# Week 7: Support Vector Machines

## Large Margin Classification

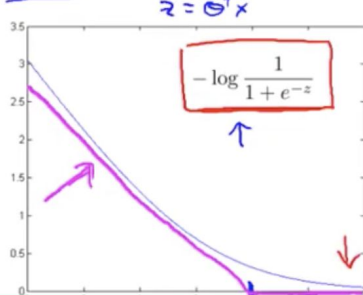### Optimization Objective



**Alternative view of logistic regression**  $(x, y)$
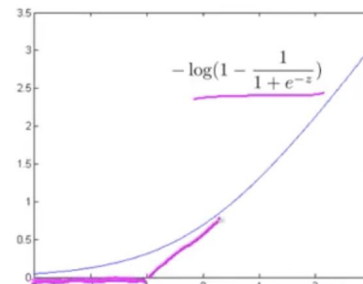
Cost of example:  $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$ ←

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log(1 - \frac{1}{1 + e^{-\theta^T x}})$$ ←

If $y = 1$ (want $\theta^T x \gg 0$):   $z = \theta^T x$   If $y = 0$ (want $\theta^T x \ll 0$):

$-\log \frac{1}{1 + e^{-z}}$   $-\log(1 - \frac{1}{1 + e^{-z}})$

7:25  /  14:47

$$\min_\theta C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

## Large Margin Intuition

**Support Vector Machine**

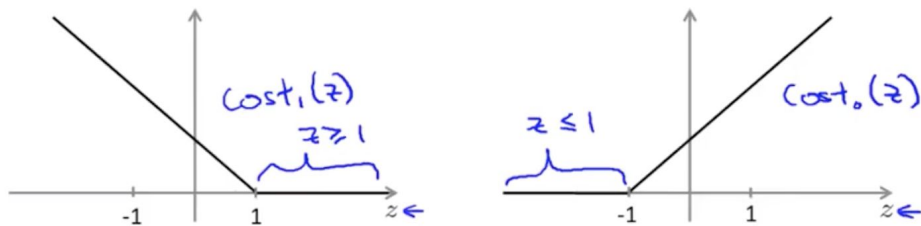$$\rightarrow \quad \min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

$cost_1(z)$    $z \geq 1$      $z \leq 1$     $cost_0(z)$

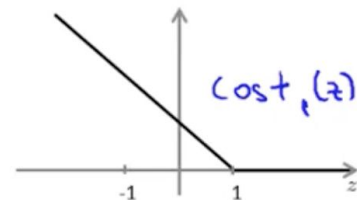$\rightarrow$ If $y = 1$, we want $\theta^T x \geq 1$ (not just $\geq 0$)

$\rightarrow$ If $y = 0$, we want $\theta^T x \leq -1$ (not just $< 0$)

**SVM Decision Boundary**

$$\min_{\theta} C \sum_{i=1}^{m} \left[ y^{(i)} cost_1(\theta^T x^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

$= 0$

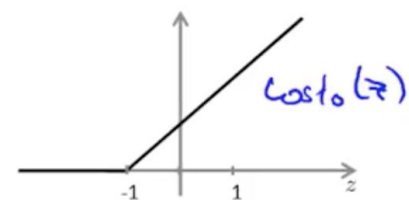Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

$cost_1(z)$

Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

$cost_0(z)$

$$\min_{\theta} \quad \cancel{C \times 0} + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$
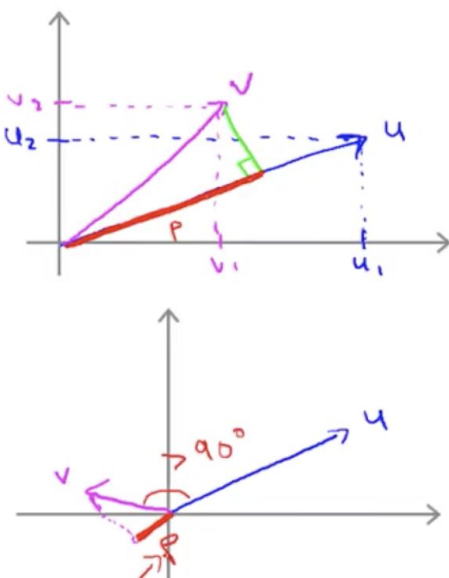
$$s.t. \quad \theta^T x^{(i)} \geq 1 \quad if \quad y^{(i)} = 1$$

$$\theta^T x^{(i)} \leq -1 \quad if \quad y^{(i)} = 0.$$

$\rightarrow$ C very large

$$\frac{1}{\lambda}$$

$\leftarrow$ C not too large

# Mathematics Behind Large Margin Classification

**Vector Inner Product**



$$\rightarrow u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \qquad \rightarrow v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$u^T v = ?$ $\qquad [u_1 \quad u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

$\|u\| = $ length of vector $u$

$\qquad = \sqrt{u_1^2 + u_2^2} \quad \in \mathbb{R}$

$p = $ length of projection of $v$ onto $u$.

