ $\begin{array}{c} (h-1) \quad himmsin \\ \hline \\ Perision \quad boundary \quad is \quad generally \quad a \quad hyperplane \\ \hline \\ \end{array}$ · W, + W, 2, + W, X, $\chi_{2} = -\frac{w_{3}}{w_{1}} + \frac{-w_{1}}{w_{2}} \pi_{1}$ (line) decison boundary herisely from MSE wo + W, N, + W2 N2 + W3 N3 affected by outliers $1 \qquad \chi_3 = \frac{-w_0}{w_3} + \frac{-w_1}{w_3} \chi_1^* + \frac{-w_2}{w_3} \chi_2$ (hyper plane). Feature Mappinz Learn new representation by creating new features from old feature. (non-linear pattern). XER^d. Non-linear Basis function. $\vec{P}(\vec{x}) = (\ell_1(\vec{x}), \ell_2(\vec{x}))$ Jeature map a new representation. $(\vec{z} = \vec{a}(\vec{x}))$ by mapping each training bata $\vec{z}^{i'}$ to feature space. fit linear produccion func. It in Cn) feature space. 1 m Hf(Z) = w, + w, Z, + ... + w& Zb. = $w_0 + w_1 \varphi_1(\tilde{x}) + \dots + w_d, \varphi_d, (\tilde{x})$ x choose polynomial x basis function.

 $\begin{array}{c} x \rightarrow (n, x^{*}, x^{3}, x^{4}) \\ 1 & 7 & 1 \\ e_{1} & e_{2} & e_{3} \end{array}$ now feature. model $H(\vec{x}) = \vec{w} \cdot Aug(\vec{z})$ where $\vec{z} = \begin{pmatrix} \varphi, (\vec{z}) \\ \vdots \\ \varphi_{\varphi}(\vec{z}) \end{bmatrix}$. nen design matrix; $X = \begin{pmatrix} 1 & \chi_1 & \chi_1^2 & \chi_1^3 & \chi_1^4 \\ 1 & \chi_2 & \chi_2^2 & \chi_2^3 & \chi_2^4 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \chi_n & \chi_n^4 & \chi_n^3 & \chi_n^4 \end{pmatrix}$ 50, $\vec{w} = (x^T x)^T (x^T \vec{y})$ Visualization Plot Decision boundary in Prediction Surface in Feature space Feature space H(Z) H(3)=0 H(8)<0 31 nonlinear Decision boundary in Prediction Surface in **Original** space $H(\vec{x}) = \omega_0 + \omega_1 \varphi_1(x)$ **Original** space twenking the surface of data via feature + W2 P2 (X) H(X) H(X)=0 decision boundary (non-linear!) H(x)>0 HUX) =0 mapping H(x)<0 X2 tent XI χ, Rabial Basis Function (RBF). A generic basis function that works for many problems: => Common choice: gaussian basis function: $\varphi(\vec{z};\vec{u},6) = e^{-||\vec{z}-\vec{u}||^2/6^*}$ center wilcoh

Contour $\mathcal{P}(\tilde{\mu}) = e^{-\sigma} = 1$ plot pl 2 $\mathcal{C}(\bar{\mu} \pm 6) = e^{-6^2/6^2} = e^{-1} = 0.369$ *U* 2 $i \left\{ \begin{array}{c} || \vec{k} - \vec{\mu} ||^2 > 2 \\ e \end{array} \right\}_{2} = e^{-\rho \sigma} = e^{-2\rho}$ 50 · if in close to pie, $\chi_1 | \mathcal{U}_1$ φ(z; μ, ε) is large.) -> intution. A · if \$ far from The, & measures how "similar" is to The. $\varphi(\vec{x};\vec{\mu},6)$ is small. (assume similar objects have close feature vectors). Procedure: # of new forcures : a' D pick pick centers, The" -- The withit: 6, -- 6 d' (usually All same) pick $(2 =) \text{ define ith basis function} \quad \mathcal{P}_i(\vec{x}) = e^{-||\vec{x} - \vec{\mu}|^{ij}||^2/6i^2}$ 3 => Last, we train a linear classifier in this new representation: $H(\bar{z}^{2}) = w_{0} + w_{1} + w_{1} + w_{2} + w_{3} + w_{4} + w_{4}$ example) Ans: Need 2 Gaussian basis func: 9, & 92 How many Gaussian basis functions would you use, goal: new feature -> linear prediction! and where would you place them to create a new representation for this data? small q, large q2 H(x)=0 H(y)=0 (32) H(3) <0 gaussian H(รี่)>ด φ, (x) (ξι) H(X) <0 H(x)=0 (Decision Boundary) place gaussian in the centor (mean) of

