



NLP and Dimensionality Reduction



Puya Vahabi

2018-August

Reminder

251.

~ 251,

251.
251.

- ▶ Results of the survey on syllabus and course complexity.
- ▶ Survey on what to do as last lesson.
- ▶ How is it going with the project?
- ▶ Each action is graded, but its voluntary.

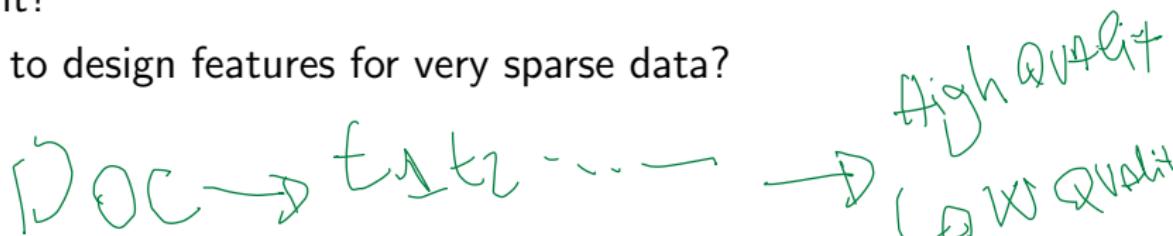
Sparse Data

$$X^{(i)} = \begin{bmatrix} \text{Age} \\ \text{height} \end{bmatrix}$$

Feature engineering: it is the process of feature types later used for classification purpose.

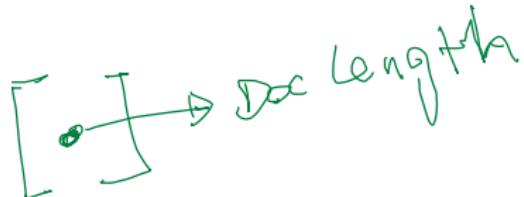
- ▶ Text Data. How would you extract features/attributes to be used in learning phase? #words? sequence of words? #phrases?
- ▶ Image Data. Pixels of a certain color? Objects inside of it?

How to design features for very sparse data?



Feature Extraction

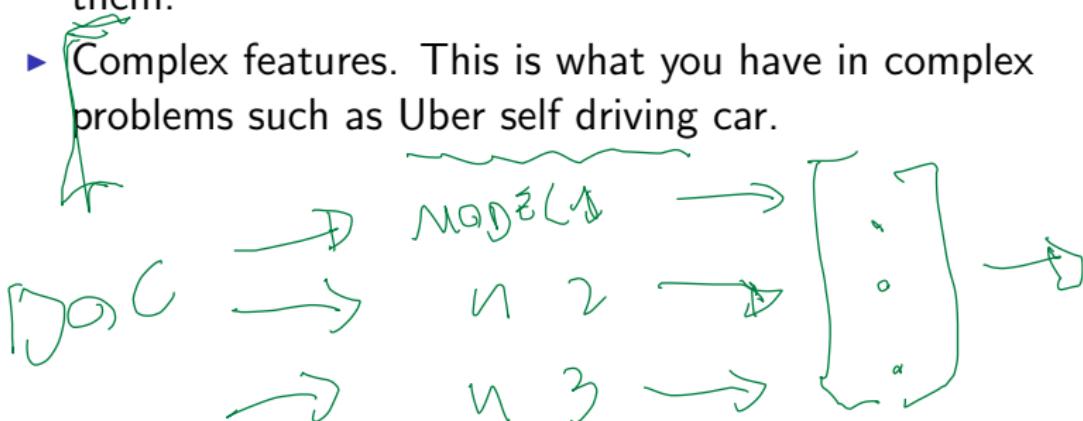
Doc



Convert the data into vector of features.

Consider *efficiency* of extracting features:

- ▶ For each instance on training and test you need to extract them.
- ▶ Complex features. This is what you have in complex problems such as Uber self driving car.



Feature Selection

Select a set of features that are important from the initial set of features.

- ▶ Can make the learning process faster.
- ▶ Can reduce the complexity of the learned model.

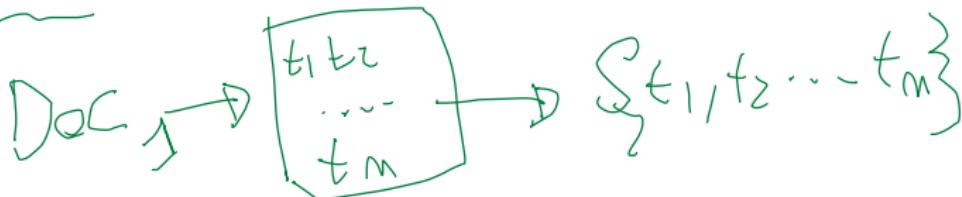
Dimensionality reduction is often used to have a reduced feature space.

Feature Extraction

Text Data

Bag of Words

"New Album RELEASED"



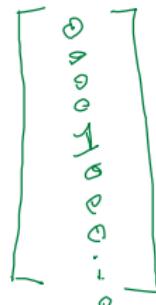
The most common feature encoding for text is called a **bag of words** representation

- ▶ Words are represented as sets, their order or position is not important.
 - ▶ If “Michael” is at the 2nd word in the document is treated the same as is in any other position in the doc.

Definitions

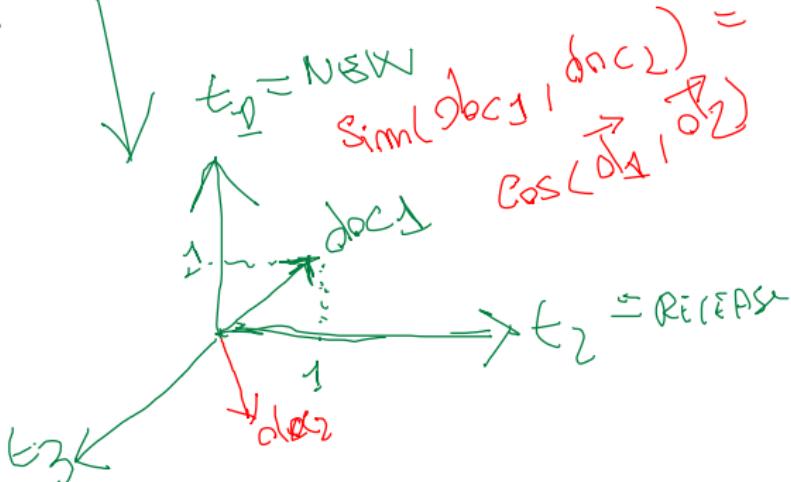
bs

- ▶ 1-Hot-Vector
- ▶ Term Vector Space
- ▶ Document as vector
- ▶ Bag of Words



$\# \text{UNIQUE TERMS}$
in your dataset

Doc →



Notes

BOW as features

$$\vec{o}_3 = \begin{bmatrix} 0 \\ 15 \\ 5 \\ 1 \\ 0 \end{bmatrix} \rightarrow t_3 \quad \rightarrow \text{Classification}$$

BOW. The features are the unique words that appears in vocabulary.

Can they have a weight?

- ▶ #times the word is appearing in a particular document.
- ▶ 0 means the word is not in the document
- ▶ You can also use binary values.

Tokenization



$w_1 = t_1$

$w_m = t_m$

Tokenization. Split a string into a set of strings or words.

- ▶ Space based? (For Chinese language would not work)
- ▶ Vocabulary matching?

What's after you find the atomic unit of language you want to work on?

