

Assignment 4

Clustering Languages

Verena Blaschke

July 04, 2018

Assignment 4

I: Feature extraction

II: K-means clustering

II: K-means clustering

III: Principal component analysis

IV: Evaluation with gold-standard labels

V: Calculating distances

VI: Hierarchical clustering

I: Feature extraction

fin s i l m æ

fin k ɔ r v a

fin n ɛ n æ

fin s u u

...

cmn t^h ʊ ŋ t ʃ i

cmn ʤ ə n n a ɪ

I: Feature extraction

fin s i l m æ

fin k ɔ r v a

fin n ɛ n æ

fin s u u

...

cmn t^h ʊ ŋ t ʂ i

cmn ʈ ə n n a ɪ

- ▶ 80 languages × 272 IPA segments

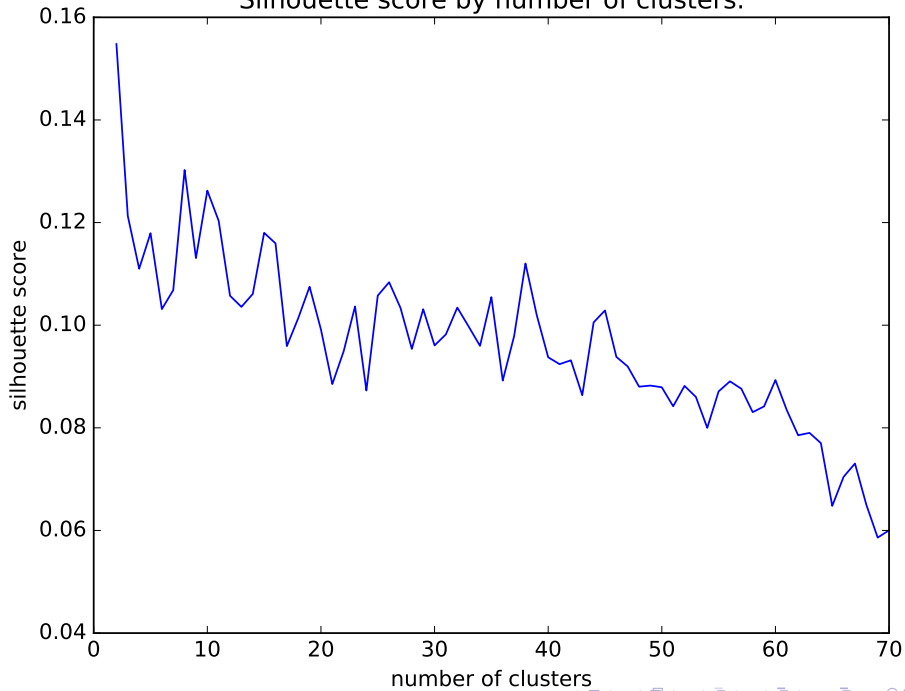
II: K-means clustering

- ▶ k-means clustering for k in $[2, 70]$
- ▶ silhouette coefficient

II: K-means clustering

- ▶ k-means clustering for k in $[2, 70]$
- ▶ silhouette coefficient
 - ▶ How close (=similar) is each data point to other points from its own cluster compared to other clusters?
 - ▶ $[-1, +1]$, higher scores are better

Silhouette score by number of clusters.



