

# Statistical Learning III

15-11-2019

1. Probability Theory Remarks
2. Bayesian Statistics
  - a) Intro
  - b) Laplace's Rule of Succession
  - c) MAP Interpretation of Ridge, Lasso
3. Plug-in Estimators
4. Naive Bayes

### 3. Probability Theory Remarks

Bayes' rule:

$$\Pr(A|B) = \frac{\Pr(B|A) \cdot \Pr(A)}{\Pr(B)}$$

Often used in Bayesian ~~and~~ standard frequentist statistics

Densities are slippery:

- \* Location of maximum depends on choice of parametrization
- \* Uniform density in one parametrization is not uniform in another parametrization

(see slides)

(2)

## 2. Bayesian Statistics

### a) Intro

Probability model:  $\mathcal{P} = \{P_\theta(x, y) | \theta \in \Theta\}$  or  $\mathcal{P} = \{P_\theta(y) | \theta \in \Theta\}$

Frequentist statistics (standard):

- optimal parameter  $\theta^*$  is fixed, but unknown
- diversity of methods
- need proofs/experiments to justify methods

Bayesian statistics:

- pretend that true parameter  $\theta^*$  is a random variable distributed according to prior distribution  $\pi(\theta)$  that we know.

$T = Y_1, \dots, Y_N$  (no features for simplicity)

$\Pr(T, \theta) = P_\theta(T) \cdot \pi(\theta)$  is joint distribution of data and parameters

Since we know the full joint distribution, can simply use probability theory to compute any probability we are interested in. ← single method

If we estimate different probabilities this way, they are beautifully consistent.

~~But~~ Bayesian premise is too strong:

- prior  $\pi$  is chosen for computational or information theoretic properties in practice, so cannot ~~blindly~~ assume  $\theta^*$  is random sample from  $\pi$ .
- need proofs/experiments to justify Bayesian methods

(3)

frequentist

## Modern motivation:

- If we choose  $\pi$  right, then often works really well (both in theory and in practice).
- Learns faster if true  $\theta^*$  has high prior probability, slower if true  $\theta^*$  has small prior probability,  
so can use  $\pi$  to express prior knowledge about our data.

## Posterior Distribution:

$$\pi(\theta|T) := \Pr(\theta|T) = \frac{\Pr(\theta, T)}{\Pr(T)} = \frac{P_\theta(T) \cdot \pi(\theta)}{\Pr(T)}$$

- Often puts its probability mass closer and closer to  $\theta^*$  as  $N \rightarrow \infty$ . (vd Vaart et al.)
- Expresses uncertainty about  $\theta$

## Predictive Distribution:

$$\Pr(Y|T) = \int_0^\infty P_\theta(Y) \cdot \pi(\theta|T) d\theta = \frac{\Pr(Y, T)}{\Pr(T)}$$

↑                      ↑  
 new sample            often better predictions  
 outside of            than frequentist  $P_{\hat{\theta}}(Y)$   
 training set

because  $\pi(\theta|T)$  keeps track of uncertainty better than fixed single choice  $\hat{\theta}$ .

## b] Example: Laplace Rule of Succession

Bernoulli model:  $P_\theta(Y) = \begin{cases} \theta & \text{for } Y=1 \\ 1-\theta & \text{for } Y=0 \end{cases} \quad \theta \in [0, 1]$

Maximum likelihood:  $\hat{\theta} = \frac{n_1}{N} \rightarrow P_{\hat{\theta}}(Y=1) = \hat{\theta} = \frac{n_1}{N}$  ← dangerous for prediction if  $n_1=0$

Suppose  $\pi(\theta)=1$  is uniform prior density, ← not the same as "no prior knowledge" because depends on parametrisation

Then  $\Pr(Y=1|T) = \frac{n_1 + 1}{N + 2}$

(Laplace, 1814)

$$\Pr(Y=0|T) = \frac{n_0 + 1}{N + 2}$$

↙ beta function

$$\text{Proof: } \frac{\Pr(Y, T)}{\Pr(T)} = \frac{\int P_\theta(Y, T) \cdot \pi(\theta) d\theta}{\int P_\theta(T) \cdot \pi(\theta) d\theta} = \frac{\int \theta^{n_1+Y} (1-\theta)^{n_0+2-Y} d\theta}{\int \theta^{n_1} (1-\theta)^{n_0} d\theta} = \frac{n_1 + 1}{N + 2} \blacksquare$$

(4)

### c) MAP interpretation of Ridge and Lasso

Maximum a Posteriori (MAP) parameters:

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \pi(\theta | T) = \underset{\theta}{\operatorname{argmax}} \frac{p_{\theta}(T) \cdot \pi(\theta)}{p_T(T)}$$

$$= \underset{\theta}{\operatorname{argmax}} p_{\theta}(T) \cdot \pi(\theta)$$

- maximizes posterior density, so for continuous parameters depends on parametrisation
- "real" Bayesians prefer prediction with predictive distribution