Text Pre-Processing

- ▶ Lower case everything. Problems: “they went to the White House” vs “they lived in a white house”. 
- ▶ Handling Abbreviations. “the dr. lives in a blue box”
 - ▶ I.O.U.
 - ▶ M.D.
 - ▶ N.B.A.
 - ▶ U.S.A. 
 - ▶ mr. 
 - ▶ .sh
 - ▶ .java 

Text Pre-Processing

- ▶ Remove punctuations.
 - ▶ Treat “blue,” differently from “blue”?
 - ▶ “?” , “!” these punctuations can be very useful.
 - ▶ Emojis :) : (?
 - ▶ Apostrophes
 - ▶ x
- ▶ Handling hyphenated words. Should I consider one space separated word or two?
 - ▶ “self-assessment”, “Los Angeles-based”
 - ▶ “the New York-based co-operative was fine-tuning forty-two K-9-like models”

Text Pre-Processing

- ▶ Numerical and special expressions
 - ▶ Email addresses
 - ▶ URLs
 - ▶ Complex enumeration of items
 - ▶ Telephone Numbers
 - ▶ 123 – 456 – 7890
 - ▶ (123) – 456 – 7890
 - ▶ Dates
 - ▶ July 10th, 2020
 - ▶ 10/07/2020
 - ▶ Time
 - ▶ Measures
 - ▶ Vehicle Licence Numbers
 - ▶ Paper and book citations

The same thing

Text Pre-Processing

- ▶ Entity Extractions
 - ▶ Person
 - ▶ Date
 - ▶ Location
 - ▶ Organization
- ▶ Few tokenizers tools
 - ▶ Apache Open NLP
 - ▶ Stanford Parser
- ▶ English Enclitics
 - ▶ He's → He has (highly present)
 - ▶ You're → You are

TEXT
I went to BARCELONA
Location

Text Pre-Processing

Stemming is a technique to convert a word to its “root” or “base” form

- ▶ This forces different forms of a word to be treated as the same feature



Word	Stem
fish	fish
fishes	fish
fished	fish
fishing	fish

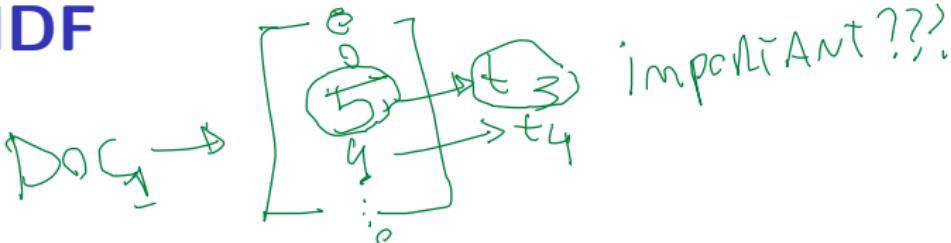
Text Pre-Processing

Stemming limitations.

- ▶ We are loosing info on the meaning of the sentence
- ▶ Stemming algorithms are making mistakes
 - ▶ university, universe, universal

Other things such as stop words removal.

TF-IDF



Common words like “the” and “and” have low information, and needs to be weighted.

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{idf}(t, d)$$
$$\text{idf}(t, d) = \log \frac{n_d}{1 + \text{df}(d, t)}$$

Some docs
diff. docs

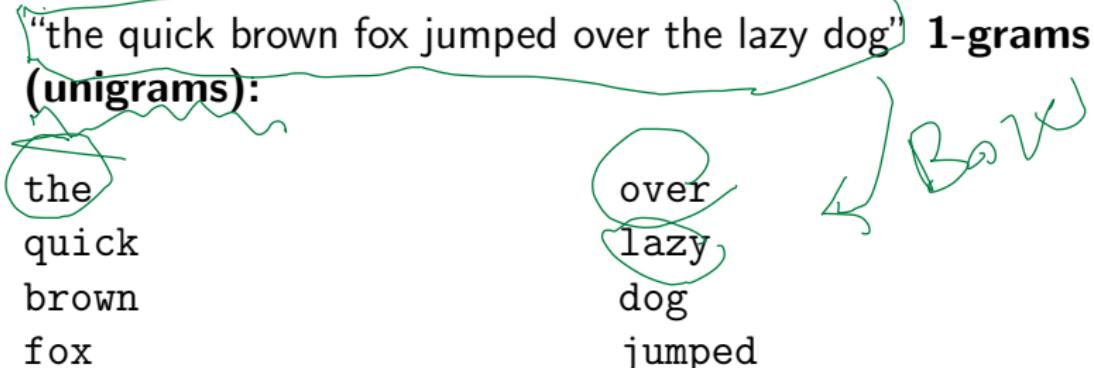
where n_d is the total number of documents, and $\text{df}(d, t)$ is the number of documents that contain the term t .

N-Grams

An extension of the the bag of words features:

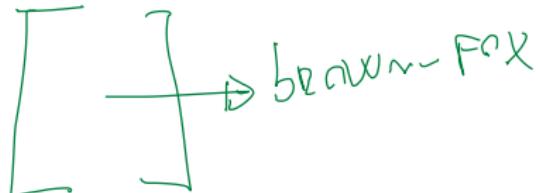
- ▶ n words in sequence considered as a single token
- ▶ 3-grams. All the 3-words in a document

The BOW is 1-gram.



2-Grams

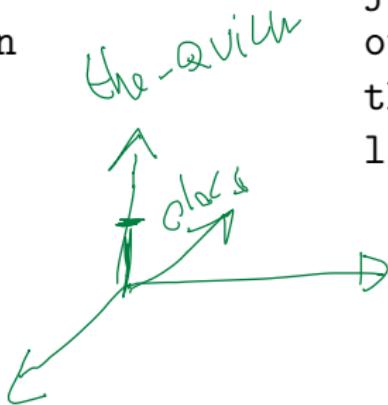
"the quick brown fox jumped over the lazy dog"



2-grams (bigrams):

the_quick
quick_brown
brown_fox
fox_jumped

jumped_over
over_the
the_lazy
lazy_dog



3-Grams

"the quick brown fox jumped over the lazy dog"

3-grams (trigrams):

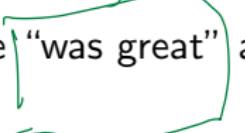
the_quick_brown
quick_brown_fox
brown_fox_jumped
fox_jumped_over

jumped_over_the
over_the_lazy
the_lazy_dog

N-Grams

There is a trade-off.

- ▶ Too short. We are missing important sequential detail
“great” vs “NOT great”
- ▶ Too long. Not enough examples.

- ▶ “the sausage was great” 
- ▶ “the cheese was great” 
- ▶ “the crust was great” 
- ▶ a bigram can capture “was great” a 4-gram no! 

Other Features in documents?

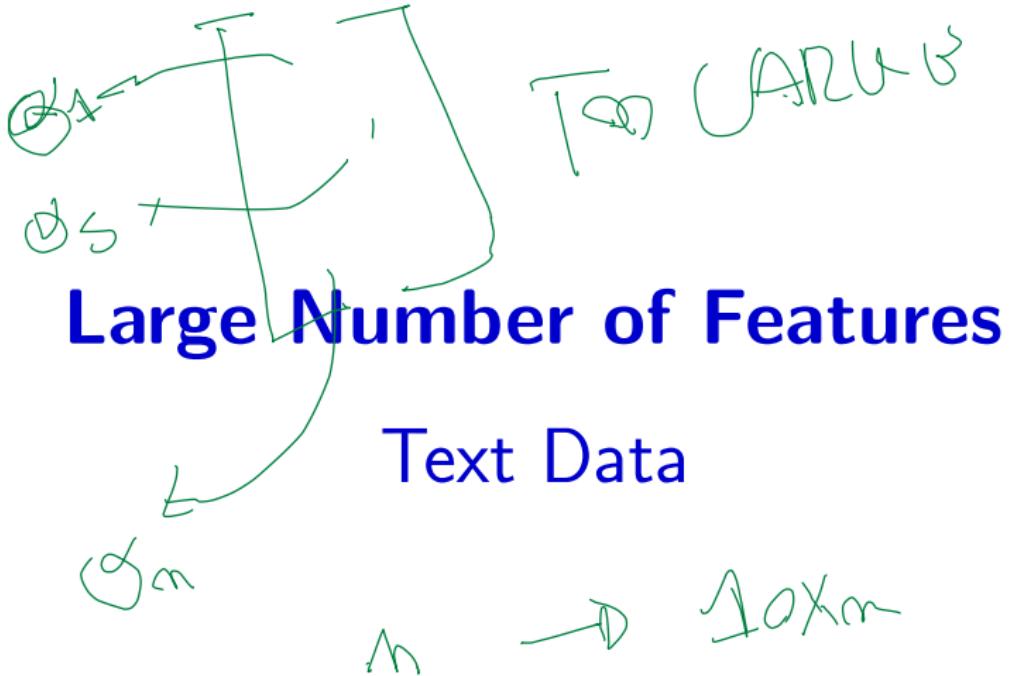
$\langle H_1 \rangle$ key ~ $\langle /H_1 \rangle$