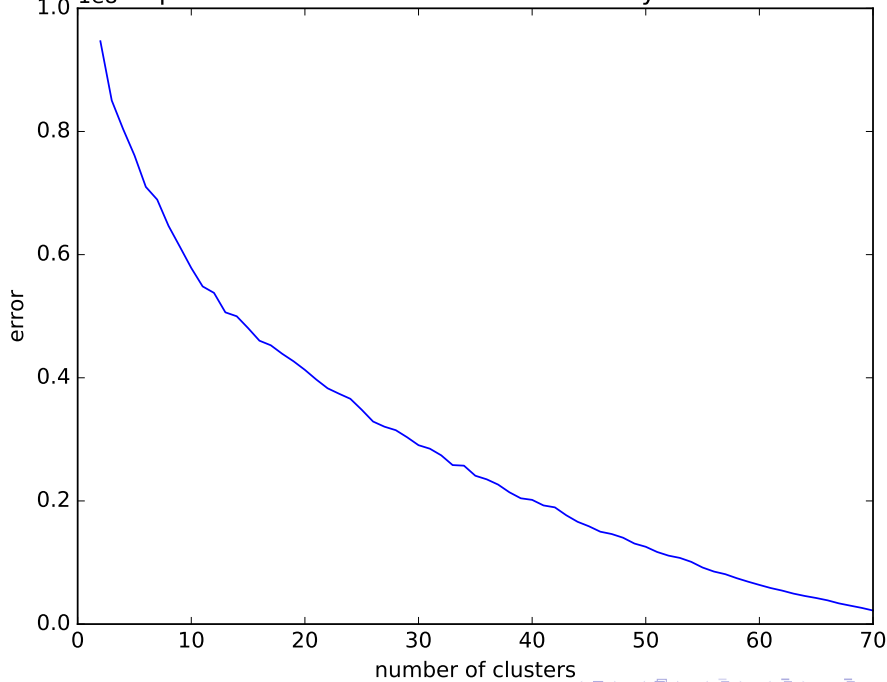
II: K-means clustering

- ▶ k-means clustering for k in $[2, 70]$
- ▶ silhouette coefficient
 - ▶ How close (=similar) is each data point to other points from its own cluster compared to other clusters?
 - ▶ $[-1, +1]$, higher scores are better
- ▶ error function
 - ▶ sum of squared distances from the closest centroids
`kmeans.inertia_`

II: K-means clustering

- ▶ k-means clustering for k in $[2, 70]$
- ▶ silhouette coefficient
 - ▶ How close (=similar) is each data point to other points from its own cluster compared to other clusters?
 - ▶ $[-1, +1]$, higher scores are better
- ▶ error function
 - ▶ sum of squared distances from the closest centroids
`kmeans.inertia_`
 - ▶ What is a good number of clusters? → elbow method

Sum of squared distances from the centroids by number of clusters.



