CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Oct 31, 2025

Administration

Reminder: HW3 is due on Nov 5th.

Outline

1 Principal Component Analysis (PCA)

Markov models

Midden Markov Model

Outline

- Principal Component Analysis (PCA)
- 2 Markov models
- 3 Hidden Markov Model

Dimensionality reduction is yet another important unsupervised learning problem.

Dimensionality reduction is yet another important unsupervised learning problem.

Goal: reduce the dimensionality of a dataset so

• it is easier to visualize and discover patterns

Dimensionality reduction is yet another important unsupervised learning problem.

Goal: reduce the dimensionality of a dataset so

- it is easier to visualize and discover patterns
- it takes less time and space to process for any applications (classification, regression, clustering, etc)

Dimensionality reduction is yet another important unsupervised learning problem.

Goal: reduce the dimensionality of a dataset so

- it is easier to visualize and discover patterns
- it takes less time and space to process for any applications (classification, regression, clustering, etc)
- noise is reduced
-

Dimensionality reduction is yet another important unsupervised learning problem.

Goal: reduce the dimensionality of a dataset so

- it is easier to visualize and discover patterns
- it takes less time and space to process for any applications (classification, regression, clustering, etc)
- noise is reduced
- <u>. . . .</u>

There are many approaches, we focus on a linear method: **Principal Component Analysis (PCA)**

Consider the following dataset:

• 17 features, each represents the average consumption of some food

Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	1102	674	957	1137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1506	1572	12 56
Sugars	156	139	147	175

Example

picture from here

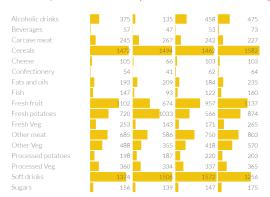
Consider the following dataset:

- 17 features, each represents the average consumption of some food
- 4 data points, each represents some country

Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	1102	674	957	1 137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1506	1572	1256
Sugars	156	139	147	175

Consider the following dataset:

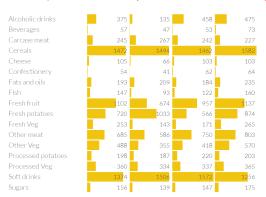
- 17 features, each represents the average consumption of some food
- 4 data points, each represents some country



What can you tell?

Consider the following dataset:

- 17 features, each represents the average consumption of some food
- 4 data points, each represents some country



What can you tell?

Hard to say anything looking at all these 17 features.

Example

picture from here

PCA can help us!

PCA can help us! The first principal component of this dataset:



i.e. we reduce the dimensionality from 17 to just 1.

PCA can help us! The first principal component of this dataset:



i.e. we reduce the dimensionality from 17 to just 1.

Now one data point is clearly different from the rest!

PCA can help us! The first principal component of this dataset:

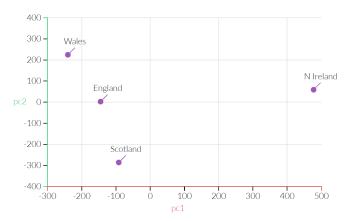


i.e. we reduce the dimensionality from 17 to just 1.

Now one data point is clearly different from the rest!

That turns out to be data from Northern Ireland, the only country not on the island of Great Britain out of the 4 samples.

PCA can find the **second (and more) principal component** of the data too:



High level idea

How does PCA find these principal components (PC)?



High level idea

How does PCA find these principal components (PC)?



The first PC is in fact **the direction with the most variance**, i.e. the direction where the data is most spread out.

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance,

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance, i.e.

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left(\boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{v} - \frac{1}{N} \sum_m \boldsymbol{x}_m^{\mathrm{T}} \boldsymbol{v}\right)^2$$

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance, i.e.

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left(\boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{v} - \frac{1}{N} \sum_m \boldsymbol{x}_m^{\mathrm{T}} \boldsymbol{v}\right)^2$$

 $oldsymbol{\circ} oldsymbol{x}_n^{
m T} oldsymbol{v}$ is exactly the projection of $oldsymbol{x}_n$ onto the direction $oldsymbol{v}$

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance, i.e.

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left(\boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{v} - \frac{1}{N} \sum_m \boldsymbol{x}_m^{\mathrm{T}} \boldsymbol{v}\right)^2$$

- $oldsymbol{\circ} oldsymbol{x}_n^{
 m T} oldsymbol{v}$ is exactly the projection of $oldsymbol{x}_n$ onto the direction $oldsymbol{v}$
- if we pre-center the data, i.e. let $x_n' = x_n \frac{1}{N} \sum_m x_m$, then the objective simply becomes

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left({\boldsymbol{x}_n'}^{\mathrm{T}} \boldsymbol{v} \right)^2$$

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance, i.e.

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left(\boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{v} - \frac{1}{N} \sum_m \boldsymbol{x}_m^{\mathrm{T}} \boldsymbol{v}\right)^2$$

- $oldsymbol{e} oldsymbol{x}_n^{
 m T} oldsymbol{v}$ is exactly the projection of $oldsymbol{x}_n$ onto the direction $oldsymbol{v}$
- if we pre-center the data, i.e. let $x_n' = x_n \frac{1}{N} \sum_m x_m$, then the objective simply becomes

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^{N} \left(\boldsymbol{x}_n^{\prime} {}^{\mathrm{T}} \boldsymbol{v}\right)^2 = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\sum_{n=1}^{N} \boldsymbol{x}_n^{\prime} \boldsymbol{x}_n^{\prime} {}^{\mathrm{T}}\right) \boldsymbol{v}$$

More formally, we want to find a direction $v \in \mathbb{R}^D$ with $||v||_2 = 1$, so that the projection of the dataset on this direction has the most variance, i.e.

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^N \left(\boldsymbol{x}_n^{\mathrm{T}} \boldsymbol{v} - \frac{1}{N} \sum_m \boldsymbol{x}_m^{\mathrm{T}} \boldsymbol{v}\right)^2$$

- $ullet x_n^{\mathrm{T}} v$ is exactly the projection of x_n onto the direction v
- if we pre-center the data, i.e. let $x_n' = x_n \frac{1}{N} \sum_m x_m$, then the objective simply becomes

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{n=1}^{N} \left(\boldsymbol{x}_n^{\prime} {}^{\mathrm{T}} \boldsymbol{v}\right)^2 = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\sum_{n=1}^{N} \boldsymbol{x}_n^{\prime} \boldsymbol{x}_n^{\prime} {}^{\mathrm{T}}\right) \boldsymbol{v}$$

ullet we will simply assume $\{x_n\}$ is centered (to avoid notation x_n')

With $oldsymbol{X} \in \mathbb{R}^{N imes D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1}\boldsymbol{v}^{\mathrm{T}}\left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}\right)\boldsymbol{v}$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X}\right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 \lambda_i,$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\mathsf{T}} \boldsymbol{v})^2 \lambda_i,$$

Note:
$$\sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 = \boldsymbol{v}^{\top} \left(\sum_{i=1}^{D} \boldsymbol{u}_i \boldsymbol{u}_i^{\top} \right) \boldsymbol{v}$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\mathsf{T}} \boldsymbol{v})^2 \lambda_i,$$

Note:
$$\sum_{i=1}^D (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 = \boldsymbol{v}^{\top} \left(\sum_{i=1}^D \boldsymbol{u}_i \boldsymbol{u}_i^{\top} \right) \boldsymbol{v} = \boldsymbol{v}^{\top} \boldsymbol{v}$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathsf{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 \lambda_i,$$

Note:
$$\sum_{i=1}^D (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 = \boldsymbol{v}^{\top} \left(\sum_{i=1}^D \boldsymbol{u}_i \boldsymbol{u}_i^{\top} \right) \boldsymbol{v} = \boldsymbol{v}^{\top} \boldsymbol{v} = 1,$$

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 \lambda_i,$$

where the columns of U^{T} are unit **eigenvectors** u_1, \ldots, u_D of $X^{\mathrm{T}}X$ with corresponding **eigenvalues** $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_D \geq 0$.

Note: $\sum_{i=1}^D (\boldsymbol{u}_i^\top \boldsymbol{v})^2 = \boldsymbol{v}^\top \left(\sum_{i=1}^D \boldsymbol{u}_i \boldsymbol{u}_i^\top\right) \boldsymbol{v} = \boldsymbol{v}^\top \boldsymbol{v} = 1$, so the maximum above is λ_1 ,

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\mathsf{T}} \boldsymbol{v})^2 \lambda_i,$$

where the columns of U^{T} are unit **eigenvectors** u_1, \ldots, u_D of $X^{\mathrm{T}}X$ with corresponding **eigenvalues** $\lambda_1 \geq \lambda_2 \geq \cdots \lambda_D \geq 0$.