the data points (same class label have large "similarity".

H(52) is sum of Guussians: $H(\vec{x}) = w_{0} + w_{1} \varphi_{1}(\vec{x}) + w_{2} \varphi_{2}(\vec{x})$ = $w_{0} + w_{1} e^{-\|\vec{x} - \vec{k}\|^{2}/6^{2}} + w_{2} e^{-\|\vec{x} - \vec{k}\|^{2}/6^{2}}$ $= W_0 + W_1 e^{-\|\vec{x} - \vec{\mu}_{\infty}^{"}\|^2 / \sigma^2} + W_2 e^{-\|\vec{x} - \vec{\mu}_{\infty}^{"}\|^2 / \sigma^2}$ **Decision Boundary** By increasing # of basis functions, we can make more complex decision surfaces. RBF Neework Gaussians are ex) of radial basis func. each basis func has lenter 2. example RBF: inverse multiquaterate value depends only on discance from center. $\varphi(\vec{n},\vec{c}) = \frac{1}{\sqrt{1 + \xi^2 (\|\vec{n} \cdot \vec{c}\|\|^2)}}$ $\mathcal{D}\left(\overline{\lambda};\overline{c}\right) = f\left(\Vert \overline{\lambda} - \overline{c} \Vert\right)$ Function to 2 distance ben input 32 & center 2 measure distance. 12 BF Neework: O choose basis function P., ..., Pd. > Transform data to new representation: $\vec{\lambda} \mapsto \left(\varphi_{1}(\vec{\lambda}), \varphi_{2}(\vec{\lambda}), \cdots, \varphi_{d}(\vec{\lambda}) \right)'.$ 3 Train . $H(\vec{n}) = w_{\sigma} + w, \varphi, (\vec{x}) + \dots + w_{d'} \varphi_{d'}(\vec{n}).$ RBF necusik has these parameters.

· parameter of each individual basis function

· ju (conter) · b (variance) · weight associated to each "new" feature: W: Training: · find parameters of RBF-s first the is 1 through optimization, clustering, ranked My, After fixing those parameters · optimize w's $\vec{H}(\vec{x},\vec{u})$ which is linear. How to finde paraneters of RBFs? A pproaches O Everz data points as a center. n basis function (one for each point). 'n features $\vec{\phi}(\vec{x}) = \left(\phi_{1}(\vec{x}), \phi_{2}(\vec{x}), \cdots, \phi_{n}(\vec{x})\right)^{T}.$ problem: overfixting 4 computationally expensive. -> & (ndi) if d=n $=> \Theta(n^3)$ 2 Rankon Sample rankonly choose K data points as center problems: unhersample / oversample a region.

3 Clustering · group data points into clusters. · cluster centers => RBi-s. Inference : • given a point \overline{k} , map it to feature space $\overline{k} \mapsto (P, (\overline{k}) \cdots P_k (\overline{k}))^{T}$. · evaluated. If in feature space. $H_f(\vec{z})$ K-means clustering for picking parameters of RBFs. i hea : compress each clustering into a single point while minimizing information loss Given laca Exist ER and a parameter k. · Fink: K clusters center, IL' ... The so that average square distance from a data point to nearest cluseer conter is small nearest center $\sum_{i=1}^{2} \left(ose\left(\overline{\mu}_{i}^{(i)} \cdots \overline{\mu}_{i}^{(k)} \right) = \frac{1}{2} \sum_{i=1}^{2} \min_{i \in \{i, k\}} \left\| \overline{\lambda}_{i}^{(i)} - \overline{\mu}_{i}^{(i)} \right\|^{2}.$ optimize cost index of index of data cluster center. which cost is smaller? cost #1 or #a? 9 larger lose Smaller Cost x of cluster center x2 3 4 5 6 PetalLengthCm