$u^T v = p \cdot \|u\| \leftarrow \qquad = v^T u$

Signed

$\qquad = u_1 v_1 + u_2 v_2 \leftarrow \qquad p \in \mathbb{R}$

$u^T v = p \cdot \|u\|$

$p < 0$

## SVM Decision Boundary

$$\omega = (\sqrt{\omega})^2$$

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{n} \theta_j^2 \;=\; \tfrac{1}{2}(\theta_1^2 + \theta_2^2) \;=\; \tfrac{1}{2}\left(\sqrt{\theta_1^2 + \theta_2^2}\right)^2 \;=\; \tfrac{1}{2}\|\theta\|^2$$

$$= \|\theta\|$$

s.t. $\theta^T x^{(i)} \geq 1$   if $y^{(i)} = 1$
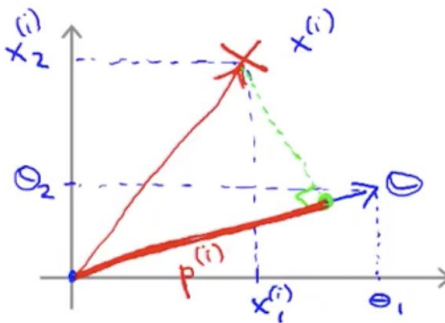
$\theta^T x^{(i)} \leq -1$   if $y^{(i)} = 0$

$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$   $\theta_0 = 0$

Simplication: $\theta_0 = 0$.    $n=2$

$\theta^T x^{(i)} = ?$

$\uparrow \quad \uparrow$
$u^T v$

$\theta^T x^{(i)} = \boxed{p^{(i)} \cdot \|\theta\|}$

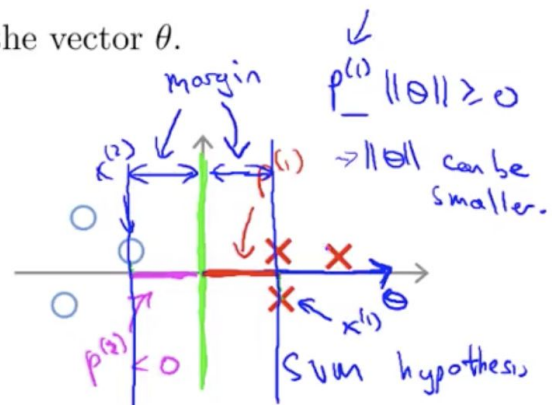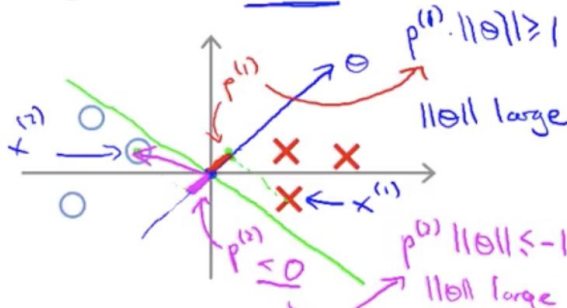$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$



Andrew Ng

## SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{n} \theta_j^2 \;=\; \tfrac{1}{2}\|\theta\|^2$$

s.t. $p^{(i)} \cdot \|\theta\| \geq 1$   if $y^{(i)} = 1$

$p^{(i)} \cdot \|\theta\| \leq -1$   if $y^{(i)} = 1$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector $\theta$.

Simplification: $\theta_0 = 0$

margin    $p^{(i)} \|\theta\| \geq 0$

$p^{(i)} \cdot \|\theta\| \geq 1$

$\|\theta\|$ large

$\|\theta\|$ can be smaller.