Ridge / Lasso:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \lambda \text{pen}(\beta)$$

$$\text{ridge: pen}(\beta) = \sum_{j=2}^p \beta_j^2 \quad \text{lasso: pen}(\beta) = \sum_{j=2}^p |\beta_j|$$

Suppose Gaussian noise:

$$y = x^T \beta + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

$$p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^T \beta)^2}{2\sigma^2}}$$

$$\hat{\beta}_{\text{MAP}} = \underset{\beta}{\operatorname{argmax}} p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) \cdot \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} -\log p_{\beta}(y_1, \dots, y_N | x_1, \dots, x_N) - \log \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} N \cdot \left( -\log \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{1}{2\sigma^2} \text{RSS}(\beta) - \log \pi(\beta)$$

$$= \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) - 2\sigma^2 \log \pi(\beta)$$

$$\begin{aligned} \log(x^a) &= a \log(x) \\ \log(ab) &= \log(a) + \log(b) \end{aligned}$$

Suppose data have been pre-processed such that we can assume that the intercept  $\hat{\beta}_0 = 0$ .

(5)

Ridge: Choose  $\pi$  s.t.

$$\beta_1 = 0 \text{ with prob. 1}$$

$$(\beta_2, \dots, \beta_p) \sim N(0, \sigma^2 I)$$

$$-\log \pi(\beta) = \sum_{j=2}^p -\log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(\beta_j - 0)^2}{2\sigma^2}} \right)$$

$$= (p-1)(-\log \frac{1}{\sqrt{2\pi\sigma^2}}) + \sum_{j=2}^p \frac{\beta_j^2}{2\sigma^2}$$

$$\hat{\beta}_{MAP} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \frac{\sigma^2}{2\sigma^2} \sum_{j=2}^p \beta_j^2$$

is ridge with  $\lambda = \frac{\sigma^2}{2\sigma^2}$

Is also posterior mean  $E[\beta]_{\pi(\beta|T)}$ .

Lasso: Choose  $\pi$  s.t.

$$\beta_1 = 0 \text{ with prob. 1.}$$

$$(\beta_2, \dots, \beta_p) \sim \pi \underset{j=2}{\overset{p}{\prod}} \frac{1}{2\sigma} e^{-\frac{|\beta_j|}{\sigma}}$$

$$-\log \pi(\beta) = \sum_{j=2}^p -\log \left( \frac{1}{2\sigma} \cdot e^{-\frac{|\beta_j|}{\sigma}} \right)$$

$$= (p-1)(-\log \frac{1}{2\sigma}) + \sum_{j=2}^p \frac{|\beta_j|}{\sigma}$$

$$\hat{\beta}_{MAP} = \underset{\beta}{\operatorname{argmin}} \text{RSS}(\beta) + \frac{2\sigma^2}{\sigma} \sum_{j=2}^p |\beta_j|$$

is Lasso with  $\lambda = \frac{2\sigma^2}{\sigma}$ .

Remarks:

- "real" Bayesian prefers predicting with predictive distribution
- MAP + CV to determine  $\lambda$  is very unBayesian.

## 1. Plug-in Estimators

Bayes' rule gives :  $P^*(y=k|x) = \frac{P^*(x|y=k) P^*(y=k)}{P^*(x)}$

Notation :

$$f_k(x) := P^*(x|y=k)$$

$$\pi_k := P^*(y=k)$$

Plug-in estimators :

- Estimate  $P^*(y=k|x)$  by  $\hat{P}(y=k|x)$

and predict with  $\underset{k}{\operatorname{argmax}} \hat{P}(y=k|x)$

- By Bayes' rule, it is sufficient to estimate  $f_k$  and  $\pi_k$  for all classes  $k$ :

$$\hat{f}(x) = \underset{k}{\operatorname{argmax}} \frac{\hat{f}_k(x) \hat{\pi}_k}{P^*(x)} = \underset{k}{\operatorname{argmax}} \hat{f}_k(x) \cdot \hat{\pi}_k$$

N.B. Although we have used Bayes' rule, there is nothing inherently Bayesian about this approach.  
(Can estimate  $\hat{f}_k$  by Bayesian or frequentist methods.)

### 3. Naive Bayes

Model: features are independent given their class  $k$

Remark: This implies that  $\Sigma_k = \begin{pmatrix} \sigma_{k,11}^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_{k,p}^2 \end{pmatrix}$

is a diagonal matrix with  $p$  parameters, because (by assumption) all covariances are  $0$ .