Examples:

- ▶ Document length
- ▶ H_1 position
- ▶ Image ratio w.r.t. text
- ▶ HTML parsing



Notes



Feature Selection


$$\text{Hand} \rightarrow + \sum_{i=1}^m \| \phi_i \parallel$$

We can use for instance L1 regularization.

- ▶ Useful for feature selection
- ▶ The “selection” is happening during training (is not a pre-processing)

Often worse than L2 in prediction, but it can improve the prediction time.

Sequential algorithms (slow):

- ▶ Remove 1 feature at a time and see the performance.
- ▶ remove the feature that is less important.

ALL FEATURES \rightarrow PRECISION
REMOVED FEATURE \rightarrow //

Feature Selection

Other methods:

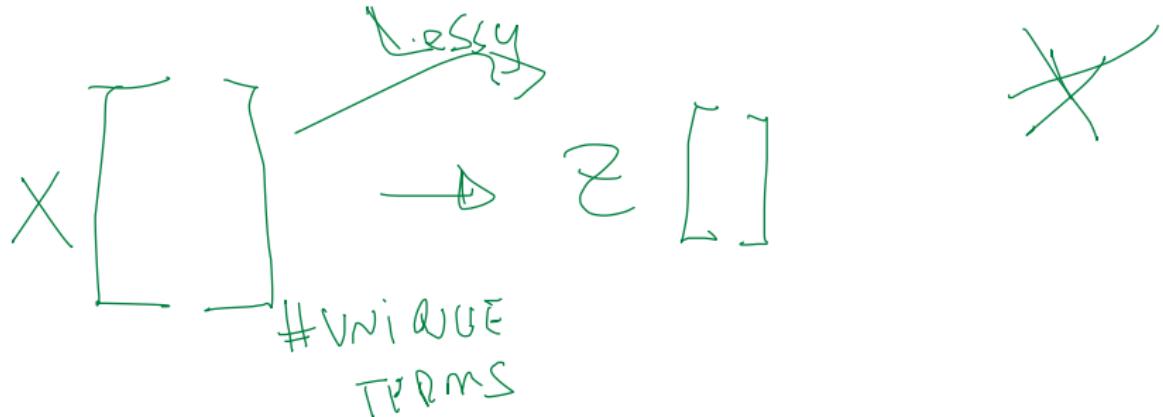
- ▶ Remove features based on statistical correlation.
- ▶ Remove the long tail.

Dimensionality Reduction

Dimensionality Reduction

Dimensionality reduction means reducing the number of features in your data

- ▶ We want to transform the feature space, creating features that are different from the original features.



Dimensionality Reduction

Example: two different types of blood pressure usually correlated

BP(S)	BP(D)	Heart Rate	Temperature
120	80	75	98.5
125	82	78	98.7
140	93	95	98.5
112	74	80	98.6

Dimensionality Reduction

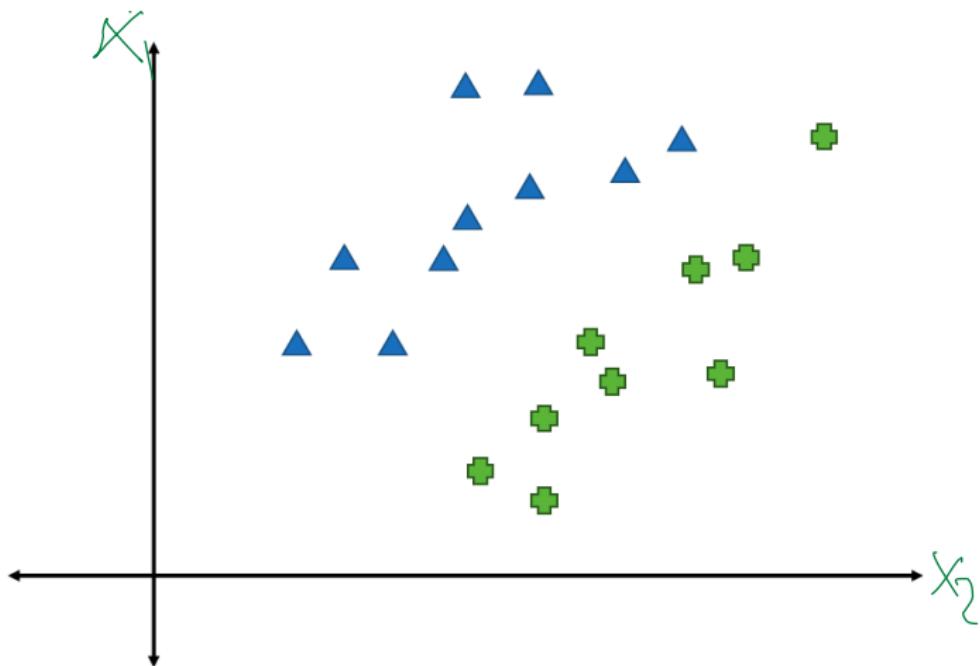
Example: two different types of blood pressure, usually correlated

BP(S)	BP(D)	BP-Avg	Heart Rate	Temperature
120	80	100	75	98.5
125	82	104	78	98.7
140	93	117	95	98.5
112	74	93	80	98.6

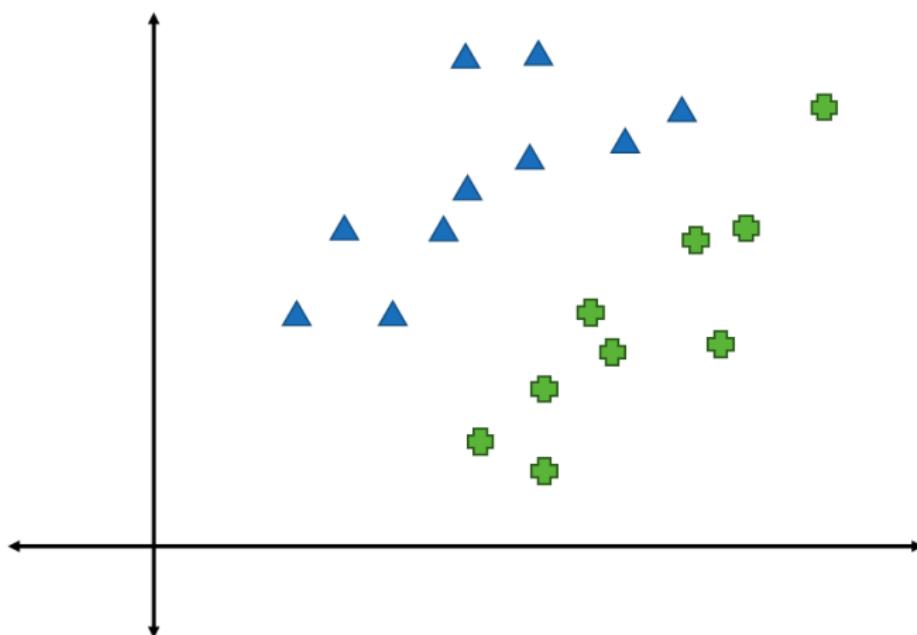
Replace the two BP features with the average of them.