Note: $\sum_{i=1}^D (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 = \boldsymbol{v}^{\top} \left(\sum_{i=1}^D \boldsymbol{u}_i \boldsymbol{u}_i^{\top} \right) \boldsymbol{v} = \boldsymbol{v}^{\top} \boldsymbol{v} = 1$, so the maximum above is λ_1 , realized by $\boldsymbol{v} = \boldsymbol{u}_1$.

With $X \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), we want

$$\max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \left(\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X} \right) \boldsymbol{v} = \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \boldsymbol{v}^{\mathrm{T}} \boldsymbol{U}^{\mathrm{T}} \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & \lambda_{\mathrm{D}} \end{bmatrix} \boldsymbol{U} \boldsymbol{v}$$
$$= \max_{\boldsymbol{v}:\|\boldsymbol{v}\|_2=1} \sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 \lambda_i,$$

where the columns of U^{T} are unit **eigenvectors** u_{1}, \ldots, u_{D} of $X^{T}X$ with corresponding **eigenvalues** $\lambda_{1} \geq \lambda_{2} \geq \cdots \lambda_{D} \geq 0$.

Note: $\sum_{i=1}^{D} (\boldsymbol{u}_i^{\top} \boldsymbol{v})^2 = \boldsymbol{v}^{\top} \left(\sum_{i=1}^{D} \boldsymbol{u}_i \boldsymbol{u}_i^{\top} \right) \boldsymbol{v} = \boldsymbol{v}^{\top} \boldsymbol{v} = 1$, so the maximum above is λ_1 , realized by $\boldsymbol{v} = \boldsymbol{u}_1$.

Conclusion: the first PC is the top eigenvector of the covariance matrix!

Finding the other PCs

If v_1 is the first PC, then the **second PC** is found via

$$\max_{\boldsymbol{v}_2:\|\boldsymbol{v}_2\|_2=1,\boldsymbol{v}_1^{\mathrm{T}}\boldsymbol{v}_2=0}\boldsymbol{v}_2^{\mathrm{T}}\left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}\right)\boldsymbol{v}_2$$

i.e. the direction that maximizes the variance among all other dimensions

Finding the other PCs

If v_1 is the first PC, then the second PC is found via

$$\max_{\boldsymbol{v}_2:\|\boldsymbol{v}_2\|_2=1,\boldsymbol{v}_1^{\mathrm{T}}\boldsymbol{v}_2=0}\boldsymbol{v}_2^{\mathrm{T}}\left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}\right)\boldsymbol{v}_2$$

i.e. the direction that maximizes the variance among all other dimensions

This is just the second top eigenvector of the covariance matrix!

Finding the other PCs

If v_1 is the first PC, then the second PC is found via

$$\max_{\boldsymbol{v}_2:\|\boldsymbol{v}_2\|_2=1,\boldsymbol{v}_1^{\mathrm{T}}\boldsymbol{v}_2=0}\boldsymbol{v}_2^{\mathrm{T}}\left(\boldsymbol{X}^{\mathrm{T}}\boldsymbol{X}\right)\boldsymbol{v}_2$$

i.e. the direction that maximizes the variance among all other dimensions

This is just the second top eigenvector of the covariance matrix!

Conclusion: the d-th principal component is the d-th eigenvector (sorted by the eigenvalue from largest to smallest).

 ${f Input}$: a dataset represented as ${m X}$, $\#{\bf components}\ p$ we want

 ${f Input}$: a dataset represented as ${m X}$, $\#{\bf components}\ p$ we want

Step 1 Center the data by subtracting the mean

Input: a dataset represented as $oldsymbol{X}$, #components p we want

Step 1 Center the data by subtracting the mean

Step 2 Find the top p eigenvectors (with unit norm) of the covariance matrix $X^{\mathrm{T}}X$, denoted by $V \in \mathbb{R}^{\mathsf{D} \times p}$

Input: a dataset represented as X, #components p we want

Step 1 Center the data by subtracting the mean

Step 2 Find the top p eigenvectors (with unit norm) of the covariance matrix $X^{\mathrm{T}}X$, denoted by $V \in \mathbb{R}^{\mathsf{D} \times p}$

Step 3 Construct the new compressed dataset $oldsymbol{X} oldsymbol{V} \in \mathbb{R}^{N imes p}$

Input: a dataset represented as X, #components p we want

Step 1 Center the data by subtracting the mean

Step 2 Find the top p eigenvectors (with unit norm) of the covariance matrix $X^{\mathrm{T}}X$, denoted by $V \in \mathbb{R}^{\mathsf{D} \times p}$

Step 3 Construct the new compressed dataset $oldsymbol{X} oldsymbol{V} \in \mathbb{R}^{N imes p}$

(Optional) Step 4 "Reconstruct" original dataset as $m{X}m{V}m{V}^{\mathrm{T}} \in \mathbb{R}^{N imes \mathsf{D}}$

One common rule: pick p large enough so it covers about 90% of the spectrum (so 90% of variance is explained),

One common rule: pick p large enough so it covers about 90% of the spectrum (so 90% of variance is explained), i.e.

$$\frac{\sum_{d=1}^{p} \lambda_d}{\sum_{d=1}^{\mathsf{D}} \lambda_d} \ge 90\%$$

where $\lambda_1 \geq \cdots \geq \lambda_N$ are sorted eigenvalues.

One common rule: pick p large enough so it covers about 90% of the spectrum (so 90% of variance is explained), i.e.

$$\frac{\sum_{d=1}^{p} \lambda_d}{\sum_{d=1}^{\mathsf{D}} \lambda_d} \ge 90\%$$

where $\lambda_1 \geq \cdots \geq \lambda_N$ are sorted eigenvalues.

Note: $\sum_{d=1}^{D} \lambda_d = \text{Tr}(\boldsymbol{X}^T \boldsymbol{X})$, so no need to actually find all eigenvalues.

One common rule: pick p large enough so it covers about 90% of the spectrum (so 90% of variance is explained), i.e.

$$\frac{\sum_{d=1}^{p} \lambda_d}{\sum_{d=1}^{\mathsf{D}} \lambda_d} \ge 90\%$$

where $\lambda_1 \geq \cdots \geq \lambda_N$ are sorted eigenvalues.

Note: $\sum_{d=1}^{D} \lambda_d = \text{Tr}(\boldsymbol{X}^T \boldsymbol{X})$, so no need to actually find all eigenvalues.

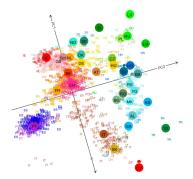
For visualization, also often pick p = 1 or p = 2.

A famous study of genetic map

• dataset: genomes of 1,387 Europeans

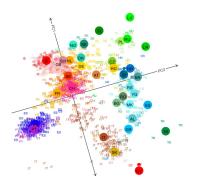
A famous study of **genetic map**

- dataset: genomes of 1,387 Europeans
- first 2 PCs shown below;



A famous study of genetic map

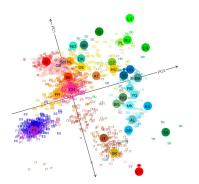
- dataset: genomes of 1,387 Europeans
- first 2 PCs shown below; looks remarkably like the geographic map





A famous study of genetic map

- dataset: genomes of 1,387 Europeans
- first 2 PCs shown below; looks remarkably like the geographic map





Genetic variation is highly structured by geography; top PCs correspond to geographic factors.

Dataset: 256×256 ($\approx 65K$ pixels) dimensional images of about 2500 faces, all framed similarly (N=2500, D $\approx 65,000$)

Dataset: 256×256 ($\approx 65K$ pixels) dimensional images of about 2500 faces, all framed similarly (N=2500, D $\approx 65,000$)

Even with $p \approx 100$, reconstructed images $\boldsymbol{X}\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} \in \mathbb{R}^{N \times D}$ look very much similar to original images

Dataset: 256×256 ($\approx 65K$ pixels) dimensional images of about 2500 faces, all framed similarly (N=2500, D $\approx 65,000$)

Even with $p \approx 100$, reconstructed images $\boldsymbol{X}\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} \in \mathbb{R}^{N \times D}$ look very much similar to original images

The principal components (called eigenfaces $\in \mathbb{R}^D$) are themselves interpretable too!



Figure 2. Seven of the eigenfaces calculated from the input images of Figure 1.

Dataset: 256×256 ($\approx 65K$ pixels) dimensional images of about 2500 faces, all framed similarly (N=2500, D $\approx 65,000$)

Even with $p \approx 100$, reconstructed images $\boldsymbol{X}\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} \in \mathbb{R}^{N \times D}$ look very much similar to original images

The principal components (called eigenfaces $\in \mathbb{R}^D$) are themselves interpretable too!