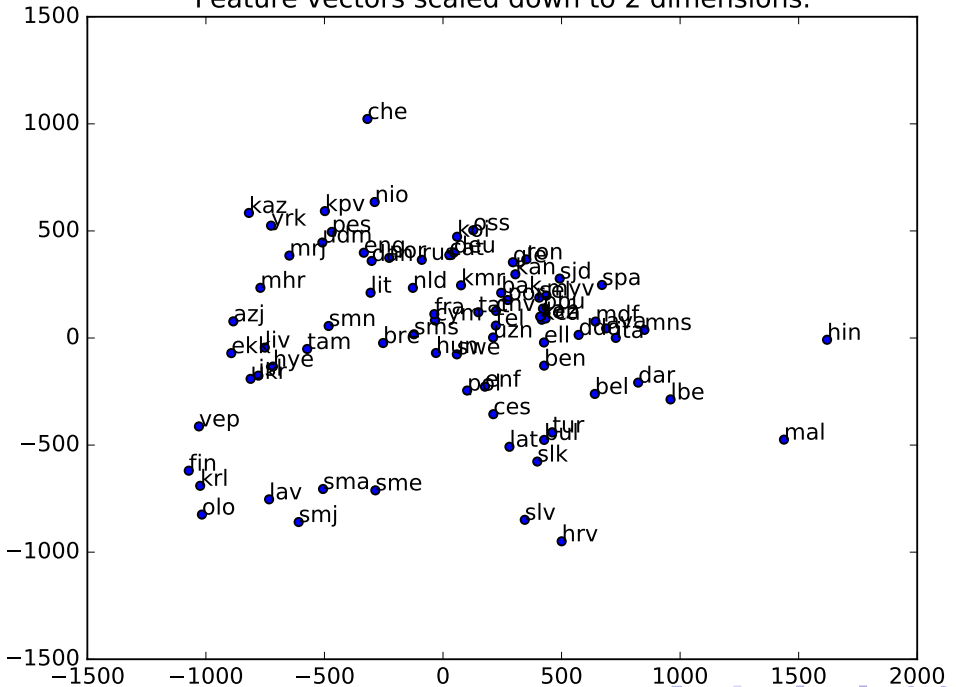
III: Principal component analysis

- ▶ remove redundant features

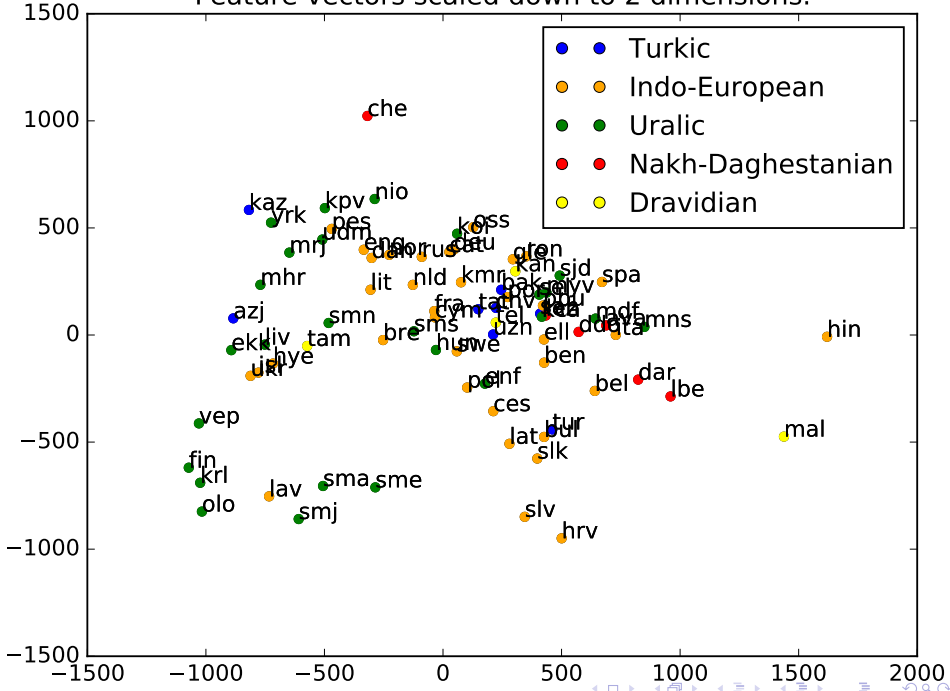
III: Principal component analysis

- ▶ remove redundant features
 - ▶ remove noise
 - ▶ train machine learning models more quickly

Feature vectors scaled down to 2 dimensions.



Feature vectors scaled down to 2 dimensions.



III: Principal component analysis

```
pca = PCA(features.shape[1])
d = 0
var_explained = 0
while var_explained < 0.9:
    var_explained += pca.explained_variance_ratio_[d]
    d += 1

featuresPCA = PCA(d).fit_transform(features)
```

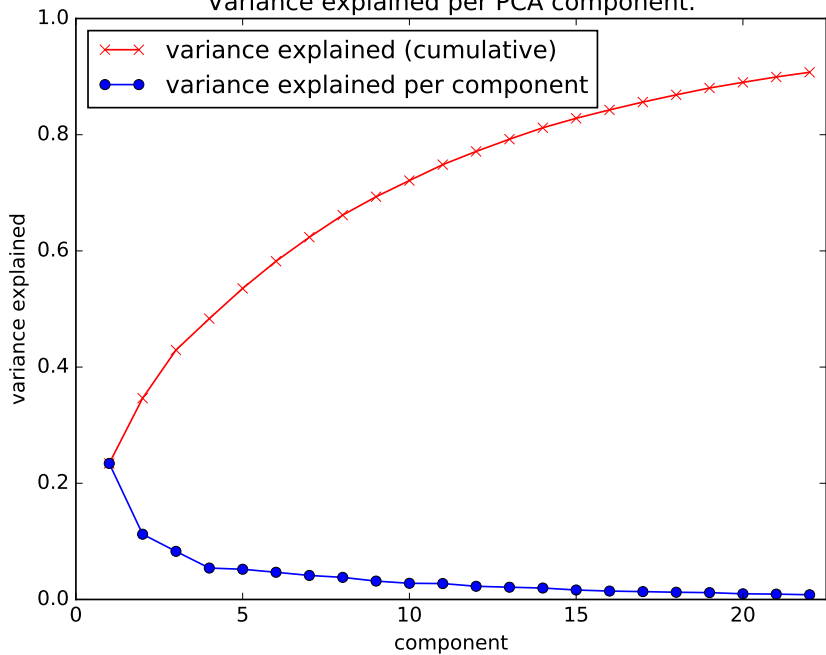
III: Principal component analysis

```
pca = PCA(features.shape[1])
d = 0
var_explained = 0
while var_explained < 0.9:
    var_explained += pca.explained_variance_ratio_[d]
    d += 1
```

```
featuresPCA = PCA(d).fit_transform(features)
```

```
pca = PCA(0.9)
print(pca.n_components_)
```

Variance explained per PCA component.



IV: Evaluation with gold-standard labels

```
n_fam = len(set(family))  
pred_all = KMeans(n_fam).fit_predict(features)  
pred_pca = KMeans(n_fam).fit_predict(featuresPCA)
```

IV: Evaluation with gold-standard labels

```
n_fam = len(set(family))  
pred_all = KMeans(n_fam).fit_predict(features)  
pred_pca = KMeans(n_fam).fit_predict(featuresPCA)
```

lang	all	pca	family
kan	2	4	Dravidian
tam	3	0	Dravidian
tel	4	0	Dravidian
mal	4	2	Dravidian
bul	0	1	Indo-European
ces	0	1	Indo-European
...			

IV: Evaluation with gold-standard labels

- ▶ **Homogeneity:** Each cluster contains data points of the same gold-standard class.

IV: Evaluation with gold-standard labels

- ▶ **Homogeneity:** Each cluster contains data points of the same gold-standard class.
- ▶ **Completeness:** All members of a gold-standard class are in the same cluster.