- ▶ We are reducing the features in a lossy way, its not selection.

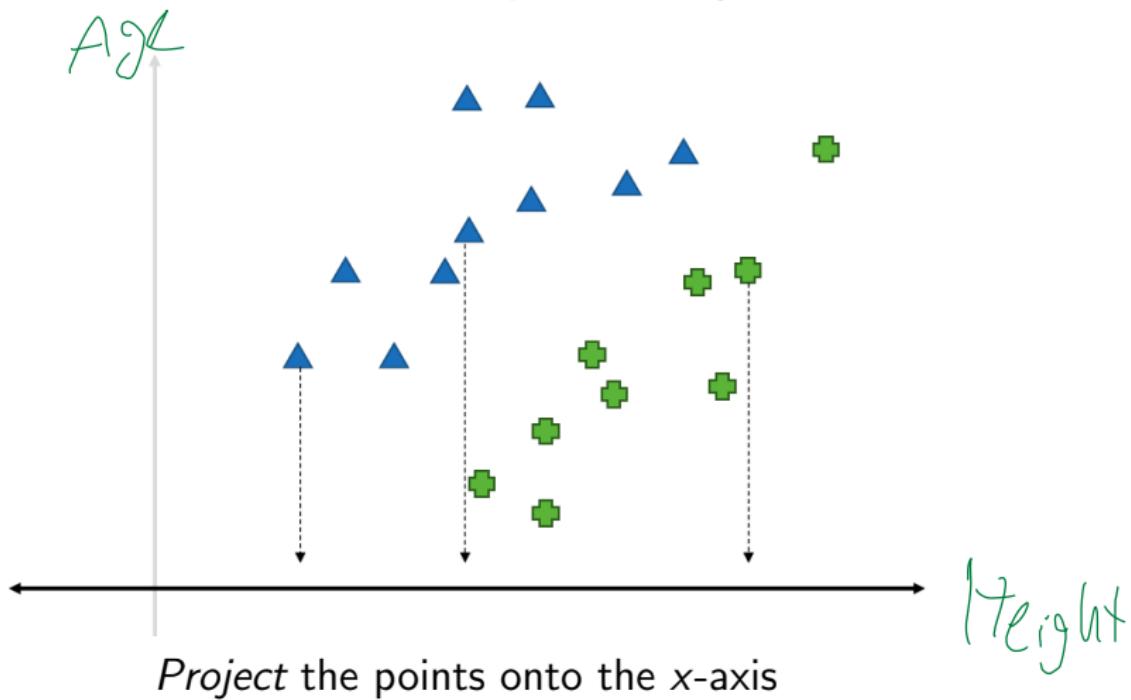
Suppose we have two dimensions (two features)



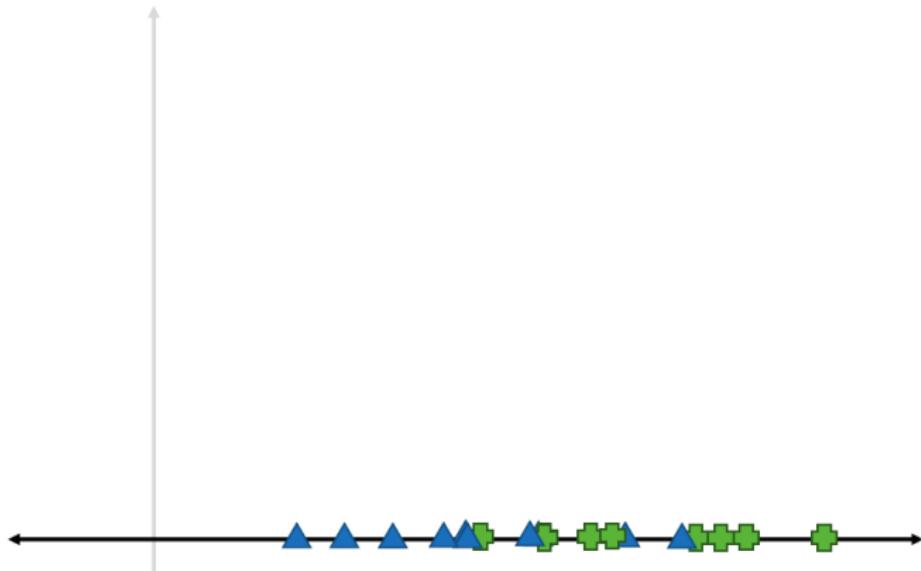
Feature selection: choose one of the two features to keep



Suppose we choose the feature represented by the x-axis

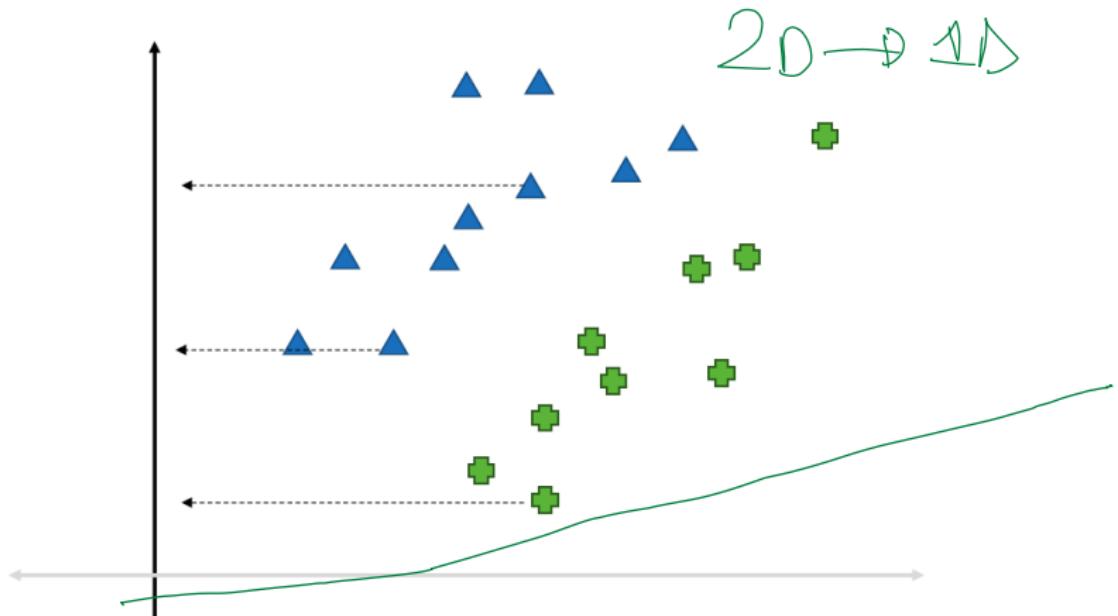


Suppose we choose the feature represented by the x-axis



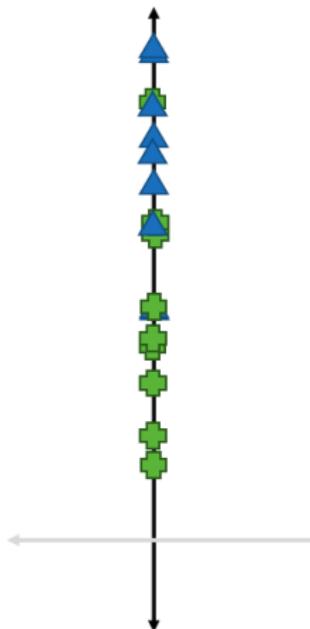
The positions along the x-axis now represent the feature values of each instance