PCA finds the subspace (out of the space of all 256×256 images) where faces lie



Figure 2. Seven of the eigenfaces calculated from the input images of Figure 1.

Word embedding via PCA

Create meaningful vector representation of words; see project.

```
såd
                                                              petroleum
                                                         gas
              country
               continent
          city
     town
  village
```

Outline

- Principal Component Analysis (PCA)
- Markov models
 - Markov chain
 - Learning Markov models
- 3 Hidden Markov Model

Sequential data (i.e., each x_n is a *sequence*) are ubiquitous:

- text or speech data
- stock market data
- gene data

Sequential data (i.e., each x_n is a *sequence*) are ubiquitous:

- text or speech data
- stock market data
- gene data

Markov models are powerful probabilistic tools to analyze them:

Sequential data (i.e., each x_n is a *sequence*) are ubiquitous:

- text or speech data
- stock market data
- gene data

Markov models are powerful probabilistic tools to analyze them:

• not the state-of-the-art for e.g. Natural Language Processing (NLP)

Sequential data (i.e., each x_n is a *sequence*) are ubiquitous:

- text or speech data
- stock market data
- gene data

Markov models are powerful probabilistic tools to analyze them:

- not the state-of-the-art for e.g. Natural Language Processing (NLP)
- but still extremely useful in many areas (in particular, it is the foundation of reinforcement learning)

A motivating example

Next word prediction





A motivating example

Next word prediction





Can be formulated as predicting the probability of the next word/state:

$$P(Z_{t+1} \mid Z_1, \dots, Z_t)$$
?

A motivating example

Next word prediction





Can be formulated as predicting the probability of the next word/state:

$$P(Z_{t+1} | Z_1, \ldots, Z_t)$$
?

abbreviated as $P(Z_{t+1} \mid \mathbf{Z}_{1:t})$ (i.e., $Z_{1:t}$ denotes the sequence Z_1, \ldots, Z_t).

A Markov chain is a stochastic process with Markov property:

A Markov chain is a stochastic process with Markov property: a sequence of random variables Z_1, Z_2, \cdots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t)$$
 (Markov property)

i.e. the current state only depends on the most recent state

A Markov chain is a stochastic process with Markov property: a sequence of random variables Z_1, Z_2, \cdots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t)$$
 (Markov property)

i.e. the current state only depends on the most recent state

We only consider the following case:

• All Z_t 's take value from the same discrete set $\{1, \ldots, S\}$

A Markov chain is a stochastic process with Markov property: a sequence of random variables Z_1, Z_2, \cdots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t)$$
 (Markov property)

i.e. the current state only depends on the most recent state

We only consider the following case:

- All Z_t 's take value from the same discrete set $\{1,\ldots,S\}$
- $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$, known as transition probability

A Markov chain is a stochastic process with Markov property: a sequence of random variables Z_1, Z_2, \cdots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t)$$
 (Markov property)

i.e. the current state only depends on the most recent state

We only consider the following case:

- All Z_t 's take value from the same discrete set $\{1, \ldots, S\}$
- $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$, known as transition probability
- $P(Z_1 = s) = \pi_s$

A Markov chain is a stochastic process with Markov property: a sequence of random variables Z_1, Z_2, \cdots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t)$$
 (Markov property)

i.e. the current state only depends on the most recent state

We only consider the following case:

- All Z_t 's take value from the same discrete set $\{1,\ldots,S\}$
- $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$, known as transition probability
- $P(Z_1 = s) = \pi_s$
- \bullet $(\{\pi_s\}, \{a_{s,s'}\}) = (\boldsymbol{\pi}, \boldsymbol{A})$ are parameters of the model

Examples

• Example 1 (Language model)

States [S] represent a dictionary of words,

$$a_{\mathsf{ice},\mathsf{cream}} = P(Z_{t+1} = \mathsf{cream} \mid Z_t = \mathsf{ice})$$

is an example of the transition probability.

Examples

• Example 1 (Language model)

States [S] represent a dictionary of words,

$$a_{\mathsf{ice},\mathsf{cream}} = P(Z_{t+1} = \mathsf{cream} \mid Z_t = \mathsf{ice})$$

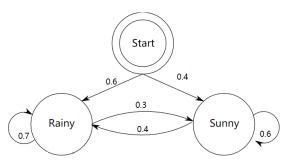
is an example of the transition probability.

• Example 2 (Weather)

States [S] represent weather at each day

$$a_{\text{sunny,rainy}} = P(Z_{t+1} = \text{rainy} \mid Z_t = \text{sunny})$$

It is intuitive to represent a Markov model as a graph



Is the Markov assumption reasonable?

Is the Markov assumption reasonable? Not completely for the language model for example:

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_t = \text{the})$ (bigram, unreasonable)

Is the Markov assumption reasonable? Not completely for the language model for example:

$$P(Z_{t+1} = ext{office} \mid Z_{1:t} = ext{I'II meet you at the})$$
 $= P(Z_{t+1} = ext{office} \mid Z_t = ext{the})$ (bigram, unreasonable)

Higher-order Markov chains make it more reasonable, e.g.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t, Z_{t-1})$$
 (second-order Markov)

Is the Markov assumption reasonable? Not completely for the language model for example:

$$P(Z_{t+1} = ext{office} \mid Z_{1:t} = ext{I'II meet you at the})$$
 $= P(Z_{t+1} = ext{office} \mid Z_t = ext{the})$ (bigram, unreasonable)

Higher-order Markov chains make it more reasonable, e.g.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t, Z_{t-1})$$
 (second-order Markov)

i.e. the current word only depends on the last two words (the **trigram** model).

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-1:t} = \text{at the})$ (trigram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'll meet you at the})$$
 $= P(Z_{t+1} = \text{office} \mid Z_{t-1:t} = \text{at the})$ (trigram)

$$\begin{split} &P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the}) \\ &= P(Z_{t+1} = \text{office} \mid Z_{t-2:t} = \text{you at the}) \end{split} \tag{4-gram}$$

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'll meet you at the})$$

$$= P(Z_{t+1} = \text{office} \mid Z_{t-1:t} = \text{at the}) \qquad \qquad \text{(trigram)}$$

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'll meet you at the})$$

$$= P(Z_{t+1} = \text{office} \mid Z_{t-2:t} = \text{you at the}) \qquad \qquad \text{(4-gram)}$$

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-3:t} = \text{meet you at the})$ (5-gram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-1:t} = \text{at the})$ (trigram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

$$= P(Z_{t+1} = \text{office} \mid Z_{t-2:t} = \text{you at the})$$
(4-gram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-3:t} = \text{meet you at the})$ (5-gram)

Learning higher-order Markov chains is similar, but more expensive.

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-1:t} = \text{at the})$ (trigram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'll meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-2:t} = \text{you at the})$ (4-gram)

$$P(Z_{t+1} = \text{office} \mid Z_{1:t} = \text{I'II meet you at the})$$

= $P(Z_{t+1} = \text{office} \mid Z_{t-3:t} = \text{meet you at the})$ (5-gram)

Learning higher-order Markov chains is similar, but more *expensive*.

For simplicity, we will only consider standard Markov chains.

Now suppose we have observed N sequences of examples:

- $z_{1,1},\ldots,z_{1,T}$
-
- $z_{n,1}, \ldots, z_{n,T}$
-
- \bullet $z_{N,1},\ldots,z_{N,T}$

Now suppose we have observed N sequences of examples:

- $z_{1,1}, \ldots, z_{1,T}$
- ...
- \bullet $z_{n,1},\ldots,z_{n,T}$
- ...
- \bullet $z_{N,1},\ldots,z_{N,T}$

where

ullet for simplicity we assume each sequence has the same length T

Now suppose we have observed N sequences of examples:

- $z_{1,1},\ldots,z_{1,T}$
-
- \bullet $z_{n,1},\ldots,z_{n,T}$
- ...
- \bullet $z_{N.1},\ldots,z_{N.T}$

where

- ullet for simplicity we assume each sequence has the same length T
- ullet lower case $z_{n,t}$ represents the value of the random variable $Z_{n,t}$

Now suppose we have observed N sequences of examples:

- $z_{1,1}, \ldots, z_{1,T}$
-
- \bullet $z_{n,1},\ldots,z_{n,T}$
-
- \bullet $z_{N,1},\ldots,z_{N,T}$

where

- ullet for simplicity we assume each sequence has the same length T
- ullet lower case $z_{n,t}$ represents the value of the random variable $Z_{n,t}$

From these observations how do we *learn the model parameters* (π, A) ?

Same story, find the MLE.