Lloydes algorithm for K-means: 1) inicialize centers 1/2 ···· 1/2 Ð repeat until convergence. assign each point it to cloest center. uplicite each the as the mean of points assigned to it. k= a.
 ^Δ
 <sup>μ⁽²⁾
 ⁽²⁾
 </sup> • new # (27 mean location of μű assigned points for • new HC17 New Miles each center fil K=2 0 25 µ12) o دنتم How to choose 1 ? Increase le always decrease objective function J Elbow Mechok: with increasing values of k. K-means repeatobly · run the values of the objective as function of k. · plut fink elbor in ploc. elbow -> good value of k. obj= cost (µ(",..., µ(")) k= 2 K= 1 A. \bigwedge K= 4 k=3 k= # of data points ¥. <u>_</u> ¥ **4**0 **6**0 **6**0 **6**0 **6**0 * ¥ × *** X X

RBF network: · k-mpan cluceering to find center · create new features using & RBFs. · least square classifier. There is other parameter of RISI-s -) choose via cross validiation. Newal Neevolk Neurons are organized into layers. V inpue layers ouxpur layers hibben layers ошерис H(x) Notation: II) With neuron index in previou layer Iayer. $W'' = \begin{pmatrix} 2 & -1 & 0 \\ 4 & 5 & 2 \end{pmatrix} \in \mathbb{R}$ $\begin{array}{c} -3 & (1) \\ b & = (3 - 2 - 2) \\ b & = (3 - 2 - 2) \\ hidden \\ hidden \\ \psi \\ W^{(2)} = \begin{pmatrix} 3 \\ 2 \\ -4 \end{pmatrix} \\ \mathcal{E} \\ \mathcal{R}^{3 \times 1} \\ \mathcal{Output} \\ \mathcal{Output$ $\frac{-2}{b}$ = $(-4)^{7} \in R'$

· # of cells in input layer determined by dim. of input feature vector. · # of cells in hidden layer beterminch by you. · output layor can have >1 neurons. These are "fully-connected, feel-formach!" neewoods with one input. \bigvee They are function $H(\overline{k}): \mathbb{R}^{d} \to \mathbb{R}^{d}$ 5 forwarde pass : comput layer (i), use as input for layer it 1. · lee z'' be ourpur of note j in lager i. $0uzputs : Z = (Z_1, Z_2, \cdots)^T$ $\cdot \frac{1}{2} = [w^{(i)}]^{\frac{1}{2}} + b^{(i)}$ Each lager is a function: enj · H"/2) = [w"] 2 + 6" $H'' : R^{\perp} \rightarrow K^{3}$ $|-|(\bar{n})| = |+|^{(2)} (|+|^{(1)} (\bar{n}))$ $\cdot H^{(2)}(\vec{z}) = [w^{(2)}]^{\vec{z}} = \frac{-2}{2} + 6$ $= w^{(2)} \overline{w^{(1)}} \overline{z^{2}} + w^{(2)} \overline{b}^{(1)} + \overline{b}^{(2)}$ $= \widetilde{U}^{T} \overline{z}^{2} + \widetilde{b}$ $H' \cdot R^{2} - 2K'$ $H(\vec{z}) = [w''']'([w''']'\vec{z} + \vec{b}'') + \vec{b}'' = \vec{w} \cdot A_{uq}(\vec{z}).$ since NINI are lingher. Non-linearity -> Activation Function $\alpha_{j}^{(i)} = g(z_{j}^{(i)})$ the actual output of neurons. activation function 9(.)