Suppose we choose the feature represented by the y -axis



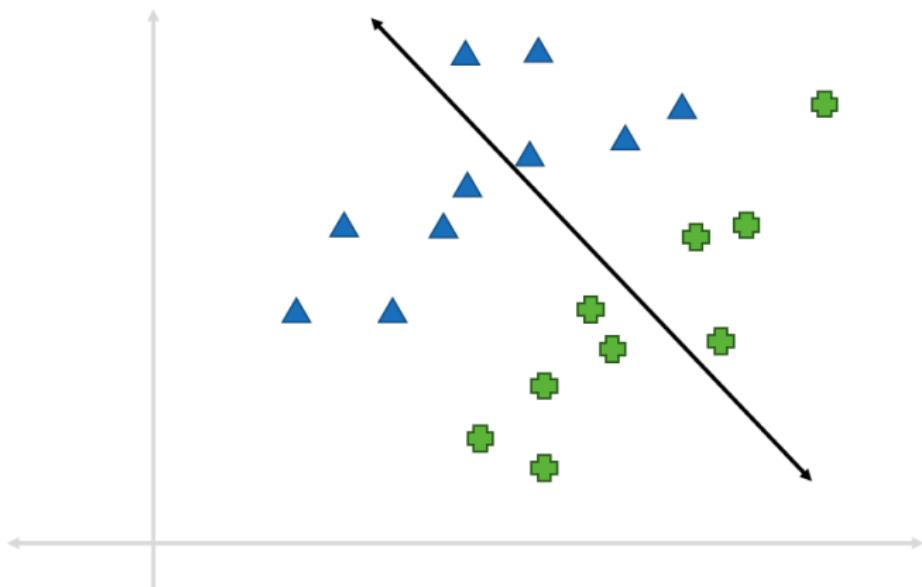
Project the points onto the y -axis

Suppose we choose the feature represented by the y -axis



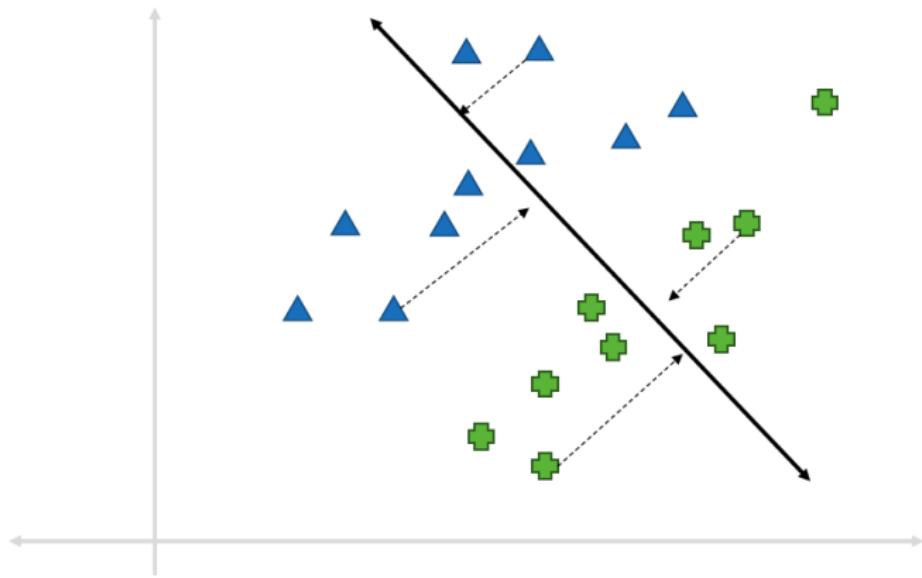
The positions along the y -axis now represent the feature values of each instance

We don't have to restrict ourselves to picking either the x -axis or y -axis



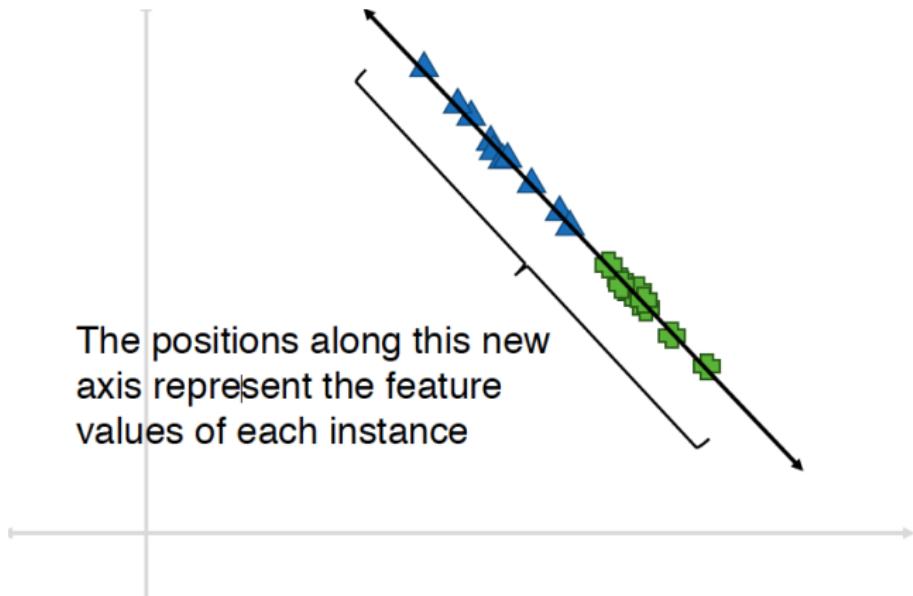
We could create a new axis!