$p^{(2)} \|\theta\| \leq -1$

$\|\theta\|$ large

$p^{(2)} < 0$

SVM hypothesis



$\theta_0 = 0$: decision boundary passes the origin.

$$||\theta|| = \frac{1}{2}$$

# Kernels

**Kernel**



Given $x$, compute new feature depending on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

$x_2$

$x_1$

Given $x$:

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(- \frac{||x - l^{(1)}||^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, l^{(1)}) = \exp\left(- \frac{||x - l^{(2)}||^2}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp(\ldots)$$

kernel (Gaussian kernels)

$||w||$

$-||x - l^{(1)}||^2$

$||x - l^{(1)}||^2$

$k(x, l^{(i)})$

## Example:

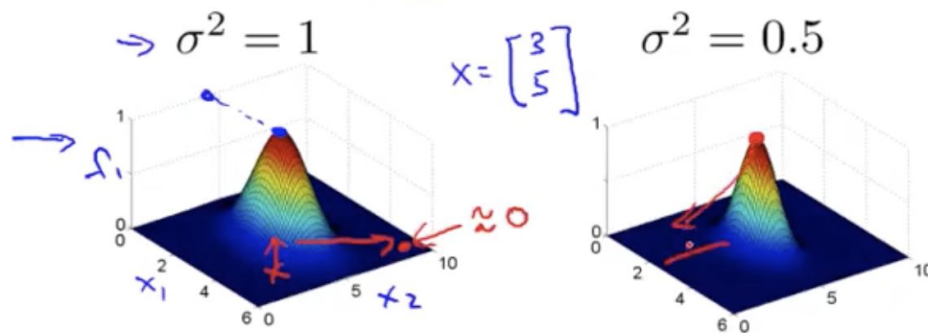$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$\sigma^2 = 1 \qquad \sigma^2 = 0.5$$

$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$



When $\sigma$ is small, the height does not change while it narrows down.

## SVM with Kernels

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})$,

choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \ldots, l^{(m)} = x^{(m)}$.

Given example $x$:

$$f_1 = \text{similarity}(x, l^{(1)})$$
$$f_2 = \text{similarity}(x, l^{(2)})$$
$$\ldots$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \qquad f_0 = 1$$

For training example $(x^{(i)}, y^{(i)})$:

$$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} \quad \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$$

$$f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp\left(-\frac{0}{2\sigma^2}\right) = 1$$

$$x^{(i)} \in \mathbb{R}^{n+1} \quad (\text{or } \mathbb{R}^n)$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

$$f_0^{(i)} = 1$$

Andrew Ng

## SVM with Kernels

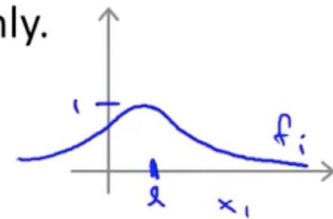Hypothesis: Given $x$, compute features $f \in \mathbb{R}^{m+1}$ $\quad \theta \in \mathbb{R}^{m+1}$
  → Predict "y=1" if $\theta^T f \geq 0$ $\qquad \theta_0 f_0 + \theta_1 f_1 + \cdots + \theta_m f_m$

Training: $\qquad n = m$

$$\min_\theta C \sum_{i=1}^{m} y^{(i)} cost_1(\theta^T f^{(i)}) + (1 - y^{(i)})cost_0(\theta^T f^{(i)}) + \frac{1}{2}\sum_{j=1}^{m}\theta_j^2$$

$\theta^T x^{(i)} \quad \theta^T f^{(i)}$

$= m$

$\to \theta_0$

Using kernels for logistic regression is very slow, because it runs very slowly.

## SVM parameters:

$C\ (=\frac{1}{\lambda})$. → Large C: Lower bias, high variance. $\quad$ (small $\lambda$)
  $\quad$ → Small C: Higher bias, low variance. $\quad$ (large $\lambda$)

$\sigma^2 \quad$ Large $\sigma^2$: Features $f_i$ vary more smoothly.
  → Higher bias, lower variance.

$\exp\left(-\frac{\|x - \ell^{(i)}\|^2}{2\sigma^2}\right)$

Small $\sigma^2$: Features $f_i$ vary less smoothly.
  Lower bias, higher variance.

Suppose you train an SVM and find it overfits your training data. Which of these would be a reasonable next step?

Decrease C and/or Increase sigma square

# SVM in Practice

Package: liblinear, libsvm,

Specify:

parameter C;

kernel (similarity function);

E.g. No kernel ("linear kernel")      $\Theta_0 + \Theta_1 x_1 +$
    Predict "y = 1" if $\underline{\theta^T x \geq 0}$      $\rightarrow \underline{n}$ large,

Gaussian kernel:
$$f_i = \exp\left(-\frac{||x - l^{(i)}||^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}.$$
Need to choose $\sigma^2$.