 $a_{j}^{(i) \leftarrow index}$ cample of activation function: neuron Trodex (non-lmear NN) • $6(2) = \frac{1}{1+e^{-2}}$ sigmoid $\left(\begin{array}{c|c} z_1^{(1)} & a_1^{(1)} \\ \hline \end{array}\right)$ activations **x**₁ $\begin{bmatrix} z_2^{(1)} & a_2^{(1)} \\ \hline \end{array}$ $z_1^{(2)}$ $a_1^{(2)}$ 0 x₂ $\alpha_{\bar{J}}^{(i)} = 9(3_{\bar{J}}^{(i)})$ z₃⁽¹⁾ a₃⁽¹⁾ • $g(z) = max \{0, z\} = \{2, 2, 2, 0, z < 0,$ $rightarrow z_i^{(i)}$ is the linear activation before g is applied. $a_i^{(i)} = g(z_t^{(i)})$ is the actual output of the neuron. g(.): activation function (sometimes o(.)) be difference than Activation of ourpur activation neurons can of hidden neurons Interpretation from feature map: Ļ hidden layer of networks Carn feature map. a α neural The hidden la Interpretation: The hidden layers of a neural predicter network learn a feature map. We b output layer's parameters for prediction. W Duphate together. ÷Σ---feature X2 predictor: ÷Σ _ RAR for feature map para neces_ predictor: R2+R ÷(Σ)-È feature map: 12=12 Training : Empirical risk minization. · a training set $\left\{\left(\overline{\mathcal{X}}^{(i)}\right), \mathcal{J}_{i}\right\}$ · Initialize prediction function H.

 $\frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}$ $\nabla \vec{z} R(\vec{w}) = \frac{1}{2} \frac{2}{2H} \nabla \vec{z} H(\vec{z}; \vec{w}) = \frac{1}{2} \frac{2}{2H} \nabla \vec{z} H(\vec{z}; \vec{w}) = \frac{1}{2} \frac{1}{2H} \frac{1}$ What is TH? suppose 1-1 is neural not with incrivation $P(H) | - | (\bar{x}, \bar{w}) = 6(w^{R})^{T} 6(w^{''})^{T} 6(w^{'''}x + b^{''}) + b^{(2)} + b^{(3)}).$ Chain rule ? hote: often useful to pack all neights into a parameter vector 23. Aside: herivative of RELU. $\overline{w} = \left(w_{i_{1}}^{(i)}, w_{i_{2}}^{(i)}, \dots, b_{i_{3}}^{(i)}, b_{2}^{(i)}, w_{i_{2}}^{(2)}, \dots, b_{i_{3}}^{(2)}, b_{2}^{(2)}, \dots, b_{i_{2}}^{(2)}, \dots, b_{i_{3}}^{(2)}, b_{2}^{(2)}, \dots, b_{i_{3}}^{(2)}, b_{2}^{(2)}, \dots, b_{i_{3}}^{(2)}, \dots, b_{i_{3}$ 0 g(z) = max {0, z} $g'(z) = \begin{cases} 0 & , z \leq 0 \\ 1 & , z > 0 \end{cases}$ chain rule: f(g(x)), $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$, $e_{\mathcal{H}}) |-| = a_1^{(2)} = f(z_1^{(2)})$ $Z_{i}^{(2)} = W_{i1}^{(2)} \cdot a_{1}^{(\prime)} + W_{2i}^{(2)} \cdot a_{2}^{(\prime)}$ $= W_{11}^{(2)} \mathcal{J} \left(W_{11}^{(1)} \times_{1} + W_{21}^{(1)} \times_{2} + W_{31}^{(1)} \times_{3} \right) +$ $W_{21}^{(2)} g(W_{12}^{(1)} \chi_{1} + W_{22}^{(1)} \chi_{2} + W_{32}^{(1)} \chi_{3})$ $\frac{\partial H}{\partial w_{ii}^{(1)}} = g'(z_{i}^{(2)}) \frac{\partial z_{i}^{(2)}}{\partial w_{ij}^{(1)}} \left(w_{ii}^{(2)} g'(z_{i}^{(1)}) \cdot z_{i} \right)$