We don't have to restrict ourselves to picking either the x -axis or y -axis

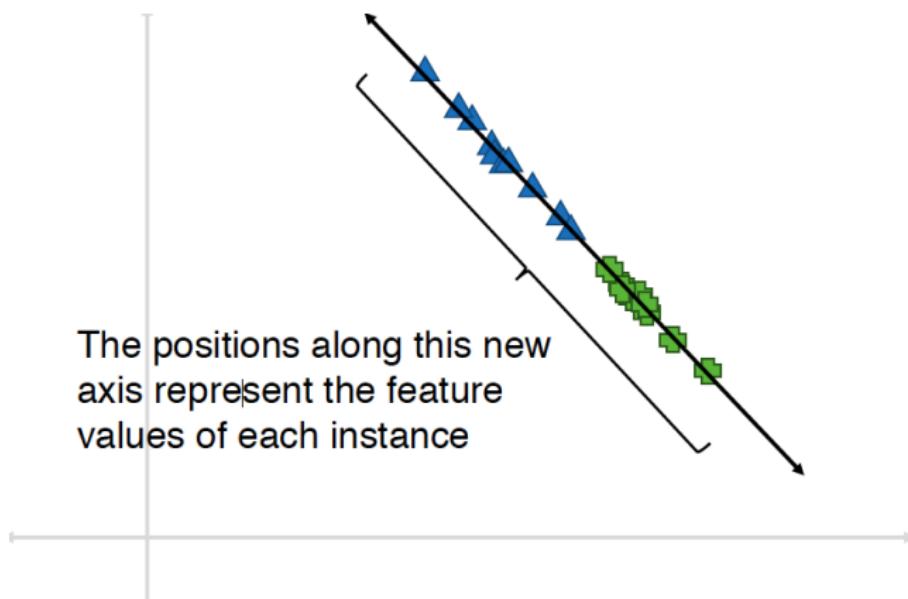


Project points onto this new axis

We don't have to restrict ourselves to picking either the x-axis or y-axis



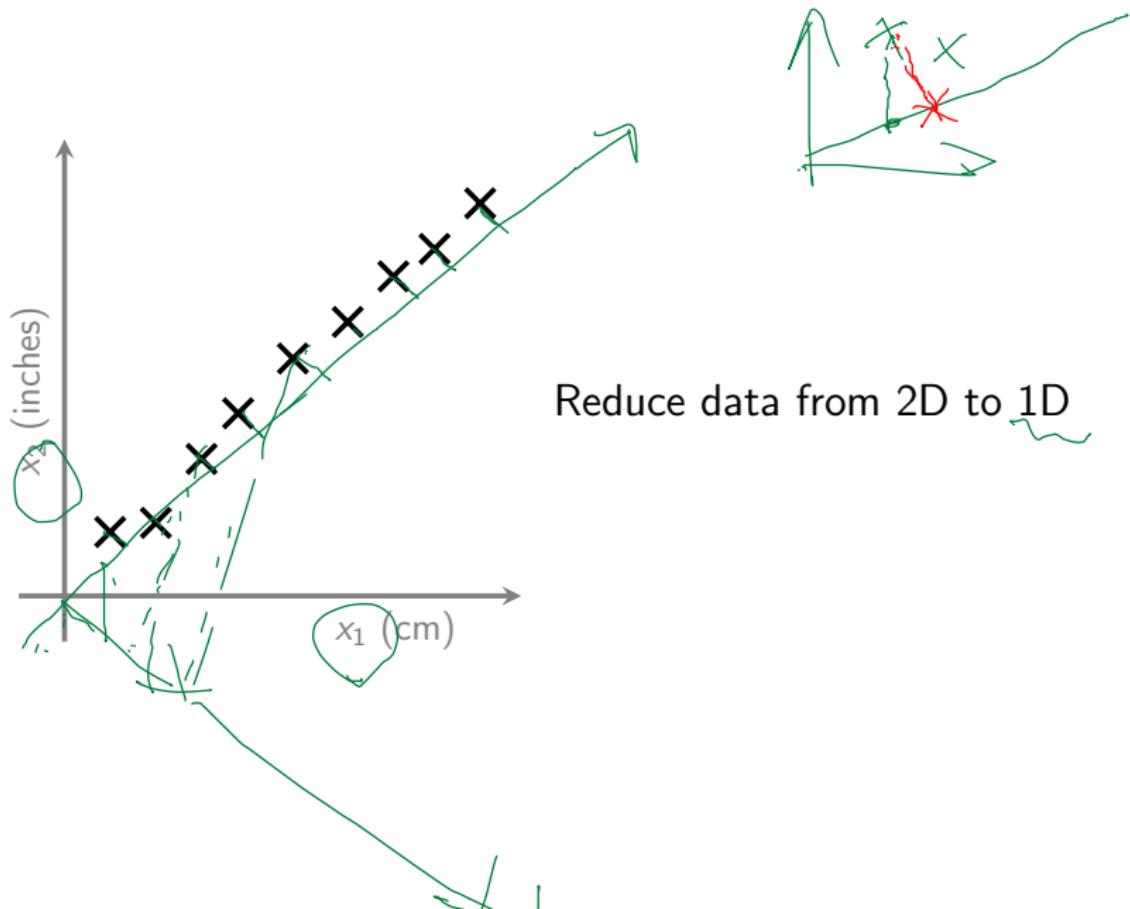
This is an example of transforming the feature space (as opposed to *selecting* a subset of features)



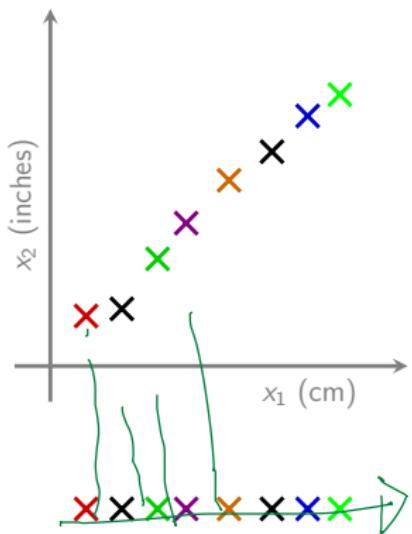
Dimensionality Reduction

Data Compression

Data Compression



Data Compression



[Age Height] \rightarrow [· J

Reduce data from 2D to 1D

$$\begin{aligned} x^{(1)} &\xrightarrow{\text{ER}} z^{(1)} \\ x^{(2)} &\rightarrow z^{(2)} \end{aligned}$$

:

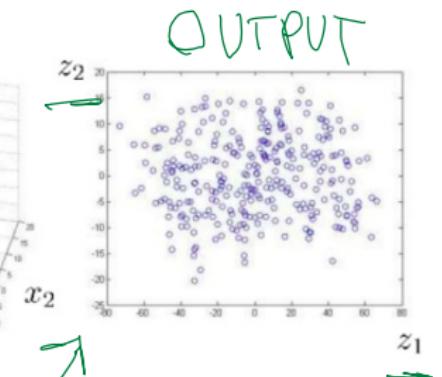
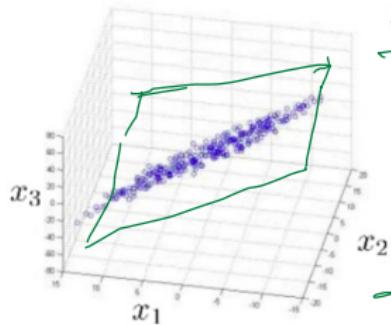
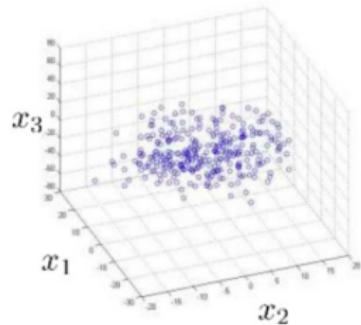
$$x^{(m)} \rightarrow z^{(m)}$$

$$\begin{aligned} x^{(1)} &\xrightarrow{\text{ER}} z^{(1)} \\ x^{(2)} &\rightarrow z^{(2)} \end{aligned}$$

$$x^{(m)} \rightarrow z^{(m)}$$

Data Compression

Reduce data from 3D to 2D



3D

→ 2D

Dimensionality Reduction

Data Visualization

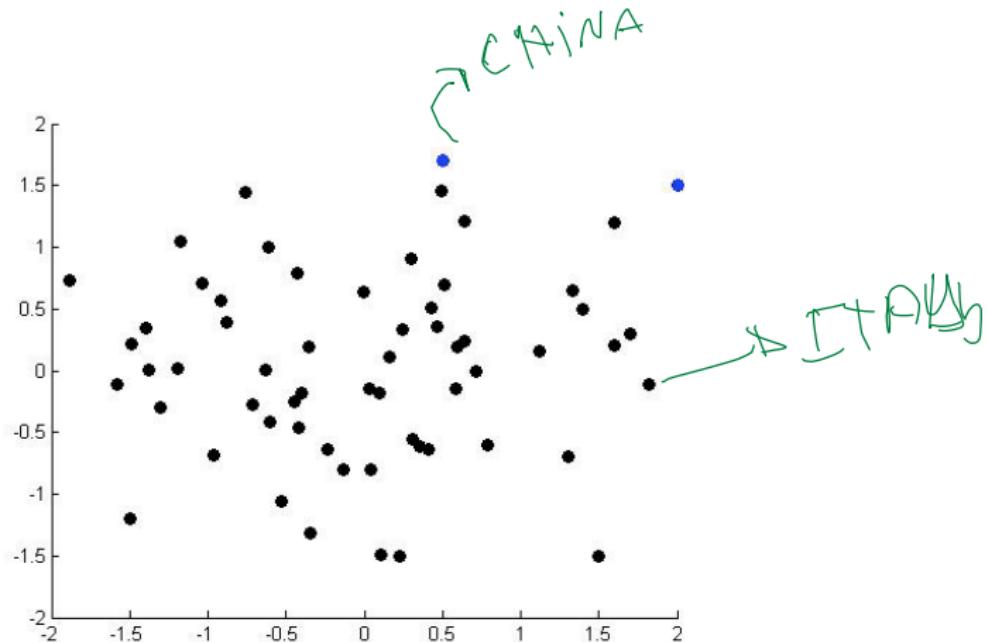
Data Visualization

Data Visualization



Country	z_1	z_2
Canada	1.6	1.2
China	1.7	0.3
India	1.6	0.2
Russia	1.4	0.5
Singapore	0.5	1.7
USA	2	1.5
...