Same story, find the MLE. The log-likelihood of a sequence z_1,\dots,z_T is

$$ln P(Z_{1:T} = z_{1:T})$$

Same story, find the MLE. The log-likelihood of a sequence z_1, \ldots, z_T is

$$\ln P(Z_{1:T} = z_{1:T})$$

$$= \sum_{t=0}^{T} \ln P(Z_t = z_t \mid Z_{1:t-1} = z_{1:t-1})$$
 (always true)

Same story, find the MLE. The log-likelihood of a sequence z_1, \ldots, z_T is

$$\ln P(Z_{1:T}=z_{1:T})$$

$$=\sum_{t=1}^T \ln P(Z_t=z_t\mid Z_{1:t-1}=z_{1:t-1})$$
 (always true)
$$=\sum_{t=1}^T \ln P(Z_t=z_t\mid Z_{t-1}=z_{t-1})$$
 (Markov property)

Same story, find the MLE. The log-likelihood of a sequence z_1, \ldots, z_T is

$$\ln P(Z_{1:T} = z_{1:T})$$

$$= \sum_{t=1}^{T} \ln P(Z_t = z_t \mid Z_{1:t-1} = z_{1:t-1}) \qquad \text{(always true)}$$

$$= \sum_{t=1}^{T} \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) \qquad \text{(Markov property)}$$

$$= \ln \pi_{z_1} + \sum_{t=2}^{T} \ln a_{z_{t-1}, z_t}$$

Same story, find the MLE. The log-likelihood of a sequence z_1,\ldots,z_T is

$$\begin{split} & \ln P(Z_{1:T} = z_{1:T}) \\ & = \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{1:t-1} = z_{1:t-1}) \\ & = \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) \\ & = \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1},z_t} \\ & = \sum_s \mathbb{I}[z_1 = s] \ln \pi_s + \sum_{s,s'} \left(\sum_{t=2}^T \mathbb{I}[z_{t-1} = s, z_t = s'] \right) \ln a_{s,s'} \end{split}$$

Summing over all N sequences, we obtain the joint log-likelihood:

$$L(\pi, \mathbf{A}) = \sum_s (\text{\#initial states with value } s) \ln \pi_s$$

$$+ \sum_{s,s'} (\text{\#transitions from } s \text{ to } s') \ln a_{s,s'}$$

Summing over all N sequences, we obtain the joint log-likelihood:

$$L(\pi, \mathbf{A}) = \sum_s (\text{\#initial states with value } s) \ln \pi_s$$

$$+ \sum_{s,s'} (\text{\#transitions from } s \text{ to } s') \ln a_{s,s'}$$

The MLE $\operatorname{argmax}_{\boldsymbol{\pi},\boldsymbol{A}}L(\boldsymbol{\pi},\boldsymbol{A})$ is:

$$\pi_s \propto \# \text{initial states with value } s$$
 $a_{s,s'} \propto \# \text{transitions from } s \text{ to } s'$

Example

Suppose we observed the following 2 sequences of length 5

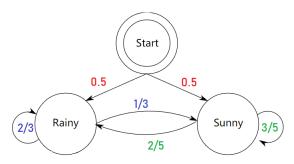
- sunny, sunny, rainy, rainy, rainy
- rainy, sunny, sunny, rainy

Example

Suppose we observed the following 2 sequences of length 5

- sunny, sunny, rainy, rainy, rainy
- rainy, sunny, sunny, rainy

MLE is the following model



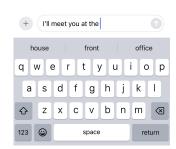
Recall: I'll meet you at the _____.



Recall: I'll meet you at the _____.

Suppose

we use a 5-gram model



Recall: I'll meet you at the _____.

Suppose

- we use a 5-gram model
- "meet you at the" appears 1000 times in the training set:
 - "meet you at the house" appears 500 times
 - "meet you at the front" appears 300 times
 - "meet you at the office" appears 200 times



Recall: I'll meet you at the _____.

Suppose

- we use a 5-gram model
- "meet you at the" appears 1000 times in the training set:
 - "meet you at the house" appears 500 times
 - "meet you at the front" appears 300 times
 - "meet you at the office" appears 200 times

Then the model predicts the next word as "house" with prob. 0.5, "front" with prob. 0.3, and "office" with prob. 0.2.



adapted from Stanford CS224n

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

adapted from Stanford CS224n

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

today the _____

adapted from Stanford CS224n

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

- today the _____
- today the price _____

company 0.153 bank 0.153 price 0.077 italian 0.039 emirate 0.039 After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

- today the _____
- today the price _____
- today the price of _____

company bank	0.153 0.153
price	0.077
italian	0.039
emirate	0.039

of	0.308	
for	0.050	
it	0.046	
to	0.046	
is	0.031	
l .		

adapted from Stanford CS224n

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

- today the _____
- today the price _____
- today the price of _____
- today the price of gold _____

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039

of	0.308	
for	0.050	
it	0.046	
to	0.046	
is	0.031	

the	0.072	
18	0.043	
oil	0.043	
its	0.036	
gold	0.018	

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

- today the _____
- today the price _____
- today the price of _____
- today the price of gold _____

final result: today the price of gold per ton, while production of shoe lasts and shoe industry, the bank intervened just after it considered and rejected an IMF demand to rebuild depleted European stocks.

company bank	0.153 0.153
price	0.077
italian	0.039
emirate	0.039

of	0.308	
for	0.050	
it	0.046	
to	0.046	
is	0.031	

the	0.072	
18	0.043	
oil	0.043	
its	0.036	
gold	0.018	

adapted from Stanford CS224n

After learning the model, can also **generate texts** by repeatedly sampling conditioning on what have been generated:

- today the _____
- today the price _____
- today the price of _____
- today the price of gold _____

final result: today the price of gold per ton, while production of shoe lasts and shoe industry, the bank intervened just after it considered and rejected an IMF demand to rebuild depleted European stocks.

Surprisingly grammatical! but incoherent...

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039

of	0.308	
for	0.050	
it	0.046	
to	0.046	
is	0.031	

the	0.072	
18	0.043	
oil	0.043	
its	0.036	
gold	0.018	

Outline

- 1 Principal Component Analysis (PCA)
- Markov models
- Hidden Markov Model
 - Inferring HMMs
 - Learning HMMs

Markov model with outcomes/observations

Now suppose each state Z_t also "emits" some **outcome/observation** $X_t \in [O]$ based on the following model

$$P(X_t = o \mid Z_t = s) = b_{s,o}$$
 (emission probability)

independent of anything else.

Markov model with outcomes/observations

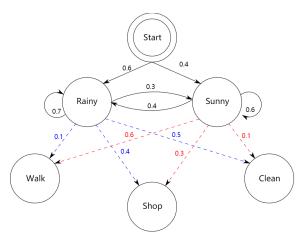
Now suppose each state Z_t also "emits" some **outcome/observation** $X_t \in [O]$ based on the following model

$$P(X_t = o \mid Z_t = s) = b_{s,o}$$
 (emission probability)

independent of anything else.

The model parameters are $(\{\pi_s\}, \{a_{s,s'}\}, \{b_{s,o}\}) = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B}).$

On each day, we also observe **Bob's activity: walk, shop, or clean**, which only depends on the weather of that day.



The joint log-likelihood of a state-outcome sequence $z_1, x_1, \dots, z_T, x_T$ is

$$\ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T})$$

The joint log-likelihood of a state-outcome sequence $z_1, x_1, \dots, z_T, x_T$ is

$$\begin{split} & \ln P(Z_{1:T}=z_{1:T},X_{1:T}=x_{1:T})\\ & = \ln P(Z_{1:T}=z_{1:T}) + \ln P(X_{1:T}=x_{1:T}\mid Z_{1:T}=z_{1:T}) \quad \text{(always true)} \end{split}$$

The joint log-likelihood of a state-outcome sequence $z_1, x_1, \ldots, z_T, x_T$ is

$$\begin{split} & \ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \\ & = \ln P(Z_{1:T} = z_{1:T}) + \ln P(X_{1:T} = x_{1:T} \mid Z_{1:T} = z_{1:T}) \quad \text{(always true)} \\ & = \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) + \sum_{t=1}^T \ln P(X_t = x_t \mid Z_t = z_t) \\ & \qquad \qquad \text{(due to all the independence)} \end{split}$$

The joint log-likelihood of a state-outcome sequence $z_1, x_1, \ldots, z_T, x_T$ is

$$\begin{split} \ln P(Z_{1:T} &= z_{1:T}, X_{1:T} = x_{1:T}) \\ &= \ln P(Z_{1:T} = z_{1:T}) + \ln P(X_{1:T} = x_{1:T} \mid Z_{1:T} = z_{1:T}) \quad \text{(always true)} \\ &= \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) + \sum_{t=1}^T \ln P(X_t = x_t \mid Z_t = z_t) \\ &\qquad \qquad \qquad \text{(due to all the independence)} \\ &= \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1}, z_t} + \sum_{t=1}^T \ln b_{z_t, x_t} \end{split}$$