Claim #1 Claim #2 Claim #3 $\frac{\partial H}{\partial z_j^{(\ell)}} = \frac{\partial H}{\partial a_j^{(\ell)}} \underbrace{g'(z_j^{\ell})'}_{\ell}$ $\frac{\partial H}{\partial W_{ij}^{(\ell)}} = \begin{bmatrix} \frac{\partial H}{\partial z_j^{(\ell)}} a_i^{(\ell-1)} \\ \vdots \\ \vdots \end{bmatrix}$ $\frac{\partial H}{\partial a_j^{(\ell)}} = \sum_{k=1}^{n_{\ell+1}} \boxed{\frac{\partial H}{\partial z_k^{(\ell+1)}}} W_{jk}^{(\ell+1)}$ $H = \phi_{1}(3^{(l+1)}) + \phi_{2}(3^{(l+1)}) + \cdots + \phi_{n_{l+1}}(3^{(n_{l+1})}) + W^{(l+1)}_{j_{2}}(3^{(l+1)}) + W^{(l+1)}_{j_{2}$ $\frac{\partial H}{\partial H} = \frac{\partial H}{\partial \chi_{k+1}^{(2+1)}} = \frac{\partial H}{\partial \chi_{k+1}^{(2+1)}} \frac{\partial \chi_{k+1}^{(2+1)}}{\partial \chi_{k+1}^{(2+1)}} + \frac{\partial H}{\partial \chi_{k+1}^{(2+1)}} \frac{\partial \chi_{k+1}^{(2+1)}}{\partial \chi_{k+1}^{(2+1)}} + \frac{\partial H}{\partial \chi_{k+1}^{$ $\begin{array}{c} z_{i}^{(q)} & a_{i}^{(q)} \\ b_{j}^{(\ell)} & W_{jn_{\ell+1}}^{(\ell+1)} & z_{n_{\ell+1}}^{(\ell+1)} \end{array}$ $=g'(3_{\bar{J}}^{(2)})$ Now, how to calculate $\frac{\partial H}{\partial a_{j}^{cb}}$? · Derivative in layer I dependes on dierivatives in layer 12+17. recursive formulos: · yiven input n and a current parameter vector ~? · evaluate the network to compute Z: and a; for all notices. · for each layer I from lose to firse. $\mathcal{L}_{\mathcal{A}_{i}}^{\mathcal{H}} = \sum_{k=1}^{n_{\ell+1}} \frac{\mathcal{L}_{i}}{\mathcal{L}_{k}}^{\mathcal{H}} \mathcal{W}_{jk} \qquad from last layer/next layer.$ $\int_{2Z_{j}}^{2H} \frac{\partial H}{\partial u_{j}} = \frac{\partial H}{\partial u_{j}} g'(Z_{j}^{(e)})$ $\frac{\partial H}{\partial w_{ij}}(\ell) = \frac{\partial H}{\partial z_{ij}}(\ell) = \frac{\partial H}{\partial z_{ij}}(\ell)$ $\cdot \frac{2H}{2b_{j}} = \frac{2H}{2Z_{j}}$

example).

Compute the entries of the gradient given: seepse no bias

$$M^{(1)} = \begin{pmatrix} 2 & -3 \\ 2 & 1 \end{pmatrix} W^{(2)} = \begin{pmatrix} 2 & -1 \\ 0 & 1 \end{pmatrix} W^{(2)} = \begin{pmatrix} 3 \\ -2 \end{pmatrix} \tilde{x} = (2, 1)^{2} g(2) = ReLU$$

$$Q = \frac{2H}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 1 = 1$$

$$\begin{cases} 2 & \frac{9}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 1 = 1 \\ \frac{1}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 1 = 2 \\ \frac{1}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 1 = 2 \\ \frac{1}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 2 \\ \frac{1}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = 2 \\ \frac{2H}{2g_{2}^{2}} = \frac{2H}{2g_{2}^{2}} g(2g_{2}^{2})^{2} = \frac{2H}{2g_{2}^{2}} g(2g_{2})^{2} = \frac{2H}{2g_{2}^{2}} g(2g_{2})$$

