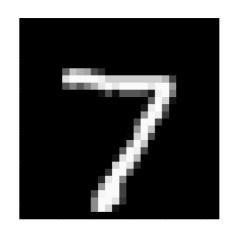
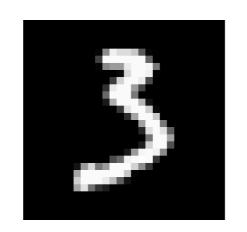
#### CS 453X: Class 16

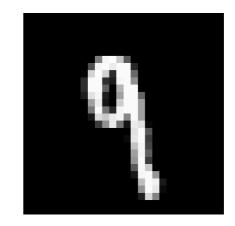
Jacob Whitehill

- Prior to choosing a particular ML model, it can sometimes be helpful to visualize your dataset.
- One of the most commonly used visualization techniques is called **principal component analysis** (**PCA**).

- Consider the MNIST dataset of hand-written digits.
- Visualizing each individual example is easy, e.g.:







 But this doesn't tell us a whole lot about the dataset as a whole, or how separable the classes are.

- Each MNIST image is 28x28 pixels => 784 dimensions.
- To show all examples in the original input space, we would need a 784-dimensional visualization.
  - But humans struggle with perception beyond 3-D.
- How can we represent a collection of many highdimensional images in just 2-D or 3-D?

- We somehow have to condense the interesting information of a 784-dim vector into just 2-3 dimensions!
- Core question:
  - How do we pick the dimensions or more generally, directions — of data to present?

- Let's give a naive try...
- For each image  $\mathbf{x}^{(i)}$  (where i=1, ..., n):

- Let's give a naive try...
- For each image  $\mathbf{x}^{(i)}$  (where i=1, ..., n):
  - Retrieve the values of just the first two pixels, i.e.,  $(\mathbf{x}^{(i)}_1, \mathbf{x}_2^{(i)})$ .

(0,0)	(1,0)	(2,0)	 (27,0)
(0,1)			
(0,2)			
(0,27)			(27,27)

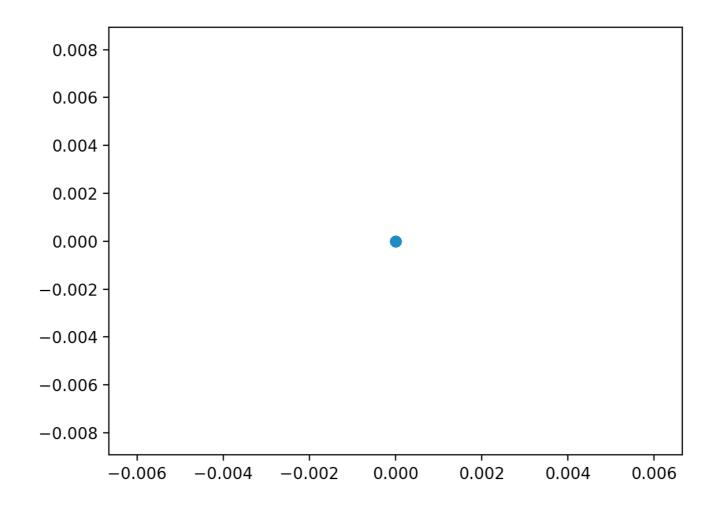
- Let's give a naive try...
- For each image  $\mathbf{x}^{(i)}$  (where i=1, ..., n):
  - Retrieve the values of just the first two pixels, i.e.,  $(\mathbf{x}^{(i)}_1, \mathbf{x}_2^{(i)})$ .
  - Plot the point  $(\mathbf{x}^{(i)}_1, \mathbf{x}_2^{(i)})$  in 2-D space.

## MNIST 2-D visualization: first two pixel values

 Let's apply this procedure to 2500 images from the MNIST test set...

## MNIST 2-D visualization: first two pixel values

 Let's apply this procedure to 2500 images from the MNIST test set...