Learning the model

If we observe N state-outcome sequences: $z_{n,1}, x_{n,1}, \ldots, z_{n,T}, x_{n,T}$ for $n=1,\ldots,N$, the MLE is again very simple (verify yourself):

```
\pi_s \propto #initial states with value s a_{s,s'} \propto #transitions from s to s' b_{s,o} \propto #state-outcome pairs (s,o)
```

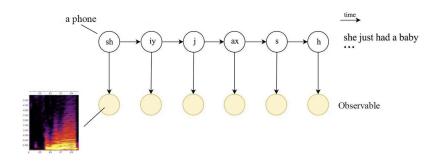
Hidden Markov models

However, *most often we do not observe the states!* Hence the name **Hidden Markov Model (HMM)**

Hidden Markov models

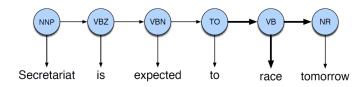
However, *most often we do not observe the states!* Hence the name **Hidden Markov Model (HMM)**

• speech recognition: observe the speech $X_{1:T}$ but not the underlying words/phones $Z_{1:T}$



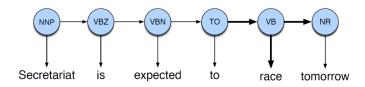
Hidden Markov models (cont.)

• Part of Speech (POS) tagging: observe the sentence $X_{1:T}$ but not the underlying categories $Z_{1:T}$ (noun, verb, adjective, adverb, etc.)



Hidden Markov models (cont.)

• Part of Speech (POS) tagging: observe the sentence $X_{1:T}$ but not the underlying categories $Z_{1:T}$ (noun, verb, adjective, adverb, etc.)



See programming project!

Learning HMMs

How to learn HMMs?

Learning HMMs

How to learn HMMs? **Roadmap**:

first discuss how to infer when the model is known (key: dynamic programming)

Learning HMMs

How to learn HMMs? **Roadmap**:

- first discuss how to infer when the model is known (key: dynamic programming)
- then discuss how to **learn** the model (key: EM)

What can we infer about an HMM?

Knowing the parameter of an HMM, we can infer

most likely hidden states path, given an observation sequence

$$\operatorname*{argmax}_{z_{1:T}} P(Z_{1:T} = z_{1:T} \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, what's the most likely weather for this week?

What can we infer about an HMM?

Knowing the parameter of an HMM, we can infer

most likely hidden states path, given an observation sequence

$$\operatorname*{argmax}_{z_{1:T}} P(Z_{1:T} = z_{1:T} \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, what's the most likely weather for this week?

• the probability of observing some sequence

$$P(X_{1:T} = x_{1:T})$$

e.g. prob. of observing Bob's activities "walk, walk, shop, clean, walk, shop, shop" for one week

What can we infer for a known HMM?

Knowing the parameter of an HMM, we can infer

• the state at some point, given an observation sequence

$$P(Z_t = s \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, how was the weather like on Wed?

What can we infer for a known HMM?

Knowing the parameter of an HMM, we can infer

• the state at some point, given an observation sequence

$$P(Z_t = s \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, how was the weather like on Wed?

• the transition at some point, given an observation sequence

$$P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, how was the weather like on Wed and Thu?

Forward and backward messages

The key to infer all these is to compute two things:

Forward and backward messages

The key to infer all these is to compute two things:

ullet forward messages: for each s and t

$$\alpha_s(t) = P(Z_t = s, X_{1:t} = x_{1:t})$$

Forward and backward messages

The key to infer all these is to compute two things:

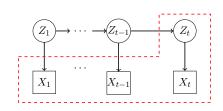
ullet forward messages: for each s and t

$$\alpha_s(t) = P(Z_t = s, X_{1:t} = x_{1:t})$$

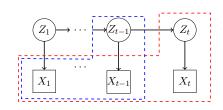
ullet backward messages: for each s and t

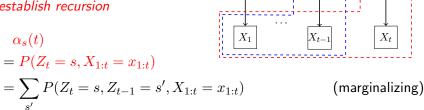
$$\beta_s(t) = P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s)$$

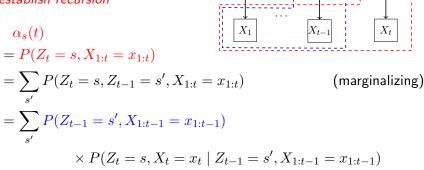
$$\alpha_s(t)$$
= $P(Z_t = s, X_{1:t} = x_{1:t})$

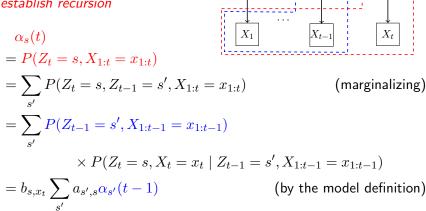


$$\alpha_s(t)$$
= $P(Z_t = s, X_{1:t} = x_{1:t})$









Key: establish recursion

establish recursion
$$\alpha_{s}(t) = P(Z_{t} = s, X_{1:t} = x_{1:t})$$

$$= \sum_{s'} P(Z_{t} = s, Z_{t-1} = s', X_{1:t} = x_{1:t})$$
 (marginalizing)
$$= \sum_{s'} P(Z_{t-1} = s', X_{1:t-1} = x_{1:t-1})$$

$$\times P(Z_{t} = s, X_{t} = x_{t} \mid Z_{t-1} = s', X_{1:t-1} = x_{1:t-1})$$

$$= b_{s,x_{t}} \sum_{s'} a_{s',s} \alpha_{s'}(t-1)$$
 (by the model definition)

Base case: $\alpha_s(1) = P(Z_1 = s, X_1 = x_1) = \pi_s b_{s,x_1}$

Forward procedure

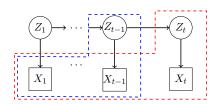
Forward procedure

For all $s \in [S]$, compute $\alpha_s(1) = \pi_s b_{s,x_1}$.

For $t = 2, \ldots, T$

• for each $s \in [S]$, compute

$$\alpha_{s}(t) = b_{s,x_t} \sum_{s'} a_{s',s} \alpha_{s'}(t-1)$$



Forward procedure

Forward procedure

For all $s \in [S]$, compute $\alpha_s(1) = \pi_s b_{s,x_1}$.

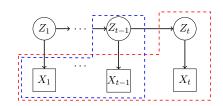
For $t = 2, \ldots, T$

ullet for each $s \in [S]$, compute

$$\alpha_{s}(t) = b_{s,x_t} \sum_{s'} a_{s',s} \alpha_{s'}(t-1)$$

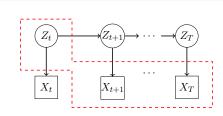
It takes

- \bullet $O(S^2T)$ time
- \bullet O(ST) space



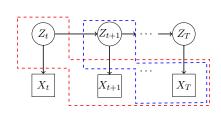
$$\beta_s(t)$$

= $P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s)$



$$\beta_s(t)$$

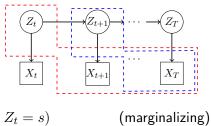
= $P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s)$

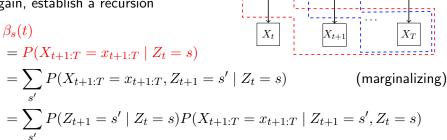


$$\beta_{s}(t)$$

$$= P(X_{t+1:T} = x_{t+1:T} \mid Z_{t} = s)$$

$$= \sum_{t} P(X_{t+1:T} = x_{t+1:T}, Z_{t+1} = s' \mid Z_{t} = s)$$





$$\beta_{s}(t) = P(X_{t+1:T} = x_{t+1:T} \mid Z_{t} = s)$$

$$= \sum_{s'} P(X_{t+1:T} = x_{t+1:T}, Z_{t+1} = s' \mid Z_{t} = s)$$
 (marginalizing)
$$= \sum_{s'} P(Z_{t+1} = s' \mid Z_{t} = s) P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s', Z_{t} = s)$$

$$= \sum_{s'} a_{s,s'} P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s') P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s')$$

gain, establish a recursion
$$\beta_{s}(t) = P(X_{t+1:T} = x_{t+1:T} \mid Z_{t} = s)$$

$$= \sum_{s'} P(X_{t+1:T} = x_{t+1:T}, Z_{t+1} = s' \mid Z_{t} = s) \qquad \text{(marginalizing)}$$

$$= \sum_{s'} P(Z_{t+1} = s' \mid Z_{t} = s) P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s', Z_{t} = s)$$

$$= \sum_{s'} a_{s,s'} P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s') P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s')$$

$$= \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$

Computing backward messages

Again, establish a recursion

gain, establish a recursion
$$\beta_{s}(t) = P(X_{t+1:T} = x_{t+1:T} \mid Z_{t} = s)$$

$$= \sum_{s'} P(X_{t+1:T} = x_{t+1:T}, Z_{t+1} = s' \mid Z_{t} = s)$$
 (marginalizing)
$$= \sum_{s'} P(Z_{t+1} = s' \mid Z_{t} = s) P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s', Z_{t} = s)$$

$$= \sum_{s'} a_{s,s'} P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s') P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s')$$

$$= \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$

Base case: $\beta_s(T) = 1$

Backward procedure

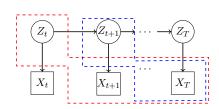
Backward procedure

For all $s \in [S]$, set $\beta_s(T) = 1$.