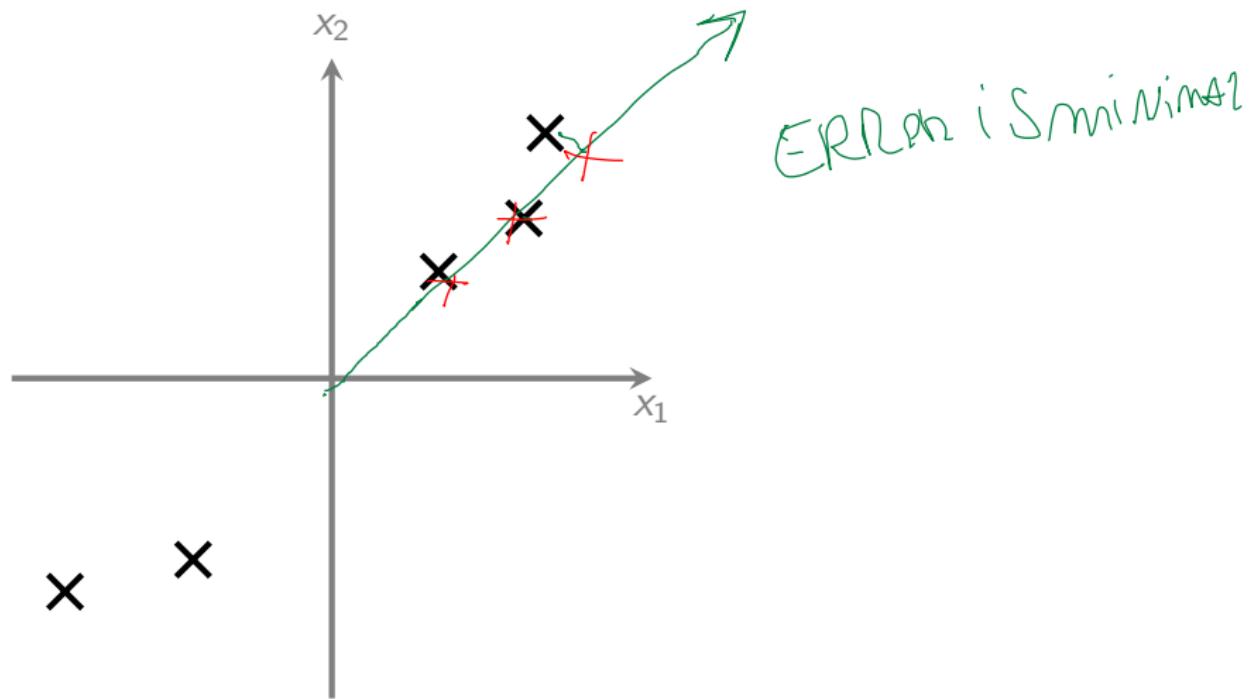
Data Visualization



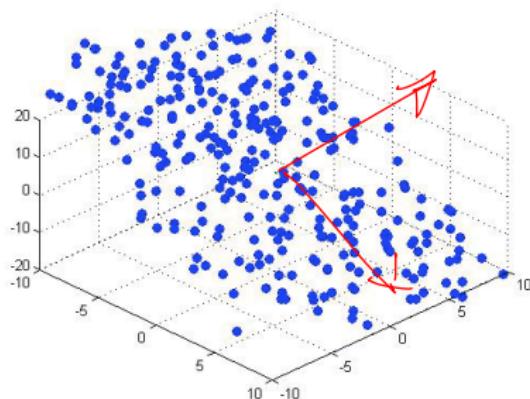
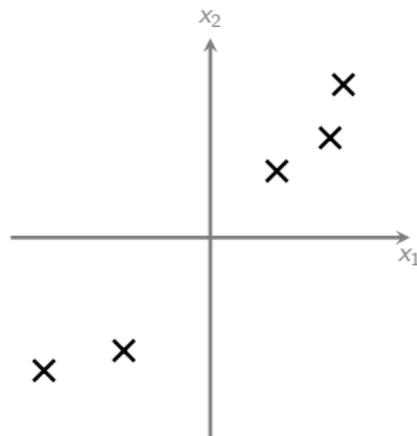
Dimensionality Reduction

Principal Component Analysis
problem formulation

Principal Component Analysis (PCA) problem formulation



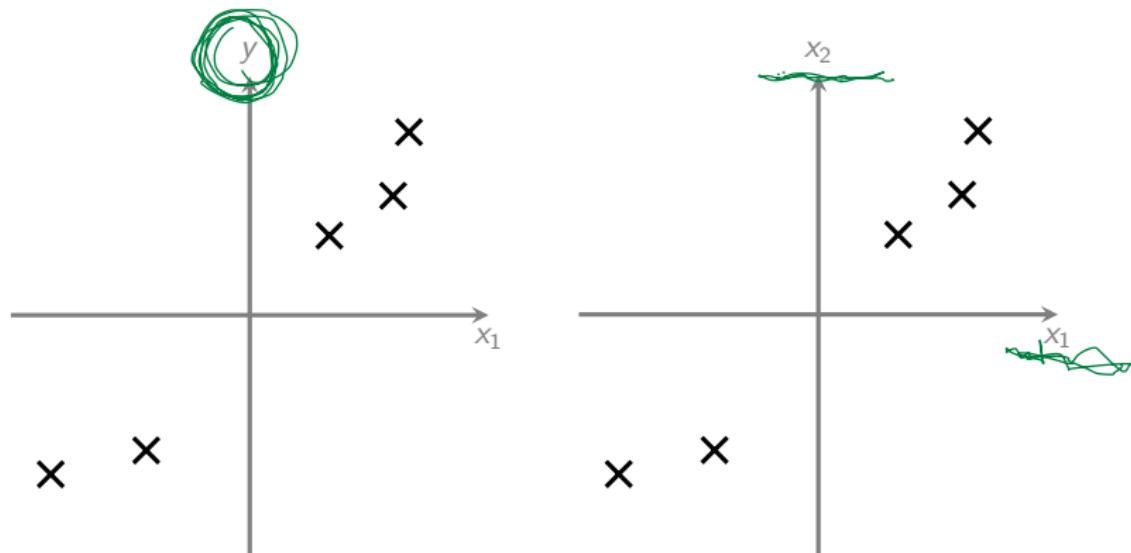
Principal Component Analysis (PCA) problem formulation



Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}$) onto which to project the data so as to minimize the projection error.

Reduce from n -dimension to k -dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

PCA is not linear regression



Dimensionality Reduction

Principal Component Analysis
algorithm

Data preprocessing

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each $x_j^{(i)}$ with $x_j - \mu_j$.

If different features on different scales (e.g., x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values.

Principal Component Analysis (PCA) algorithm

$$K < m$$

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$$

The diagram shows a column vector x_i with elements x_1, \dots, x_m . This vector is multiplied by its transpose x_i^T , which has elements x_1, \dots, x_n . The result is a scalar value representing the covariance between the i -th dimension and itself.