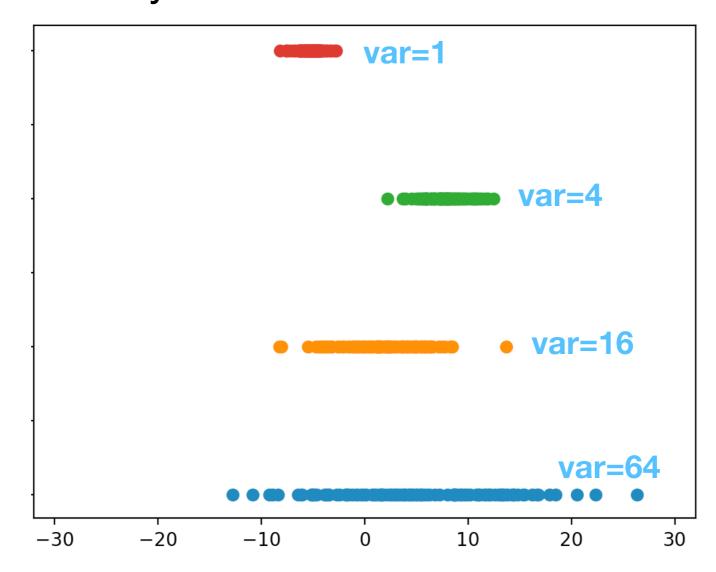


What happened?

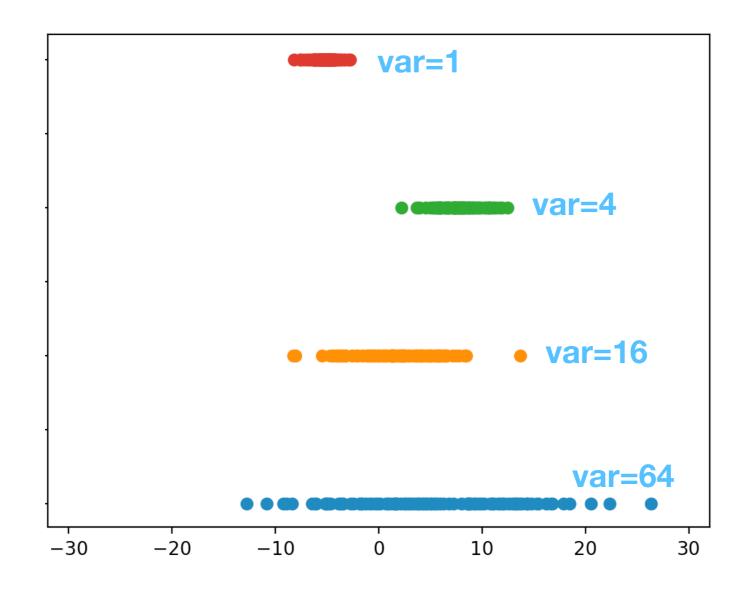
### MNIST 2-D visualization: first two pixel values

- The problem is that, for all *n* images in our dataset, the value of the first two pixels is 0!
- There was very little (actually, 0) variance across each of these two dimensions.

 Intuitively, the variance of a vector of n numbers is how "spread out" they are:



Note that the variance is independent of the mean!



- From basic statistics:
  - **Mean** of an *n*-dimensional vector: sum and divide by *n*:

$$\mathbb{E}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{p}_i$$

- From basic statistics:
  - **Mean** of an *n*-dimensional vector: sum and divide by *n*:

$$\mathbb{E}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{p}_i$$

 Variance of an n-dimensional vector: the mean squared distance from the mean:

$$\mathbb{V}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{p}_i - \mathbb{E}(\mathbf{p}))^2$$

- Alternatively, we can compute the variance in two steps:
  - First subtract the mean from each element of p:

$$ilde{\mathbf{p}} = \mathbf{p} - \mathbb{E}(\mathbf{p}) \left[ egin{array}{c} 1 \ dots \ 1 \end{array} 
ight]$$

- Alternatively, we can compute the variance in two steps:
  - First subtract the mean from each element of p:

$$\tilde{\mathbf{p}} = \mathbf{p} - \mathbb{E}(\mathbf{p}) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