For
$$t = T - 1, ..., 1$$

• for each $s \in [S]$, compute

$$\beta_s(t) = \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$



Backward procedure

Backward procedure

For all $s \in [S]$, set $\beta_s(T) = 1$.

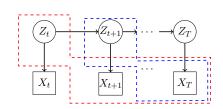
For
$$t = T - 1, ..., 1$$

• for each $s \in [S]$, compute

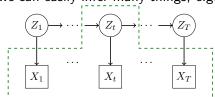
$$\beta_s(t) = \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$

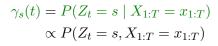
Again it takes

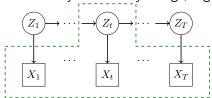
- \bullet $O(S^2T)$ time
- \bullet O(ST) space

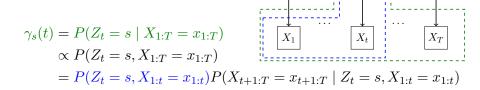


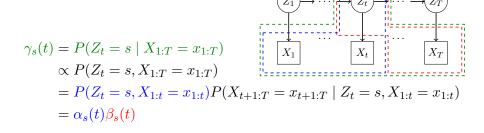
$$\gamma_s(t) = P(Z_t = s \mid X_{1:T} = x_{1:T})$$



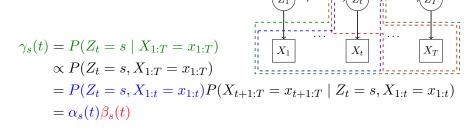






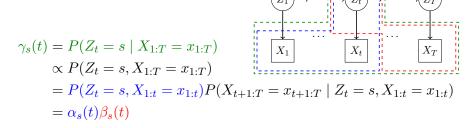


With forward and backward messages, we can easily infer many things, e.g.



What constant are we omitting in " \propto "?

With forward and backward messages, we can easily infer many things, e.g.

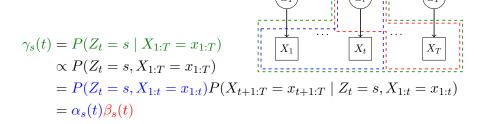


What constant are we omitting in " \propto "? It is exactly

$$P(X_{1:T} = x_{1:T}) = \sum_{s} \alpha_s(t)\beta_s(t),$$

the probability of observing the sequence $x_{1:T}$.

With forward and backward messages, we can easily infer many things, e.g.



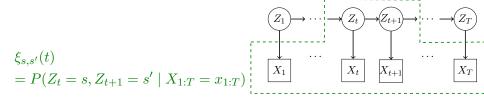
What constant are we omitting in " \propto "? It is exactly

$$P(X_{1:T} = x_{1:T}) = \sum_{s} \alpha_s(t)\beta_s(t),$$

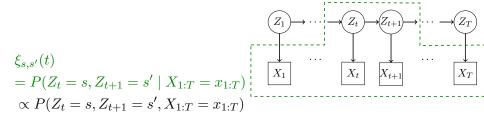
the probability of observing the sequence $x_{1:T}$.

This is true for any t; a good way to check correctness of your code.

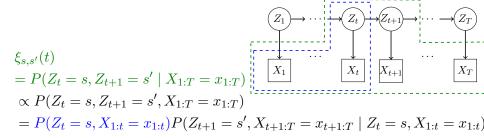
Another example: the conditional probability of transition s to s' at time t



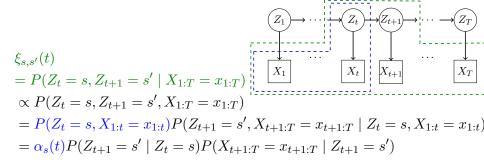
Another example: the conditional probability of transition s to s^\prime at time t



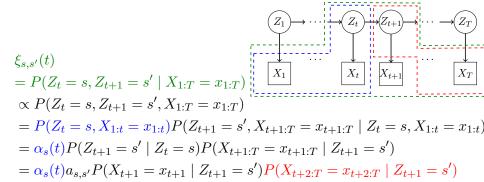
Another example: the conditional probability of transition s to s' at time t



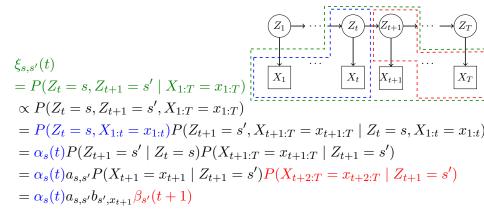
Another example: the conditional probability of transition s to s' at time t



Another example: the conditional probability of transition s to s^\prime at time t



Another example: the conditional probability of transition s to s^\prime at time t



Another example: the conditional probability of transition s to s^\prime at time t

$$\xi_{s,s'}(t) = P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T})$$

$$\propto P(Z_t = s, Z_{t+1} = s', X_{1:T} = x_{1:T})$$

$$= P(Z_t = s, Z_{t+1} = s', X_{1:T} = x_{1:T})$$

$$= P(Z_t = s, X_{1:t} = x_{1:t}) P(Z_{t+1} = s', X_{t+1:T} = x_{t+1:T} \mid Z_t = s, X_{1:t} = x_{1:t})$$

$$= \alpha_s(t) P(Z_{t+1} = s' \mid Z_t = s) P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s')$$

$$= \alpha_s(t) a_{s,s'} P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s') P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s')$$

$$= \alpha_s(t) a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$

The normalization constant is in fact again $P(X_{1:T} = x_{1:T})$

Finding the most likely path

Can't use forward and backward messages directly to find the most likely path, but it is very similar to the forward procedure.

Finding the most likely path

Can't use forward and backward messages directly to find the most likely path, but it is very similar to the forward procedure. Key: compute

$$\delta_s(t) = \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t})$$

the probability of the most likely path for time 1:t ending at state s

$$\delta_s(t) = \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t})$$

$$\begin{split} \delta_s(t) &= \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} \max_{z_{1:t-2}} P(Z_t = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t}) \end{split}$$

$$\delta_{s}(t) = \max_{z_{1:t-1}} P(Z_{t} = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t})$$

$$= \max_{s'} \max_{z_{1:t-2}} P(Z_{t} = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t})$$

$$= \max_{s'} P(Z_{t} = s \mid Z_{t-1} = s') P(X_{t} = x_{t} \mid Z_{t} = s) \cdot$$

$$\max_{z_{1:t-2}} P(Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t-1} = x_{1:t-1})$$

$$\begin{split} \delta_{s}(t) &= \max_{z_{1:t-1}} P(Z_{t} = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} \max_{z_{1:t-2}} P(Z_{t} = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} P(Z_{t} = s \mid Z_{t-1} = s') P(X_{t} = x_{t} \mid Z_{t} = s) \cdot \\ &\qquad \qquad \max_{s'} P(Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t-1} = x_{1:t-1}) \\ &= b_{s,x_{t}} \max_{s'} a_{s',s} \delta_{s'}(t-1) \end{split} \tag{recursion!}$$

$$\begin{split} \delta_s(t) &= \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} \max_{z_{1:t-2}} P(Z_t = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} P(Z_t = s \mid Z_{t-1} = s') P(X_t = x_t \mid Z_t = s) \cdot \\ &\qquad \qquad \max_{s'} P(Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t-1} = x_{1:t-1}) \\ &= b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1) \end{split} \tag{recursion!}$$

Base case:
$$\delta_s(1) = P(Z_1 = s, X_1 = x_1) = \pi_s b_{s,x_1}$$

Observe

$$\begin{split} \delta_s(t) &= \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} \max_{z_{1:t-2}} P(Z_t = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} P(Z_t = s \mid Z_{t-1} = s') P(X_t = x_t \mid Z_t = s) \cdot \\ &\qquad \qquad \max_{s'} P(Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t-1} = x_{1:t-1}) \\ &= b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1) \end{split} \tag{recursion!}$$

Base case:
$$\delta_s(1) = P(Z_1 = s, X_1 = x_1) = \pi_s b_{s,x_1}$$

Exactly the same as forward messages except replacing "sum" by "max"!