For each class  $k$  and feature vector  $x_i$ :

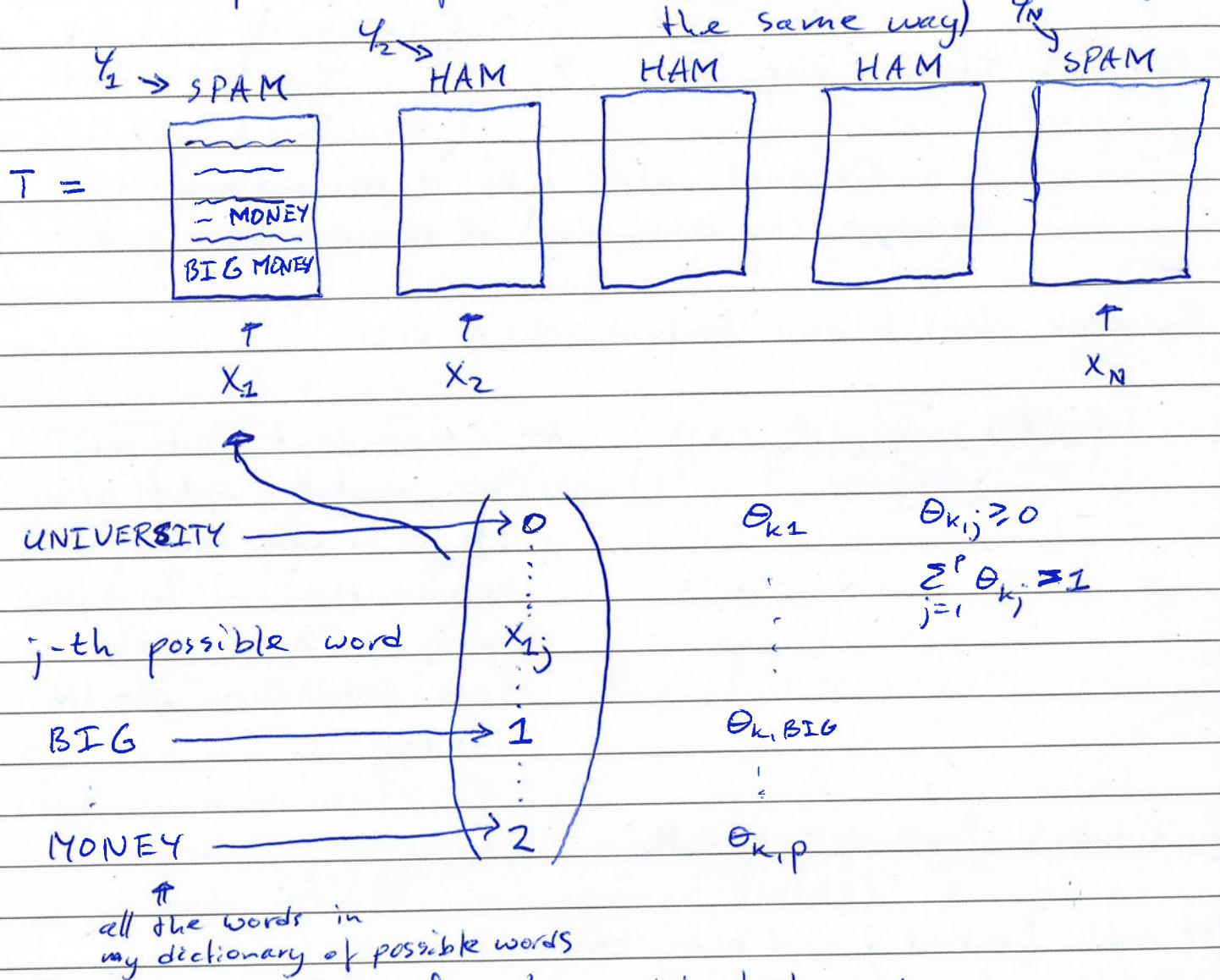
$$f_k(x_i) = \prod_{j=1}^p f_{kj}(x_{ij}) \quad \begin{matrix} \text{* independent features;} \\ \text{very useful simplification} \\ \text{if } p \text{ is very large,} \\ \text{because then we have no} \\ \text{hope of learning the} \\ \text{true distribution exactly.} \\ \text{anyway.} \end{matrix}$$

Continuous features:  $x_{ij} \in \mathbb{R}$

For each feature  $j$  and every class  $k_0$  estimate  $f_{kj}$  by fitting a univariate Gaussian (similar to LDA).

## Discrete features:

Consider spam classification. (other text classification goes the same way)



- ① Forget position of each word in text, and represent each e-mail by counts of how many times words occur.

- ② Use separate multinomial model per class  $k$ :

$$f_k(x_i) = \prod_{j=1}^p \theta_{k,j}^{n_{ij}}, \text{ where } n_{ij} \text{ is the number of times the } j\text{-th word occurs in the } i\text{-th e-mail.}$$

$$\frac{1}{n_k} = \frac{\# \text{e-mails in class } k}{N}$$

$$\hat{\theta}_{k,\text{MONEY}} = \frac{\#\text{MONEY in e-mails from class } k + \text{small num}}{\#\text{words in all e-mails from class } k + p \cdot \text{small num}}$$

Without adding the "small numbers", if a spam e-mail contains one word that we have never seen in a spam e-mail before, then  $f_{\text{SPAM}}(x_i) = 0$ , and the e-mail will be classified as HAM. So a single word can throw off classification.

### Too simple/naive?

- Only need  $\hat{P}(y=k_1|X) > \hat{P}(y=k_2|X)$  whenever  $P^*(y=k_1|X) > P^*(y=k_2|X)$
- Possible even if  $\hat{P}$  is poor estimate of  $P^*$ !