• Since the mean of  $\tilde{\mathbf{p}}$  is 0, we can compute its variance as:  $\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{N} \sum_{\tilde{\mathbf{p}} \in \mathbb{F}(\tilde{\mathbf{p}})}^{n} \mathbb{P}(\tilde{\mathbf{p}})^{2}$ 

$$\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - \mathbb{E}(\tilde{\mathbf{p}}))^2$$

- Alternatively, we can compute the variance in two steps:
  - First subtract the mean from each element of p:

$$\tilde{\mathbf{p}} = \mathbf{p} - \mathbb{E}(\mathbf{p}) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

• Since the mean of  $\tilde{\mathbf{p}}$  is 0, we can compute its variance as:  $\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - \mathbb{E}(\tilde{\mathbf{p}}))^2$ 

$$= \frac{1}{n} \sum_{i=1}^{n} (\mathbf{\tilde{p}}_i - \mathbf{\tilde{p}}_i)^2$$

- Alternatively, we can compute the variance in two steps:
  - First subtract the mean from each element of p:

$$\tilde{\mathbf{p}} = \mathbf{p} - \mathbb{E}(\mathbf{p}) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

• Since the mean of  $\tilde{\mathbf{p}}$  is 0, we can compute its variance as:  $\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{\tilde{\mathbf{p}}} \mathbb{E}(\tilde{\mathbf{p}})^2$ 

$$\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - \mathbb{E}(\tilde{\mathbf{p}}))^2$$
$$= \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - 0)^2$$
$$= \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i)^2$$

- Alternatively, we can compute the variance in two steps:
  - First subtract the mean from each element of p:

$$\tilde{\mathbf{p}} = \mathbf{p} - \mathbb{E}(\mathbf{p}) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

• Since the mean of  $\tilde{\mathbf{p}}$  is 0, we can compute its variance as:  $\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{k=1}^{n} (\tilde{\mathbf{p}}_{k} - \mathbb{E}(\tilde{\mathbf{p}}))^{2}$ 

$$\mathbb{V}(\tilde{\mathbf{p}}) = \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - \mathbb{E}(\tilde{\mathbf{p}}))^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i - 0)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\tilde{\mathbf{p}}_i)^2$$

$$= \frac{1}{n} \tilde{\mathbf{p}}^\top \tilde{\mathbf{p}}$$

- Let's search for the two pixel dimensions along which the images vary the most.
- Selecting a particular pixel from each image x is equivalent to projecting each x onto a unit vector.

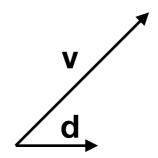
$$\mathbf{x}^{ op}\mathbf{d}$$

where **d** is a vector of *n*-1 zeros and 1 one (whose location corresponds to the particular pixel location):

$$\mathbf{d} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

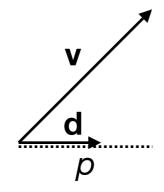
### Projection

 Recall that a projection of a vector v onto a direction (unit vector) d is given by p = v<sup>T</sup>d:



### Projection

 Recall that a (scalar) projection of a vector v onto a direction (unit vector) d is given by p = v<sup>T</sup>d:



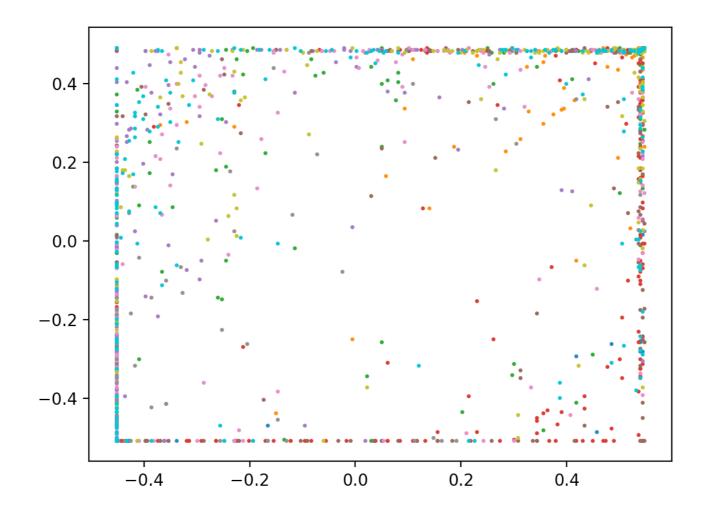
The scalar projection p measures the distance of v along
 d.