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

Viterbi Algorithm

For each
$$s \in [S]$$
, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each
$$t = 2, \dots, T$$
,

ullet for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1),$$

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \ldots, T$,

• for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1),$$

$$\Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \ldots, T$,

• for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1),$$

$$\Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \ldots, T$,

• for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1),$$

$$\Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \ldots, T$,

• for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1),$$

$$\Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each $t=T,\ldots,\hat{2}$: set $z_{t-1}^*=\Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \ldots, z_T^* .

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

Viterbi Algorithm

For each
$$s \in [S]$$
, compute $\delta_s(1) = \pi_s b_{s,x_1}$

$$\delta_{\text{sunny}}(1) = \pi_{\text{sunny}} b_{\text{sunny,clean}}$$

Viterbi Algorithm

For each
$$s \in [S]$$
, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

$$\delta_{\text{sunny}}(1) = \pi_{\text{sunny}} b_{\text{sunny,clean}}$$

$$\delta_{\text{rainy}}(1) = \pi_{\text{rainy}} b_{\text{rainy,clean}}$$

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

$$\delta_{\text{sunny}}(1) = \pi_{\text{sunny}} b_{\text{sunny,clean}}$$

$$\delta_{\text{rainy}}(1) = \pi_{\text{rainy}} b_{\text{rainy,clean}}$$

$$\delta_{\rm sunny}(1) = 0.25$$

$$\delta_{\text{rainy}}(1) = 0.4$$

(numbers are made up here)

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\rm sunny}(1) = 0.25$$

 $\delta_{\text{rainv}}(1) = 0.4$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(2) = b_{\text{sunny,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(1)$$

$$\delta_{\text{sunny}}(1) = 0.25$$

$$\delta_{\text{rainy}}(1) = 0.4$$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(2) = b_{\text{sunny,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(1)$$

$$\delta_{\text{rainy}}(2) = b_{\text{rainy,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(1)$$

$$\delta_{\text{sunny}}(1) = 0.25$$

 $\delta_{\text{rainy}}(1) = 0.4$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(2) = b_{\text{sunny,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(1)$$

$$\delta_{\text{rainy}}(2) = b_{\text{rainy,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(1)$$

$$\delta_{\text{sunny}}(1) = 0.25$$

$$\delta_{\text{sunny}}(2) = 0.1$$

$$\delta_{\text{rainy}}(1) = 0.4$$

$$\delta_{\text{rainy}}(2) = 0.19$$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(2) = b_{\text{sunny,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(1)$$

$$\delta_{\text{rainy}}(2) = b_{\text{rainy,shop}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(1)$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$

$$\delta_{\text{rainy}}(1) = 0.4$$
 $\delta_{\text{rainy}}(2) = 0.19$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(3) = b_{\text{sunny,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(2)$$

$$\delta_{\text{rainy}}(3) = b_{\text{rainy,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(2)$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$

$$\delta_{\text{rainy}}(1) = 0.4$$

$$\delta_{\text{rainy}}(2) = 0.19$$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(3) = b_{\text{sunny,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(2)$$

$$\delta_{\text{rainy}}(3) = b_{\text{rainy,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(2)$$

$$\delta_{\text{sunny}}(1) = 0.25 \quad \qquad \qquad \delta_{\text{sunny}}(2) = 0.1$$

$$\delta_{\text{sunny}}(3) = 0.04$$

$$\delta_{\text{rainy}}(1) = 0.4 \quad \leftarrow \quad \delta_{\text{rainy}}(2) = 0.19$$

$$\delta_{\rm rainy}(3) = 0.042$$

Viterbi Algorithm

. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(3) = b_{\text{sunny,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(2)$$

$$\delta_{\text{rainy}}(3) = b_{\text{rainy,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(2)$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$
 $\delta_{\text{sunny}}(3) = 0.04$
 $\delta_{\text{rainy}}(1) = 0.4$
 $\delta_{\text{rainy}}(2) = 0.19$

Viterbi Algorithm

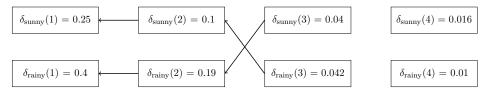
. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(4) = b_{\text{sunny,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(3)$$

$$\delta_{\text{rainy}}(4) = b_{\text{rainy,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(3)$$



Viterbi Algorithm

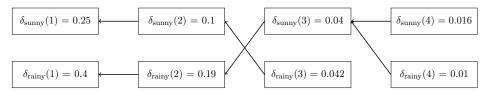
. . .

For each $t=2,\ldots,T$ and each $s\in[S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1), \quad \Delta_s(t) = \operatorname*{argmax}_{s'} a_{s',s} \delta_{s'}(t-1).$$

$$\delta_{\text{sunny}}(4) = b_{\text{sunny,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{sunny}} \delta_{s'}(3)$$

$$\delta_{\text{rainy}}(4) = b_{\text{rainy,walk}} \max_{s' \in \{\text{sunny,rainy}\}} a_{s',\text{rainy}} \delta_{s'}(3)$$



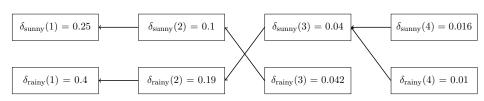
Viterbi Algorithm

. . .

Backtracking: let
$$z_T^* = \operatorname{argmax}_s \delta_s(T)$$
.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \ldots, z_T^* .



The most likely path is "

Viterbi Algorithm

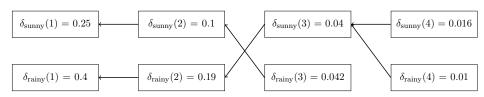
. . .

Backtracking: let
$$z_T^* = \operatorname{argmax}_s \delta_s(T)$$
.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_4^* = \underset{s \in \{\text{sunny,rainy}\}}{\operatorname{argmax}} \delta_s(4)$$



The most likely path is "

Viterbi Algorithm

. . .

Backtracking: let
$$z_T^* = \operatorname{argmax}_s \delta_s(T)$$
.

For each $t=T,\ldots,2$: set $z_{t-1}^*=\Delta_{z_t^*}(t).$

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_4^* = \underset{s \in \{\text{sunny,rainy}\}}{\operatorname{argmax}} \delta_s(4) = \text{sunny},$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$
 $\delta_{\text{sunny}}(3) = 0.04$
 $\delta_{\text{sunny}}(4) = 0.016$
 $\delta_{\text{rainy}}(1) = 0.4$
 $\delta_{\text{rainy}}(2) = 0.19$
 $\delta_{\text{rainy}}(3) = 0.042$

The most likely path is "

sunny"

Viterbi Algorithm

. . .

Backtracking: let
$$z_T^* = \operatorname{argmax}_s \delta_s(T)$$
.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_4^* = \underset{s \in \{\text{sunny, rainy}\}}{\operatorname{argmax}} \delta_s(4) = \text{sunny}, \quad z_3^* = \Delta_{\text{sunny}}(4)$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$
 $\delta_{\text{sunny}}(3) = 0.04$
 $\delta_{\text{sunny}}(4) = 0.016$
 $\delta_{\text{rainy}}(1) = 0.4$
 $\delta_{\text{rainy}}(2) = 0.19$
 $\delta_{\text{rainy}}(3) = 0.042$

The most likely path is "

sunny"

Viterbi Algorithm

. . .

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each $t=T,\ldots,2$: set $z_{t-1}^*=\Delta_{z_t^*}(t)$.

(just follow the arrow!)

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_4^* = \underset{s \in \{\text{sunny, rainy}\}}{\operatorname{argmax}} \delta_s(4) = \text{sunny}, \quad z_3^* = \Delta_{\text{sunny}}(4) = \text{sunny}$$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$
 $\delta_{\text{sunny}}(3) = 0.04$
 $\delta_{\text{sunny}}(4) = 0.016$
 $\delta_{\text{rainy}}(1) = 0.4$
 $\delta_{\text{rainy}}(2) = 0.19$
 $\delta_{\text{rainy}}(3) = 0.042$

The most likely path is "

sunny, sunny"

Viterbi Algorithm

. . .

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each $t = T, \dots, 2$: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

(just follow the arrow!)

