

# Class-based Identification of 'Deviant' Semantic Features in Historical Corpora

Kristoffer L. Nielbo

kln@cas.au.dk

DTL|IMC|AU-CAS  
Aarhus University, Denmark



# statistical learning

## **goal**

build a machine that can learn from data and automatically make the right decisions

## **supervised**

infer mapping between data & class-information  $\rightarrow$  theoretical 'ground truth'

## **unsupervised**

identify latent classes in the data  $\rightarrow$  lack theoretical 'ground truth'

# application to ctext corpus

## goal

combine statistical learning and information theory in order to explore semantic relations between data (*Shangshu*) and theoretical classes (document dating)

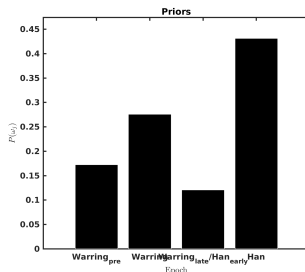
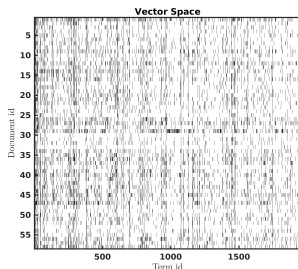
## naïve bayes

simple and well-performing Bayesian model for supervised learning, but *too constrained*

## latent dirichlet allocation

simple and popular Bayesian model for unsupervised learning, but *too unconstrained*

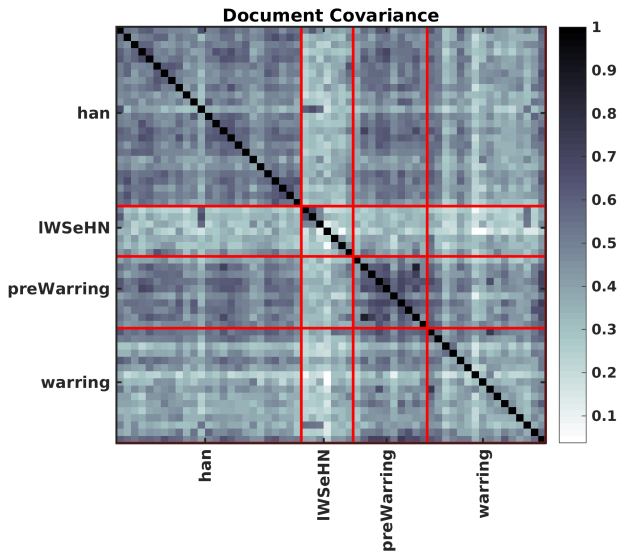
supervised/nb

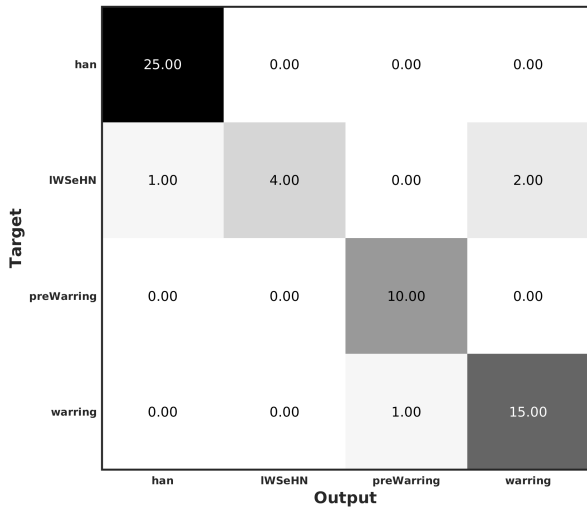


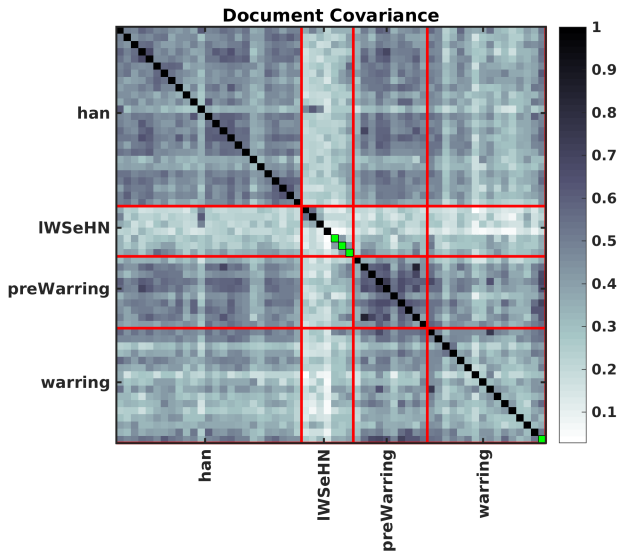
$\mathbf{d}_i$  is the  $2K$ -dimensional feature vectors,  $d_{i1} = P(c_1)$ ,  $d_{i2} = P(c_2)$ , ...  
 $d_{i1958} = P(c_{1958})$