For each possible pixel dimension d, we can obtain the n-vector of all scalar projections (over all n images) as:

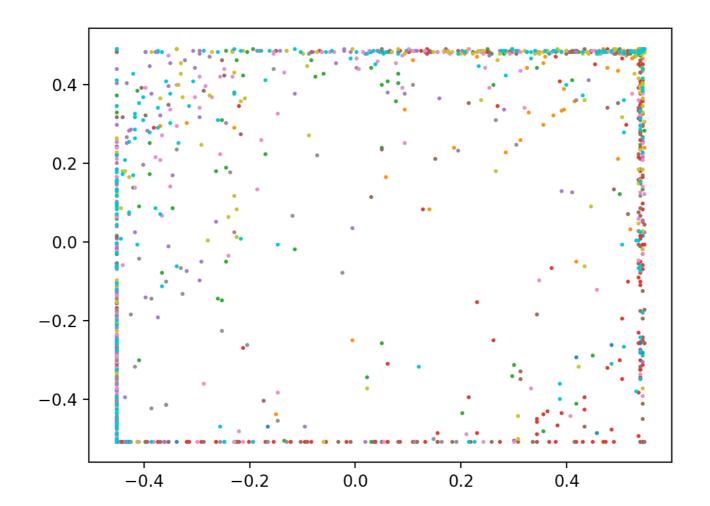
$$\mathbf{p} = \mathbf{X}^{\top} \mathbf{d}$$

- We can then calculate the variance of the scalar projections by calculating Var(p).
- By searching over all 784 possible dimensions, we can find the two dimensions along which variance is maximized.

- When we apply this to MNIST, we get:
  - Pixel dimension of 1st-highest variance: (r,c) = (13,14)
  - Pixel dimension of 2nd-highest variance: (r,c) = (14,14)

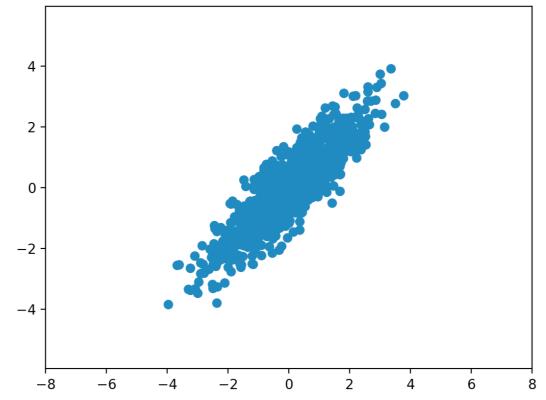


- Certainly much better!
- However, many of the values overlap each other due to saturation — in many images, these pixels' values are maximized/minimized.



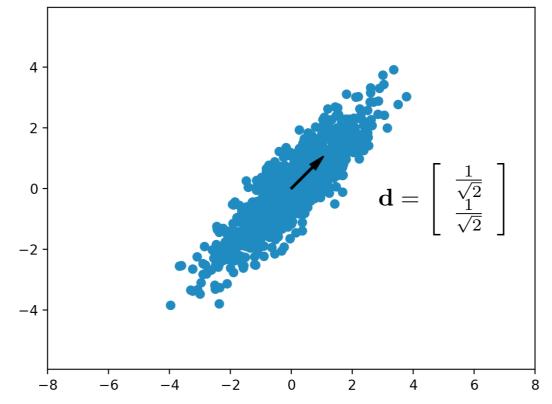
### Beyond axis-aligned directions

- But why constrain ourselves to only axis-aligned unit vectors, i.e., vectors with m-1 zeros and 1 one?
- Consider the following dataset in which each image contains just 2 pixels.
- Along which direction d is variance maximized?