i dea: f(g(n))inpur ouepur f(g(n)) x = g = f $goal: \frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$ Backpropagation via gradient descent. gradient of loss W.r.t. the weights. by iterative optimization. $f: R^{a} \rightarrow R^{a}$ gradience of faz Z $\vec{z}f(\vec{x}) = \left(\frac{\partial f}{\partial x_1}(\vec{x}) \cdots \frac{\partial f}{\partial x_k}(\vec{x})\right)^T$ $\vec{x} = \begin{bmatrix} x \\ \vdots \\ x \\ x \\ \end{bmatrix} \in \mathbb{R}^{4}$ also a function mapping from idea from taylor expansion: Rd -> Rd. $f(x_0 + 2n, y_0 + 2y) \approx f(x_0, y_0) + \overline{Z} \cdot \overline{\nabla} f(x_0, y_0).$ 0 2 3 $0 < 2 \implies 3$ should be negative: $= \overline{z}^2 \cdot \overline{z}^2 f(x_0, y_0) < 0.$ which Z can we decreese the value of @ the mose? $= \frac{1}{2} = \frac{1}{2} = \frac{1}{2} + \frac{$ · direction prehogonal to \$f(x), f does not change!

So gratiene descene as follow: pick learning race n70. repeat until convergence. $\vec{x}^{(i+1)} = \vec{x}^{(i)} - n \vec{z} f(\vec{x}^{(i)}).$ steepest descent direction. when $\| \overline{\mathcal{X}}^{(i)} - \overline{\mathcal{X}}^{(i+1)} \|$ is small. stop So, in Neural Network Training. consilier square loss: $\nabla_{\vec{\omega}} R(\vec{\omega}) = \frac{1}{n} \sum_{i=1}^{2} \frac{\partial \ell}{\partial H} \nabla_{\vec{\omega}} H(\vec{n}); \vec{\omega}.$ $= \frac{2}{n} \sum_{i=1}^{n} (H(\vec{x}^{(i)}) - \vec{y}_{i}) \nabla_{\vec{u}} H(\vec{x}^{(i)}, \vec{u}))$ continue backprop the errors through the layers. Difficulties of training NN: . when activation are non-linear, risk of NIN. is highly non-convex. · vanishing gradient problems grabile e of neight the at earlier layers can be very small rosult the update of n is very slow by gradient discent." one mitigation: use ReLU instead of signoid: 1 gratiene 1 accumulucere be by nultiplication

· computational lose is high.

possible solution Batch learning for barch size m << n the # diata points.

compute stockastic gradient based on the batch of data points (m).

Convolucion Neural Neurol (image recognizion by line detector exampk) shape lictector 0 0 255 face detection 0 0 255 0 0 255 (*) => result Convolution with an edge filter. 0 -1 1 (255×3) (detece edges in the image) 0 -1 1 vertical / horizoneal / diagonal. 0 -1 1 $\uparrow \tilde{j}$ detece transition (from no elize to elize) example) move the filter over the entire image, repeat procedure. 000000 20 - activation map. 0 -1 1 0 0 0.9 0 0 0.7 0 -1 1 20 - convolution of filter with the image. 0 0 0.3 0 0 0.9 0-111 000,7000

Typically, 3×3 or 5×5. - hyperparameter in conduction. Variations: different strike, image pakking (Stribe = 1 => move 1 pixel) controls how far to move Filter muse have 3 channels: en). 3×3×3 filter (same channel as image) oueput still 20 3 channels CNN use convolution in early layers to create new feature representation. filters are learned 3 flatten 2 hidden hidden layer#1 layer #2 . OUT PUT ENPUT filter: parameter filter: parameter FC Inyer: Parameter #3 k filters will result in k chanels in the next layer. ① Input Convolutional Layer activation func. applich entrywise. filter: 3x3×1, 4filters Ck

