

بسم الله الرحمن الرحيم

الدكتور هيتكم انهارده علي naïve Bayes و هنتكلم علي الفرق ما بين ال Generative models ... حاجة زي naïve bayes ... و ال Discriminative models .. حاجة زي ال logistic regression .. في سؤال .. ليه بنستخدم Naïve bayes .. الاجابه ان ليه ميزة وهو ان ال training ببقا سريع جدا ... اللي تعملو بس هو انك تحسب ال parameters in closed form عن طريق انك تعد بس .. في بقيت الطرق انتحتاج وقت كبير عشان تعمل ال training و بتفضل تروح وتعيد في ال epochs وكدا .. طيب يلا بيينا

Naïve bayes model ده بيوصف Generative model .. بمعنى .. كلمة generative هنا معناها هو وصف لإزاي احنا نقدر generate the data we are working wiith .. قيل كدا انت مكتش عندك الحوار ده ... انت كان آخرك في ال discriminative model هو ان الإكسات بتقا عندك .. وبعدين ال y بقا هنا كان ممكن تقول ال y ديه generative . جات منين يعني بس ال y بتتقا on x هنا انت عندك في ال generative ازاى ال x and y بيبقى generated .. ففي الصورة .. انت بت flip the coin .. خد القيمة بتاعته اللي هي هتتقا ال y .. وبعدين روج اعمل flipping to the left/right coins .. علي حسب قيمة الواي لو 1 روج اقلب اللي لونهم أزرق .. لو صفر روج اقلب اللي لونهم احمر ... الفكرة ان في Naïve Bayes assumption بتتقا independent ... اللي عاوزين تعملو ناو اننا نعمل ملخص كدا سريع لل closed form MLE applied to NVBayes model

generation of y values

② Log-likelihood =

$$\ell(\phi, \theta) = \log \prod_{i=1}^N p(x^{(i)}, y^{(i)} | \phi, \theta)$$

$$= \sum_{i=1}^N \left[\log p(y^{(i)} | \phi) + \sum_{m=1}^M \log p(x_m^{(i)} | y^{(i)}, \theta) \right]$$

$N_{y=H} = \#(y^{(i)} = H)$
 $N_{y=T} = \#(y^{(i)} = T)$

$N_{x_m=1, y=H} = \# \text{ times that } x_m^{(i)} = 1 \text{ and } y^{(i)} = H$

$N_{y=H} \log \phi + N_{y=T} \log(1-\phi)$
 $+ \sum_{m=1}^M N_{x_m=1, y=H} \log(\theta_{H,m}) + N_{x_m=0, y=H} \log(1-\theta_{H,m})$
 $+ \sum_{m=1}^M N_{x_m=1, y=T} \log(\theta_{T,m}) + N_{x_m=0, y=T} \log(1-\theta_{T,m})$

M features
 $x_m^{(i)} = 1$
 $x_m^{(i)} = 0$
 $y^{(i)} = H$
 $y^{(i)} = T$

generation of x values

Case a) ϕ b) θ elements

③a) Take partial derivatives w.r.t. ϕ

$$\frac{\partial \ell(\phi, \theta)}{\partial \phi} = \frac{N_{y=1}}{\phi} - \frac{N_{y=0}}{1-\phi}$$

big gold coin 

④a) Set to zero and solve for ϕ

$$\phi^{MLE} = \frac{N_{y=1}}{N_{y=1} + N_{y=0}} = \frac{N_{y=1}}{N}$$

No. of times it was heads

(3b) Take partial derivative wrt $\Theta_{H,m}$ (case for $y=H$ and feature m)

$$\frac{\partial \ell(\phi, \theta)}{\partial \Theta_{H,m}} = \frac{N_{x_m=1, y=H}}{\Theta_{H,m}} - \frac{N_{x_m=0, y=H}}{1 - \Theta_{H,m}}$$

for red coins, $\frac{H}{I}$
and
blue coins, $\frac{T}{I}$

(4b) Solve for $\Theta_{H,m}$

$$\Theta_{H,m}^{MLE} = \frac{N_{x_m=1, y=H}}{N_{x_m=1, y=H} + N_{x_m=0, y=H}} = \frac{N_{x_m=1, y=H}}{N_{y=H}}$$

weights of coins by counting

طبيب دلوقت احنا معانا ال Naïve bayes max likelihood estimate .. احنا هنفكر شويه في تفاصيل الموديل بتاع Bernoulli Naïve bayes