Compute "eigenvectors" of matrix Σ :

$$[U, S, V] = \text{svd}(\Sigma)$$

The diagram illustrates the singular value decomposition (SVD) of matrix Σ . It shows a large oval containing a smaller oval labeled Σ , with arrows pointing to the matrices U , S , and V in the equation above. Below the equation is a bracketed empty set symbol [].

$$\begin{bmatrix} x_1 x_1 & x_1 x_2 & \cdots & x_1 x_m \\ x_2 x_1 & x_2 x_2 & \cdots & x_2 x_m \\ \vdots & \vdots & \ddots & \vdots \\ x_m x_1 & x_m x_2 & \cdots & x_m x_m \end{bmatrix}$$

Principal Component Analysis (PCA) algorithm

V_{RESULT}

$\mathbb{R}^{n \times n}$

$x^{(i)} \in \mathbb{R}^n$

$z^{(i)} \in \mathbb{R}^k$

From $[U, S, V] = \text{svd}(\text{Sigma})$, we get

$$U = \begin{bmatrix} | & | & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(n)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$
$$z^{(i)} = \begin{bmatrix} \vdash & \vdash & \vdash \\ \vdash & \vdash & \vdash \\ \vdash & \vdash & \vdash \end{bmatrix} \begin{bmatrix} u^{(1)} \\ \vdash \\ u^{(n)} \end{bmatrix} + x^{(i)}$$

Principal Component Analysis (PCA) algorithm summary

$$T \cdot \text{ } \otimes = \begin{bmatrix} U & S & V \end{bmatrix}$$

After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

$[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\Sigma)$; 2nd step

$$\mathbf{U}_{\text{reduce}} = \mathbf{U}(:, 1:k)$$

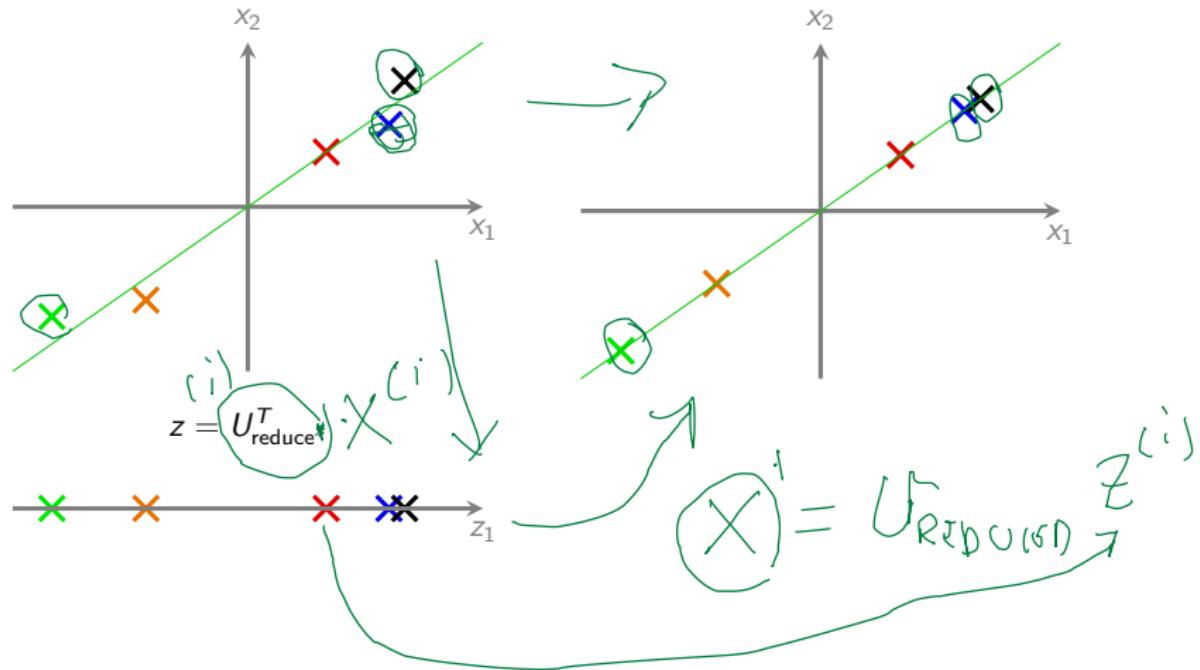
$$\mathbf{z} = \mathbf{U}_{\text{reduce}}' * \mathbf{x}$$

$$\mathbf{z}^{(i)} = [\star]$$

Dimensionality Reduction

Reconstruction from compressed
representation

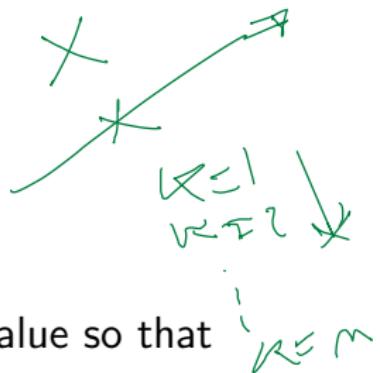
Reconstruction from compressed representation



Dimensionality Reduction

Choosing the number of principal components

Choosing k (number of principal components)



Average squared projection error:

Total variation in the data:

Typically, choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

"99% of variance is retained"

0.99
0.01

Choosing k (number of principal components)

Algorithm:

$$[U, S, V] = \text{svd}(Sigma)$$

Try PCA with $k = 1$

Compute U_{reduce} , $z^{(1)}$, $z^{(2)}$,
 $\dots, z^{(m)}$, $x_{\text{approx}}^1, \dots, x_{\text{approx}}^m$

Check if

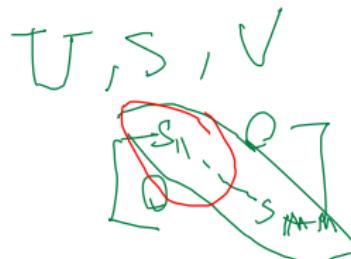
$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

\Rightarrow

$$1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^m s_{ii}}$$

Choosing k (number of principal components)

$$X = Q \begin{pmatrix} 0.0 \\ & \ddots \end{pmatrix}$$



$$[U, S, V] = \text{svd}(\Sigma)$$

Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

(99% of variance retained)

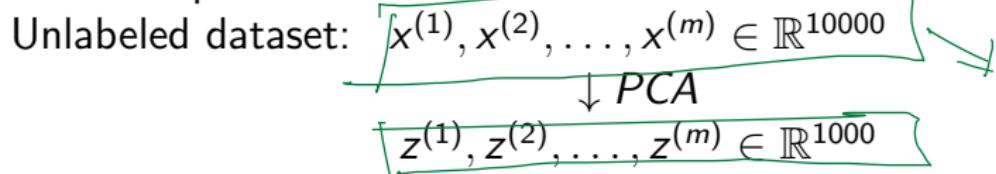
Dimensionality Reduction

Advice for applying PCA

Supervised learning speedup

$$\underbrace{(x^{(1)}, y^{(1)})}_{\text{Input}}, \underbrace{(x^{(2)}, y^{(2)})}_{\text{Input}}, \dots, \underbrace{(x^{(m)}, y^{(m)})}_{\text{Input}}$$

Extract inputs:



New training set:

$$\underbrace{(z^{(1)}, y^{(1)})}_{\text{Input}}, \underbrace{(z^{(2)}, y^{(2)})}_{\text{Input}}, \dots, \underbrace{(z^{(m)}, y^{(m)})}_{\text{Input}}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

Application of PCA

- ▶ Compression
 - ▶ Reduce memory/disk needed to store data
 - ▶ Speed up learning algorithm
- ▶ Visualization

Bad use of PCA: To prevent overfitting



Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features to $k < n$.

Thus, fewer features, less likely to overfit.

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

PCA is sometimes used where it shouldn't be

Design of ML system:

- ▶ Get training set
 $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- ▶ Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- ▶ Train logistics regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- ▶ Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_\theta(z)$ on
 $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

How about doing the whole thing without using PCA?

Before implementing PCA first try running whatever you want to do with the original/raw data $x^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $z^{(i)}$.