### Beyond axis-aligned directions

- But why constrain ourselves to only axis-aligned unit vectors, i.e., vectors with m-1 zeros and 1 one?
- Consider the following dataset in which each image contains just 2 pixels.
- Along which direction d is variance maximized?



- In general, given a dataset  $\mathbf{X}$  of training examples, we want to find the direction d that maximizes  $Var(\mathbf{X}^{\mathsf{T}}\mathbf{d})$ .
- For simplicity, let's assume that the mean of X, along each pixel dimension j, is 0, i.e.:

$$\frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_j^{(i)} = 0$$

 (If this is not the case, then just subtract off the mean vector from each example x<sup>(i)</sup>.)

$$\mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}^{\top}\mathbf{d} \quad \text{by applying the definition of mean.}$$

$$\begin{split} \mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) &= \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}^{(i)}^{\top}\mathbf{d} \\ &= \frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{m}\mathbf{x}_{j}^{(i)}\mathbf{d}_{j} \quad \text{by applying the definition} \\ &\text{of inner product.} \end{split}$$

$$\begin{split} \mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) &= \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}^{(i)}^{\top}\mathbf{d} \\ &= \frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{m}\mathbf{x}^{(i)}_{j}\mathbf{d}_{j} \\ &= \frac{1}{n}\sum_{j=1}^{m}\sum_{i=1}^{n}\mathbf{x}^{(i)}_{j}\mathbf{d}_{j} \quad \text{since we can swap the order of the summations.} \end{split}$$

$$\mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}^{\top}\mathbf{d}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} \mathbf{x}_{j}^{(i)}\mathbf{d}_{j}$$

$$= \frac{1}{n} \sum_{j=1}^{m} \sum_{i=1}^{n} \mathbf{x}_{j}^{(i)}\mathbf{d}_{j}$$

$$= \frac{1}{n} \sum_{i=1}^{m} \mathbf{d}_{j} \sum_{i=1}^{n} \mathbf{x}_{j}^{(i)}$$

$$\begin{split} \mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) &= \frac{1}{n}\sum_{i=1}^{n}\mathbf{x}^{(i)}^{\top}\mathbf{d} \\ &= \frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{m}\mathbf{x}^{(i)}_{j}\mathbf{d}_{j} \\ &= \frac{1}{n}\sum_{j=1}^{m}\sum_{i=1}^{n}\mathbf{x}^{(i)}_{j}\mathbf{d}_{j} \\ &= \frac{1}{n}\sum_{j=1}^{m}\mathbf{d}_{j}\sum_{i=1}^{n}\mathbf{x}^{(i)}_{j} \\ &= \frac{1}{n}\sum_{j=1}^{m}\mathbf{d}_{j} \underbrace{\sum_{i=1}^{n}\mathbf{x}^{(i)}_{j}}_{\text{since X has zero mean}} \\ &= \frac{1}{n}\sum_{j=1}^{m}\mathbf{d}_{j} \underbrace{0}_{j} \end{split}$$

$$\mathbb{E}\left(\mathbf{X}^{\top}\mathbf{d}\right) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}^{\top}\mathbf{d}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} \mathbf{x}_{j}^{(i)}\mathbf{d}_{j}$$

$$= \frac{1}{n} \sum_{j=1}^{m} \sum_{i=1}^{n} \mathbf{x}_{j}^{(i)}\mathbf{d}_{j}$$

$$= \frac{1}{n} \sum_{j=1}^{m} \mathbf{d}_{j} \sum_{i=1}^{n} \mathbf{x}_{j}^{(i)}$$

$$= \frac{1}{n} \sum_{j=1}^{m} \mathbf{d}_{j}0$$

$$= \mathbf{0}$$

Therefore, the variance of X<sup>T</sup>d (for any d) is just:

$$\frac{1}{n} (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d})$$

• We thus want to find the d that maximizes:

$$(\mathbf{X}^{\mathsf{T}}\mathbf{d})^{\mathsf{T}}(\mathbf{X}^{\mathsf{T}}\mathbf{d})$$

subject to the constraint that d is a unit vector, i.e.:

$$\mathbf{d}^{\mathsf{T}}\mathbf{d} = 1$$

$$L(\mathbf{d}, \alpha) = (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d}) - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$
Objective
Constraint

$$L(\mathbf{d}, \alpha) = (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d}) - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$
$$= \mathbf{d}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{d} - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$

$$L(\mathbf{d}, \alpha) = (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d}) - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$
$$= \mathbf{d}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{d} - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$
$$\frac{\partial L}{\partial \mathbf{d}} = 2\mathbf{X} \mathbf{X}^{\top} \mathbf{d} - 2\alpha \mathbf{d} = 0$$

$$L(\mathbf{d}, \alpha) = (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d}) - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$

$$= \mathbf{d}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{d} - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$

$$\frac{\partial L}{\partial \mathbf{d}} = 2\mathbf{X} \mathbf{X}^{\top} \mathbf{d} - 2\alpha \mathbf{d} = 0$$

$$\Rightarrow \mathbf{X} \mathbf{X}^{\top} \mathbf{d} = \alpha \mathbf{d}$$

$$L(\mathbf{d}, \alpha) = (\mathbf{X}^{\top} \mathbf{d})^{\top} (\mathbf{X}^{\top} \mathbf{d}) - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$

$$= \mathbf{d}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{d} - \alpha (\mathbf{d}^{\top} \mathbf{d} - 1)$$

$$\frac{\partial L}{\partial \mathbf{d}} = 2\mathbf{X} \mathbf{X}^{\top} \mathbf{d} - 2\alpha \mathbf{d} = 0$$

$$\Rightarrow \mathbf{X} \mathbf{X}^{\top} \mathbf{d} = \alpha \mathbf{d}$$

- In other words, d is an eigenvector of XX<sup>T</sup>.
- Since we want to maximize the variance of the projections, we want the eigenvector with largest associated eigenvalue.

### Eigenvector

An eigenvector v of a (square) matrix A satisfies:

$$\mathbf{A}\mathbf{v} = \alpha\mathbf{v}$$

for some scalar eigenvalue  $\alpha$ .

- For an n x n matrix A, there are n eigenvectors v and associated eigenvalues α.
- Eigenvectors/eigenvalues can be computed (in  $O(n^3)$  time) using many standard linear algebra libraries (e.g., numpy).

### Positive semi-definite matrices

- Recall that the eigenvalues of every PSD matrix A are always non-negative.
- Since XX<sup>T</sup> is PSD (as shown previously in class), its eigenvalues are non-negative.

#### PCA

- The direction **d** along which the dataset **X** varies the most is the **principal eigenvector** of **XX**<sup>T</sup>, i.e., the eigenvector with largest associated eigenvalue.
- It can be shown that the second-most highly varying direction of X is the eigenvector with second-largest associated value, etc.

#### PCA

- Algorithm:
  - 1. From design matrix  $\mathbf{X}$ , compute the mean vector  $\overline{\mathbf{x}}$ :

$$\overline{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

3. Subtract  $\overline{\mathbf{x}}$  from each example  $\mathbf{x}^{(i)}$ , and then form matrix  $\widetilde{\mathbf{X}}$  (same size as  $\mathbf{X}$ ), which should have a mean (over all n examples) of 0 along each dimension j.

$$\widetilde{\mathbf{X}} = \left[ \begin{array}{c|c} | & | & | \\ |\mathbf{x}^{(1)} - \overline{\mathbf{x}}) & \dots & |\mathbf{x}^{(n)} - \overline{\mathbf{x}}) \end{array} \right]$$