3×3×4 Second Convolutional Layer 2)` muse march the channel of input K INPUT OUTPUT filter: 3×3×4, K2 oneque -ill be 3 channels. last lager to fully connected. 3 Flattening: stretch the te May proling : OUTPUT Response map: 6×6 OUTPUT convolution 3x3: 4x4 choose the nax 3) Wij (3) 48×12 max ουτρυτ (2×2) pooling 3×3 Auto encoler representation learning is finks on encoding function: encode (x): R > R k · captures useful aspect of data distribution. · encode can decrease dimensionality. preserve as much info about in as possible.

by enoking, we hope to becoke and reconstruct original bater. $\mathcal{R}_{cconstruction} = \operatorname{rror} : \sum_{i=1}^{n} || \overline{\mathcal{R}}^{ij} - \operatorname{lecobe}(\operatorname{encobe}(\overline{\mathcal{R}}^{ij})) ||^{n}$ Trivial solution. encode $(\vec{n}) = \vec{n} = decode(\vec{n})$, which is not use ful. (rample) PCA encolie by projecting onto top k eizenvectors. $encolac(\vec{n}) = U^{T}\vec{n}$, where columns of V be top k eigenvators of $becobe(\vec{z}) = V\vec{z}$ covariance matix. Encober/decober as Neural Ner devoke (encoke (z2)) as a NN. (non-linear activation). R^d -> 12^k -> R^d input output During training, minimiz reconstruction error. $\frac{1}{2} \|\vec{x}^{(i)} - H(\vec{z}^{(i)})\|^{2} = \sum_{j=1}^{n} \frac{d}{2} (\vec{x}^{(i)}_{j} - (H(\vec{z}^{(i)}))_{j})^{n}$

Autoencolier is generalization of PCA

such that performs orthogonal projection. 50 PCA minimizes reconstruction error $encouver(\vec{x}) = \sqrt{\vec{x}}$ subject to conservine these columns of U are or the normal. decode $(\overline{z}) = U\overline{z}$ However, aucoensker learns (non-orthogonal) projection into sume space of PC17 Use of Auroencober: · dimensionality reduction (non-linear) (k < d.) · denoising auto encoller (k>4). greater dimension which rankow noise to each \$2'' to get \$2''. erain IVIV S.Z. $H(\hat{x}^{(i)}) \approx \overline{x}^2$.

(ausclith

Association: X & Y are associated iff $\exists x_1 \neq x_2, P(Y|X=x_1) \neq P(Y|X=x_2)$ $(aus_n litg: \exists x, \neq n_x, P(Y | lo(X=n,)) \neq P(Y | lo(X=n_z))$ X is manipulated/intervened to x, "bo" operation. Directed graph can represente causal relationship. Genotype MIL can be booseek Smoking Lung cancer and nore incorprecable if we finth coustility_ smoking (from "Causation, Prediction, and Search" by SGS) 3 Questions in AI. rined quebelono in current 11. X_{l} × 3 (0 42 h. • Association: Would the person cough if we find he/she has yellow fingers? Yellow f.o.g.ris $P(X3 \mid X2=1)$ Seeing Yellow fingers С • Intervention: Would the person cough if we only relas ML make sure that he/she has yellow fingers? association with no alrect $P(X3 \mid do(X2=1))$ Doing/Intervening **Causality is requi** • Counterfactual: Would George cough had Causalizz he had yellow fingers, given that he does not have yellow fingers and coughs? $P(X3_{X2=1} | X2 = 0, X3 = 1)$ Imagining/Retrospection/ Creativity

causal discovera & informate. (discover causal relationship & assimate causal effects from observational decree).

Causal discovery in complex scenarios:

- Nonlinearities; Mixed Continuous & Discrete Variables
- Latent Confounding; Latent Causal Representation Learning
- Subsampling / Temporally Aggregation in time series
- Measurement Error; Selection Bias; Missing Values ٠
- Nonstationarity/Heterogeneity



Big 5: openness (O); conscientiousness (C); extraversion (E); agreeableness (A); neuroticism (N)