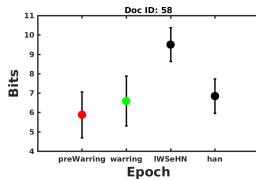
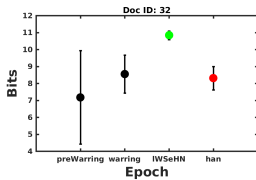
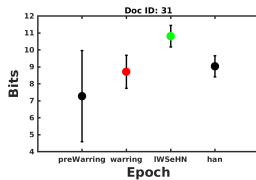
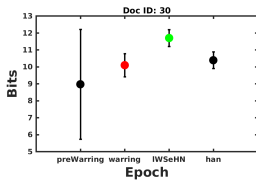
$w_j \in \{Warring_{pre}, Warring, Warring_{late}/Han_{early}, Han\}$

explore category boundaries  $\sim$  **misclassification semantics**

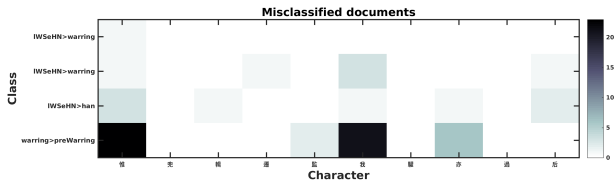
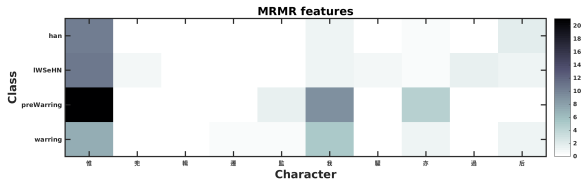


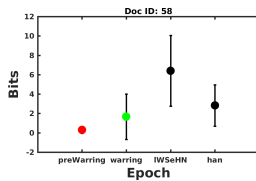
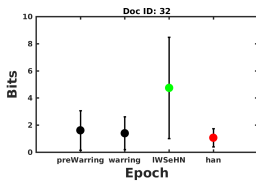
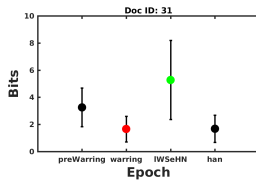
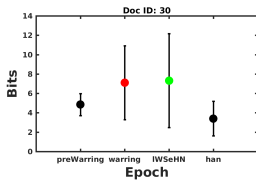




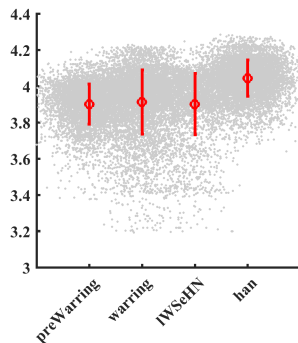
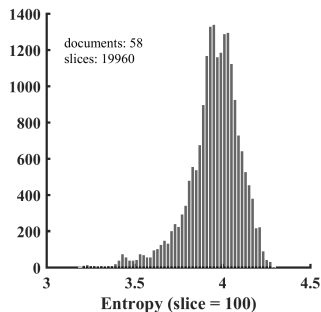








# unsupervised/lda

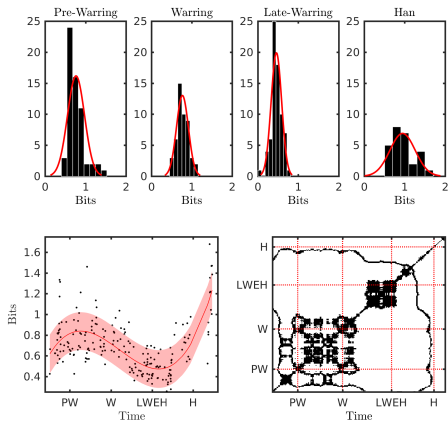


“a rose is a rose is a rose” is less lexically dense than “a rose is red and thorny”

lexical density  $\sim$  text predictability:  $H(X) = \sum_{i=1}^n p_i \log_2 p_i$

$$H(\text{a rose is a rose is a rose}) < H(\text{a rose is red and thorny})$$

$$H(\text{a rose is a rose is a rose}) = H(\text{erea oiasessar oiors})$$



Ida to the rescue:  $\theta_i$  probability distribution of  $k$  latent variables in document  $i$

disruption between document is the relative entropy:  $D_{KL}(P \parallel Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$

## in summary

### **goal**

creative use of statistical learning to support humanistic inquiry

### **supervised**

study semantics on the boundaries of theoretical classes relying on less constrained human interpretation

### **unsupervised**

use theoretical classes to study semantic evolution of cultural system without unconstrained human interpretation