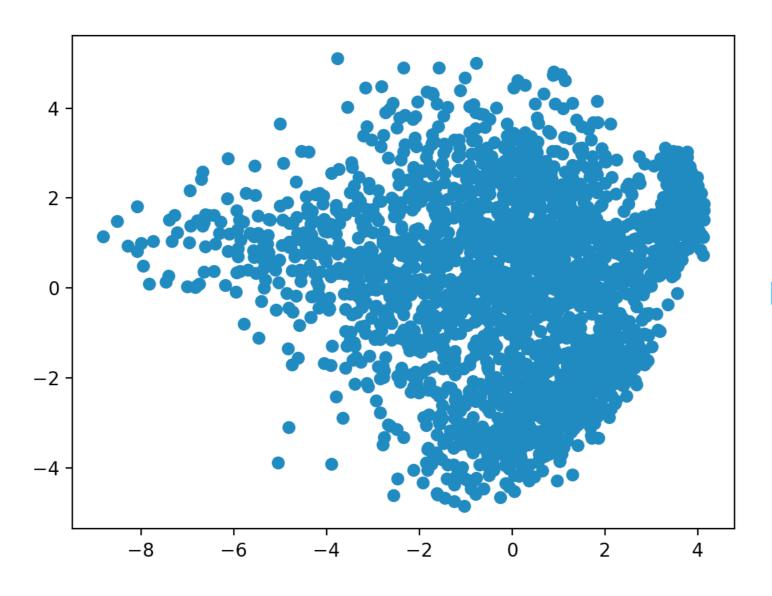
#### PCA

- Algorithm:
  - 1. From design matrix  $\mathbf{X}$ , compute the mean vector  $\overline{\mathbf{x}}$ :

$$\overline{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}^{(i)}$$

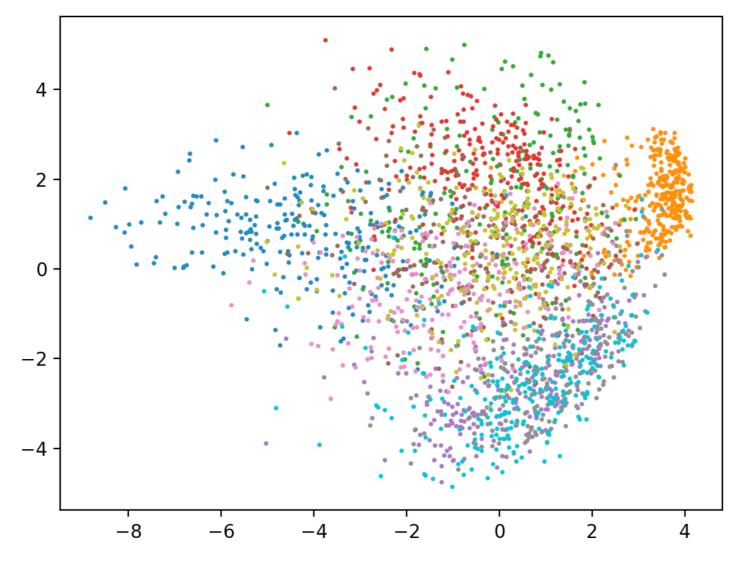
- 3. Subtract  $\overline{\mathbf{x}}$  from each example  $\mathbf{x}^{(i)}$ , and then form matrix  $\widetilde{\mathbf{X}}$  (same size as  $\mathbf{X}$ ), which should have a mean (over all n examples) of 0 along each dimension j.
- 4. Compute the eigenvectors & eigenvalues of XXT.
- 5. The  $k^{th}$  principal component (**PC**) of **X** is the eigenvector **v** of  $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{T}$  with the  $k^{th}$ -largest eigenvalue.

#### PCA on MNIST



For all classes

#### PCA on MNIST



For each class with its own color

### Unsupervised learning

- PCA is an example of an unsupervised machine learning algorithm.
- Unsupervised we never looked at the training labels!
  - In some settings, the data might not even be labeled.

### Unsupervised learning

- PCA is an example of an unsupervised machine learning algorithm.
- Unsupervised we never looked at the training labels!
  - In some settings, the data might not even be labeled.
- Note that there are other visualization methods (e.g., Linear Discriminant Analysis (LDA)) that are supervised:
  - Project data onto directions that best linearly separate the data classes.