IV: Evaluation with gold-standard labels

- ▶ **Homogeneity:** Each cluster contains data points of the same gold-standard class.
- ▶ **Completeness:** All members of a gold-standard class are in the same cluster.
- ▶ **V-measure:** Harmonic mean of homogeneity and completeness.
- ▶ $[0,1]$ higher is better

IV: Evaluation with gold-standard labels

- ▶ **Homogeneity:** Each cluster contains data points of the same gold-standard class.
- ▶ **Completeness:** All members of a gold-standard class are in the same cluster.
- ▶ **V-measure:** Harmonic mean of homogeneity and completeness.
- ▶ $[0,1]$ higher is better

all H: 0.1707 C: 0.1461 V: 0.1575

PCA H: 0.1728 C: 0.1572 V: 0.1646

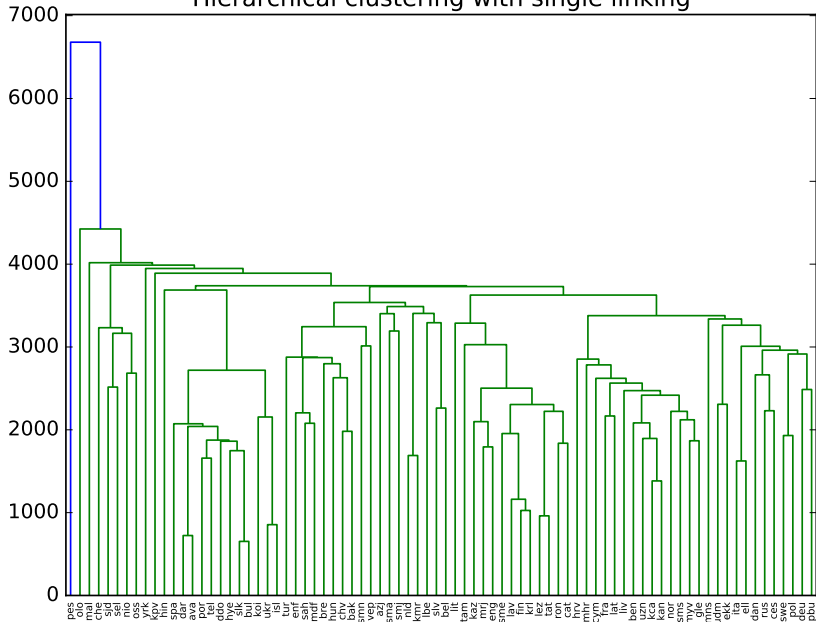
V: Calculating distances

A	B	C	D	E	
	123	452	10	572	A
		342	370	908	B
			127	754	C
				23	D
					E

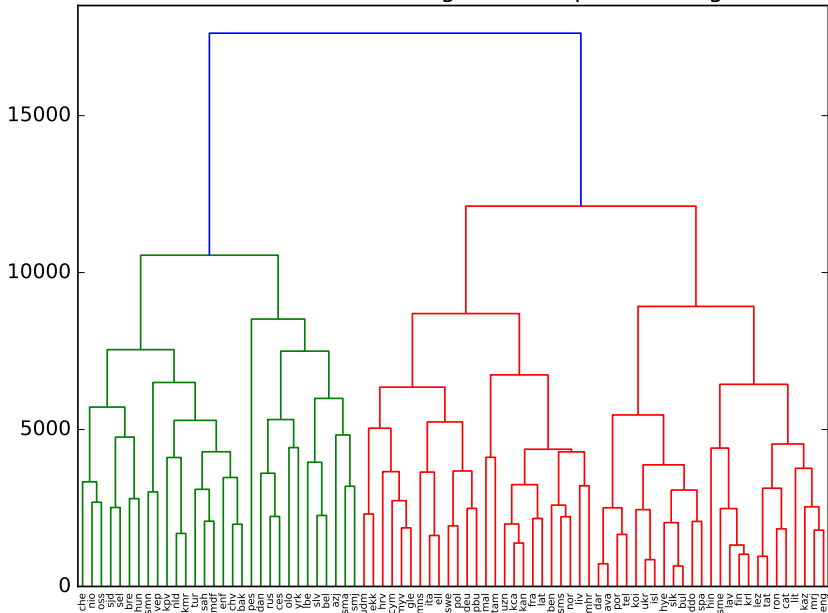
VI: Hierarchical clustering

```
for m in ['single', 'complete', 'average']:
    fig, ax = plt.subplots()
    z = scipy.cluster.hierarchy.linkage(dist, method=m)
    scipy.cluster.hierarchy.dendrogram(z, labels=languages)
    fig.savefig('dendrogram-{}.pdf'.format(method))
```

Hierarchical clustering with single linking



Hierarchical clustering with complete linking



Hierarchical clustering with average linking