Other kernels: polynomial kernel, string kernel, chi-square kernel, histogram intersection kernel…

Note: not all similarity functions made valid kernel, need to satisfy "mercer's theorem" to make sure SVM packages' optimizations run correction, and do not diverge.

## Logistic regression vs. SVMs

$n =$ number of features ($x \in \mathbb{R}^{n+1}$), $m =$ number of training examples
> If $n$ is large (relative to $m$):  (e.g. $n \geq m$, $n = 10,000$, $m = 10 \cdots 1000$)
> Use logistic regression, or SVM without a kernel ("linear kernel")

> If $n$ is small, $m$ is intermediate:    ($n = 1 - 1000$, $m = 10 - 10,000$)
    → Use SVM with Gaussian kernel

If $n$ is small, $m$ is large:   ($n = 1 - 1000$, $m = 50,000+$)
    → Create/add more features, then use logistic regression or SVM without a kernel

Neural network likely to work well for most of these settings, but may be slower to train; SVM does not need to worry about local minimum.

Quiz:

Suppose you have 2D input examples (ie, $x^{(i)} \in \mathbb{R}^2$). The decision boundary of the SVM (with the linear kernel) is a straight line.

—— ✕ ——

# Week 8:

## Clustering

Application: market segmentation; Social network analysis; organize computing clusters; astronomical data analysis

## K-Means Algorithm

## K-means algorithm

Input:
- $K$ (number of clusters) ←
- Training set $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

## K-means algorithm

$\mu_1 \quad \mu_2$
$\times \quad \times$

Randomly initialize $K$ cluster centroids $\mu_1, \mu_2, \ldots, \mu_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step
> for $i = 1$ to $m$
> $\quad c^{(i)} :=$ index (from 1 to $K$) of cluster centroid closest to $x^{(i)}$
>
> $\min_k \|x^{(i)} - \mu_k\|^2 \hookrightarrow c^{(i)}$

Move centroid
> for $k = 1$ to $K$
> $\rightarrow \mu_k :=$ average (mean) of points assigned to cluster $k$
> $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \rightarrow c^{(1)}=2, c^{(5)}=2, c^{(6)}=2, c^{(10)}=2$

$\mu_2 = \frac{1}{4}\left[x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}\right] \in \mathbb{R}^n$

}

K: total number of clusters; k: the index of cluster

# Optimization Objective

Usage:

debug the learning algorithm to ensure the K-Means is running correctly;

Find better costs for this and avoid the local ultima

# Random Initialization

- Should have K < m
- Randomly pick K training examples
- Set $\mu_1, \dots \mu_K$ equal to these K examples. $\mu_1, \dots \mu_K$ equal to these K examples.

**Random initialization**

For i = 1 to 100 {                          So – 1000

→ Randomly initialize K-means.
Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$.
Compute cost function (distortion)
→ $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$
}
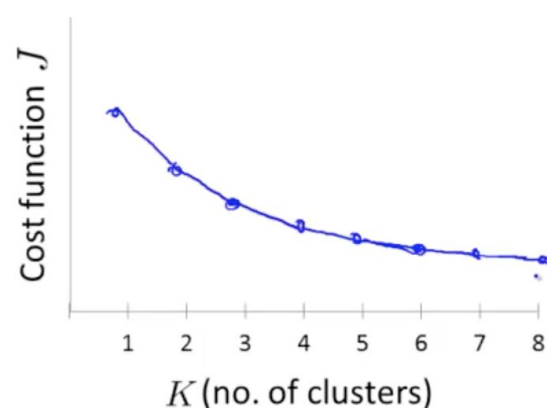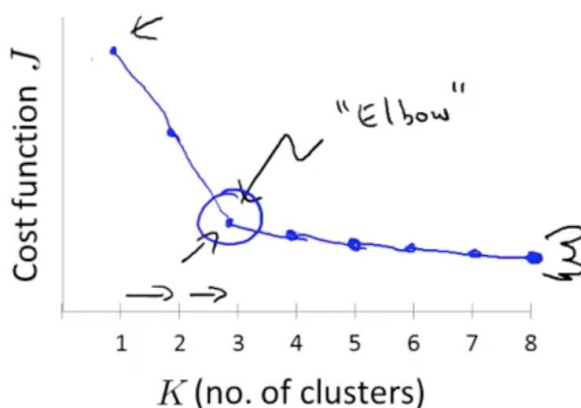
Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

## Choosing the number of Clusters

**Choosing the value of K**

Elbow method:



# Dimensionality Reduction Motivation

## Data Compression

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

## Visualization

# Principal Component Analysis

## PCA Problem Formulation

Minimize the projection error

The left is linear regression; and the right is PCA

**PCA is not linear regression**



**Data preprocessing**

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

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.