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_4^* = \underset{s \in \{\text{sunny, rainy}\}}{\operatorname{argmax}} \delta_s(4) = \text{sunny}, \quad z_3^* = \Delta_{\text{sunny}}(4) = \text{sunny}$$

 $z_2^* = \Delta_{\text{sunny}}(3) = \text{rainy},$

$$\delta_{\text{sunny}}(1) = 0.25$$
 $\delta_{\text{sunny}}(2) = 0.1$
 $\delta_{\text{sunny}}(3) = 0.04$
 $\delta_{\text{sunny}}(4) = 0.016$
 $\delta_{\text{rainy}}(1) = 0.4$
 $\delta_{\text{rainy}}(2) = 0.19$
 $\delta_{\text{rainy}}(3) = 0.042$

The most likely path is " rainy, sunny, sunny"

Viterbi Algorithm

. . .

Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each t = T, ..., 2: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

(just follow the arrow!)

Output the most likely path z_1^*, \ldots, z_T^* .

$$z_{2}^{*} = \Delta_{\text{sunny}}(3) = \text{rainy}, \qquad z_{1}^{*} = \Delta_{\text{rainy}}(2) = \text{rainy}$$

$$\delta_{\text{sunny}}(1) = 0.25 \qquad \delta_{\text{sunny}}(2) = 0.1 \qquad \delta_{\text{sunny}}(3) = 0.04 \qquad \delta_{\text{sunny}}(4) = 0.016$$

$$\delta_{\text{rainy}}(1) = 0.4 \qquad \delta_{\text{rainy}}(2) = 0.19 \qquad \delta_{\text{rainy}}(3) = 0.042 \qquad \delta_{\text{rainy}}(4) = 0.01$$

 $z_4^* = \operatorname{argmax} \delta_s(4) = \operatorname{sunny}, \quad z_3^* = \Delta_{\operatorname{sunny}}(4) = \operatorname{sunny}$

The most likely path is "rainy, rainy, sunny, sunny"

All previous inferences depend on knowing the parameters (π, A, B) .

All previous inferences depend on knowing the parameters (π, A, B) .

How do we learn the parameters based on N observation sequences $x_{n,1}, \ldots, x_{n,T}$ for $n = 1, \ldots, N$?

All previous inferences depend on knowing the parameters (π, A, B) .

How do we learn the parameters based on N observation sequences $x_{n,1}, \ldots, x_{n,T}$ for $n = 1, \ldots, N$?

MLE is intractable due to the hidden variables $Z_{n,t}$'s (similar to GMMs)

All previous inferences depend on knowing the parameters (π, A, B) .

How do we learn the parameters based on N observation sequences $x_{n,1}, \ldots, x_{n,T}$ for $n = 1, \ldots, N$?

MLE is intractable due to the hidden variables $Z_{n,t}$'s (similar to GMMs)

Need to apply EM again! Known as the Baum-Welch algorithm.

Recall in the E-Step we fix the parameters and find the **posterior** distributions q of the hidden states (for each sample n),

Recall in the E-Step we fix the parameters and find the **posterior** distributions q of the hidden states (for each sample n), which leads to the complete log-likelihood:

$$\mathbb{E}_{z_{1:T} \sim q} \left[\ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \right]$$

Recall in the E-Step we fix the parameters and find the **posterior** distributions q of the hidden states (for each sample n), which leads to the complete log-likelihood:

$$\mathbb{E}_{z_{1:T} \sim q} \left[\ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \right]$$

$$= \mathbb{E}_{z_{1:T} \sim q} \left[\ln \pi_{z_1} + \sum_{t=1}^{T-1} \ln a_{z_t, z_{t+1}} + \sum_{t=1}^{T} \ln b_{z_t, x_t} \right]$$

Recall in the E-Step we fix the parameters and find the **posterior** distributions q of the hidden states (for each sample n), which leads to the complete log-likelihood:

$$\begin{split} &\mathbb{E}_{z_{1:T} \sim q} \left[\ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \right] \\ &= \mathbb{E}_{z_{1:T} \sim q} \left[\ln \pi_{z_1} + \sum_{t=1}^{T-1} \ln a_{z_t, z_{t+1}} + \sum_{t=1}^{T} \ln b_{z_t, x_t} \right] \\ &= \sum_{s} \gamma_s(1) \ln \pi_s + \sum_{t=1}^{T-1} \sum_{s, s'} \xi_{s, s'}(t) \ln a_{s, s'} + \sum_{t=1}^{T} \sum_{s} \gamma_s(t) \ln b_{s, x_t} \end{split}$$

We have discussed how to compute

$$\gamma_s(t) = P(Z_t = s \mid X_{1:T} = x_{1:T})$$

$$\xi_{s,s'}(t) = P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T})$$

The maximizer of complete log-likelihood is simply doing **weighted counting** (compared to the unweighted counting on Slide 36):

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1) = \mathbb{E}_q \left[\text{ \#initial states with value } s \right]$$

$$a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t) = \mathbb{E}_q \left[\text{ \#transitions from } s \text{ to } s' \right]$$

$$b_{s,o} \propto \sum_n \sum_{t: x_t = o} \gamma_s^{(n)}(t) = \mathbb{E}_q \left[\text{ \#state-outcome pairs } (s,o) \right]$$

The maximizer of complete log-likelihood is simply doing **weighted counting** (compared to the unweighted counting on Slide 36):

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1) = \mathbb{E}_q \left[\text{ \#initial states with value } s \right]$$

$$a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t) = \mathbb{E}_q \left[\text{ \#transitions from } s \text{ to } s' \right]$$

$$b_{s,o} \propto \sum_n \sum_{t: x_t = o} \gamma_s^{(n)}(t) = \mathbb{E}_q \left[\text{ \#state-outcome pairs } (s,o) \right]$$

where

$$\begin{split} \gamma_s^{(n)}(t) &= P(Z_{n,t} = s \mid X_{n,1:T} = x_{n,1:T}) \\ \xi_{s,s'}^{(n)}(t) &= P(Z_{n,t} = s, Z_{n,t+1} = s' \mid X_{n,1:T} = x_{n,1:T}) \end{split}$$

The maximizer of complete log-likelihood is simply doing **weighted counting** (compared to the unweighted counting on Slide 36):

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1) = \mathbb{E}_q \left[\text{ \#initial states with value } s \right]$$

$$a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t) = \mathbb{E}_q \left[\text{ \#transitions from } s \text{ to } s' \right]$$

$$b_{s,o} \propto \sum_n \sum_{t:x_t=o} \gamma_s^{(n)}(t) = \mathbb{E}_q \left[\text{ \#state-outcome pairs } (s,o) \right]$$

where

$$\gamma_s^{(n)}(t) = P(Z_{n,t} = s \mid X_{n,1:T} = x_{n,1:T})$$

$$\xi_{s,s'}^{(n)}(t) = P(Z_{n,t} = s, Z_{n,t+1} = s' \mid X_{n,1:T} = x_{n,1:T})$$

(Recall how EM for GMM updates parameters using weighted averages.)

Step 0 Initialize the parameters (π, A, B)

Step 0 Initialize the parameters $(m{\pi}, m{A}, m{B})$

Step 1 (E-Step) Fixing the parameters, compute forward and backward messages for all sample sequences, then use these to compute $\gamma_s^{(n)}(t)$ and $\xi_{s\,s'}^{(n)}(t)$ for each n,t,s,s' (see Slides 47 and 48).

Step 0 Initialize the parameters $(m{\pi}, m{A}, m{B})$

Step 1 (E-Step) Fixing the parameters, compute forward and backward messages for all sample sequences, then use these to compute $\gamma_s^{(n)}(t)$ and $\xi_{s,s'}^{(n)}(t)$ for each n,t,s,s' (see Slides 47 and 48).

Step 2 (M-Step) Update parameters:

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1), \quad a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t), \quad b_{s,o} \propto \sum_n \sum_{t:x_t=o} \gamma_s^{(n)}(t)$$

Step 0 Initialize the parameters $(m{\pi}, m{A}, m{B})$

Step 1 (E-Step) Fixing the parameters, compute forward and backward messages for all sample sequences, then use these to compute $\gamma_s^{(n)}(t)$ and $\xi_{s,s'}^{(n)}(t)$ for each n,t,s,s' (see Slides 47 and 48).

Step 2 (M-Step) Update parameters:

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1), \quad a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t), \quad b_{s,o} \propto \sum_n \sum_{t:x_t=o} \gamma_s^{(n)}(t)$$

Step 3 Return to Step 1 if not converged

Summary

Very important models: Markov chains, hidden Markov models

Summary

Very important models: Markov chains, hidden Markov models

Several algorithms:

- forward and backward procedures
- inferring HMMs based on forward and backward messages
- Viterbi algorithm and variants (see today's discussion)
- Baum–Welch algorithm