# PCA Algorithm

**Principal Component Analysis (PCA) algorithm**

Reduce data from $n$-dimensions to $k$-dimensions

Compute "covariance matrix":

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

$n \times 1$    $1 \times n$    $n \times n$    Sigma

Compute "eigenvectors" of matrix $\Sigma$:

$$[U,S,V] = svd(Sigma);$$

$n \times n$   matrix.

$\rightarrow$ Singular value decomposition

eig (Sigma)

Since $\Sigma$ is symmetrical, we can get the same answer from both svd and eig.

## Principal Component Analysis (PCA) algorithm

From `[U,S,V] = svd(Sigma)`, we get:

$$\to U = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(n)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$k$

$x \in \mathbb{R}^n \longrightarrow z \in \mathbb{R}^k$

$$z = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(k)} \\ | & | & & | \end{bmatrix}^T \times = \begin{bmatrix} \text{---} (u^{(1)})^T \text{---} \\ \vdots \\ \text{---} (u^{(k)})^T \text{---} \end{bmatrix} \underset{n \times 1}{\overset{\downarrow x}{}}$$

$z \in \mathbb{R}^k$

$\underbrace{\qquad}_{n \times k}$ $\underbrace{\qquad}_{k \times n}$

$U_{reduce}$ $\underbrace{\qquad}_{k \times 1}$

## Principal Component Analysis (PCA) algorithm summary

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

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

$$X = \begin{bmatrix} \text{---} x^{(1)T} \text{---} \\ \vdots \\ \text{---} x^{(m)T} \text{---} \end{bmatrix}$$

$\to$ `[U,S,V] = svd(Sigma);`

$\to$ `Ureduce = U(:,1:k);`

$\to$ `z = Ureduce'*x;`

$\to$ $\text{Sigma} = (1/m) * X' * X;$

$x \in \mathbb{R}^n$ $x_0 \neq 1$

# Applying PCA

Reconstruction from Compressed Representation

$$z \in \mathbb{R} \longrightarrow x \in \mathbb{R}^2$$

$$\begin{bmatrix} x_{(i)} \\ x_{approx}^{(i)} = U_{reduce} \cdot z^{(i)} \\ \downarrow \mathbb{R}^n \end{bmatrix}$$

$\underbrace{n \times k}$ $\underbrace{k \times 1}$

$\underbrace{\qquad}_{n \times 1}$

## Choosing the Number of Principal Components

**Choosing $k$ (number of principal components)**

Average squared projection error: $\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} - x_{approx}^{(i)} \|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^{m} \| x^{(i)} \|^2$

Typically, choose $k$ to be smallest value so that

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

"99% of variance is retained"

**Choosing $k$ (number of principal components)**

Algorithm:

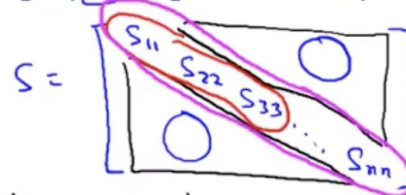Try PCA with $k = 1$ $\quad k=2 \quad k=3$ $\quad k=4$

Compute $U_{reduce}, z^{(1)}, z^{(2)},$

$\ldots, z^{(m)}, x_{approx}^{(1)}, \ldots, x_{approx}^{(m)}$

Check if

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

$k = 17$

$\rightarrow$ `[U,S,V] = svd(Sigma)`

$S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & S_{33} & \\ & & & \ddots & \\ & & & & S_{nn} \end{bmatrix}$

For given $k$

$k = 3$

$1 - \dfrac{\sum_{i=1}^{k} S_{ii}}{\sum_{i=1}^{n} S_{ii}} \leq 0.01$

## Advice for Applying PCA

**Bad use of PCA: To prevent overfitting**

$\rightarrow$ Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features to $k < n$. — 10000 1000

Thus, fewer features, less likely to overfit.

Bad !

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

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

**PCA is sometimes used where it shouldn't be**

Design of ML system:
- $\rightarrow$ Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$
- $\rightarrow$ Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- $\rightarrow$ Train logistic regression on $\{(z^{(1)}, y^{(1)}), \ldots, (z^{(m)}, y^{(m)})\}$
- $\rightarrow$ Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_\theta(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \ldots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

$\rightarrow$ How about doing the whole thing without using PCA?

$\rightarrow$ 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 $z^{(i)}$.

Quiz:

If the input features are on very different scales, it is a good idea to perform feature scaling before PCA