Model 1: Bernoulli Naïve Bayes

Data: Binary feature vectors, Binary labels

$\mathbf{x} \in \{0, 1\}^M$ $y \in \{0, 1\}$

Generative Story:

$y \sim \text{Bernoulli}(\phi)$

$x_1 \sim \text{Bernoulli}(\theta_{y,1})$

$x_2 \sim \text{Bernoulli}(\theta_{y,2})$

\vdots

$x_M \sim \text{Bernoulli}(\theta_{y,M})$

Model:

$p_{\phi, \theta}(\mathbf{x}, y) = p_{\phi, \theta}(x_1, \dots, x_M, y)$

$= p_{\phi}(y) \prod_{m=1}^M p_{\theta}(x_m | y)$

$= \left[(\phi)^y (1 - \phi)^{(1-y)} \right]$

$\prod_{m=1}^M (\theta_{y,m})^{x_m} (1 - \theta_{y,m})^{(1-x_m)}$

\mathbf{x} is binary feature vector
 y is binary label

Model 1: Bernoulli Naïve Bayes

Maximum Likelihood Estimation

عشان فصل

Training: Find the class-conditional MLE parameters

Count Variables:

$$N_{y=1} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0, x_m=1} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 0 \wedge x_m^{(i)} = 1)$$

...

Maximum Likelihood Estimators:

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = \frac{N_{y=0, x_m=1}}{N_{y=0}}$$

$$\theta_{1,m} = \frac{N_{y=1, x_m=1}}{N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$

algebraic Representation of How we can define

17

A Shortcoming of MLE

For Naïve Bayes, suppose we **never** observe the word “unicorn” in a **real** news article.

In this case, what is the MLE of the following quantity?

$$p(x_{\text{unicorn}} = 1 | y = \text{real}) = 0$$

y=Real & x_unicorn=1

Recall: $\theta_{k,0} = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 0 \wedge x_k^{(i)} = 1)}{\sum_{i=1}^N \mathbb{I}(y^{(i)} = 0)}$

total # of y, Real

Now suppose we observe the word “unicorn” at test time. What is the posterior probability that the article was a **real** article?

$$p(y = \text{real} | \mathbf{x}) = \frac{p(\mathbf{x} | y = \text{real}) p(y = \text{real})}{p(\mathbf{x})}$$

$$= \frac{\prod_{m=1}^M p(x_m | y = \text{real}) p(y = \text{real})}{p(\mathbf{x})}$$

هنا posterior prob
احتمال ان الكلمة في اي فئة

Posterior Prob $\rightarrow y = \text{Real} | \mathbf{x}$

23

الدكتور في السلايد ديه بيقول بعد ما لقي ان ال posterior هتبقا بصفر عشان احنا بس ملاحظناش الكلمة اللي اسمها unicorn قبل كدا في أي حجة Real .. فال posterior بقا بصفر .. فال document مبقتش real .. في الحالة ديه في حجة مش منطقيه .. انك تقول ان document isn't real عشان مجرد ان في كلمة نادرأ ما كانت موجودة في كل ال articles اللي شغلناها في التريننج تايم ... بيقا كان احسن ان ال $p(x|y=\text{real})$ تبقا قيمتها صغيره مش تبقي بصفر علي طول كدا .. ايه الي هيحصل لو كل الترمات الي في ال $p(x|y=\text{real})$ اللي هي عبارته عن ضرب من واحد لحد M ال $p(x_m|y=\text{real})$.. ايه اللي هيحصل لو ان كل الترمات ديه كانت ليهم high probability ماعدا الترم الي اسمو unicorn .. دلوقت معناه ان عندك مشكله .. المشكله ديه بتبقا ان ال article للي انت ماسكها ف ايدك ديه ... كان شكلها فعلاً خلاص a real news article ... ماعدا لكلمة unicorn ... طيب اللي هنعملو دلوقت هو ازاوي هنحل المشكله ديه .. ب ال MAP estimation .. وديه ال recipe بتاعتها:

Recipe for Closed-form MAP Estimation

1. Assume data was generated i.i.d. from some model (i.e. write the generative story)
 $\theta \sim p(\theta)$ and then for all i : $x^{(i)} \sim p(x|\theta)$ *sampled*
2. Write log-likelihood *prior*
 $\ell_{\text{MAP}}(\theta) = \log p(\theta) + \log p(x^{(1)}|\theta) + \dots + \log p(x^{(N)}|\theta)$
3. Compute partial derivatives (i.e. gradient)
 $\partial \ell_{\text{MAP}}(\theta) / \partial \theta_1 = \dots$
 $\partial \ell_{\text{MAP}}(\theta) / \partial \theta_2 = \dots$
 \dots
 $\partial \ell_{\text{MAP}}(\theta) / \partial \theta_M = \dots$
4. Set derivatives to zero and solve for θ
 $\partial \ell_{\text{MAP}}(\theta) / \partial \theta_m = 0$ for all $m \in \{1, \dots, M\}$
 $\theta^{\text{MAP}} = \text{solution to system of } M \text{ equations and } M \text{ variables}$
5. Compute the second derivative and check that $\ell(\theta)$ is concave down at θ^{MAP}

بدل ما بنشتغل باللوغ لايلي هوود بتاع الداتا .. بنشتغل بال L_{MAP} اللي بيضيف الترم بتاع ال $\log(\text{prior})$ وبعديه بقيت ترمات ال $\log \text{likelihood}$..
 وبعدين نحسب ال partial derivatives

$\phi \sim \text{Uniform}([0,1])$

Model 1: Bernoulli Naïve Bayes

MAP Estimation (Beta Prior)

1. Generative Story:

The parameters are drawn once for the entire dataset.

for $m \in \{1, \dots, M\}$:

for $y \in \{0, 1\}$:

$\theta_{y,m} \sim \text{Beta}(\alpha, \beta)$

for $i \in \{1, \dots, N\}$:

$y^{(i)} \sim \text{Bernoulli}(\phi)$

for $m \in \{1, \dots, M\}$:

$x_m^{(i)} \sim \text{Bernoulli}(\theta_{y^{(i)},m})$

$$N_{y=1} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 1)$$

$$N_{y=0} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 0)$$

$$N_{y=0, x_m=1} = \sum_{i=1}^N \mathbb{I}(y^{(i)} = 0 \wedge x_m^{(i)} = 1)$$

...

2. Likelihood:

$$\ell_{MAP}(\phi, \theta) = \log [p(\phi, \theta | \alpha, \beta) p(\mathcal{D} | \phi, \theta)]$$

$$= \log \left[\left(p(\phi | \alpha, \beta) \prod_{m=1}^M p(\theta_{0,m} | \alpha, \beta) \right) \left(\prod_{i=1}^N p(\mathbf{x}^{(i)}, y^{(i)} | \phi, \theta) \right) \right]$$

3. MAP Estimators: $(\phi^{MAP}, \theta^{MAP}) = \arg\max_{\phi, \theta} \ell_{MAP}(\phi, \theta)$

Take derivatives, set to zero and solve...

$$\phi = \frac{N_{y=1}}{N}$$

$$\theta_{0,m} = \frac{(\alpha - 1) + N_{y=0, x_m=1}}{(\alpha - 1) + (\beta - 1) + N_{y=0}}$$

$$\theta_{1,m} = \frac{(\alpha - 1) + N_{y=1, x_m=1}}{(\alpha - 1) + (\beta - 1) + N_{y=1}}$$

$$\forall m \in \{1, \dots, M\}$$

$$\alpha = 2$$

$$\beta = 2$$

كدا انت البس
كنزل كحروما
لهيقي بعين

كدا انت لما هتقابل كلمة unicon تاني .. هيفي ليها low probability بس هتبقا overwhelmed ببقيت الكلمات اللي فعلاً بتخلي ال Doc اللي معاك دلوقت .. real news article .. وهو ده الهدف من ان احنا ليه بنعوز prior distribution of the parameters .. عشان ال prior بتقلك

"It's extremely unlikely that one of the parameters should be an exact zero"

Model 2: Multinomial Naïve Bayes

Support:

Option 1: Integer vector (word IDs)

+ or num of features

$\mathbf{x} = [x_1, x_2, \dots, x_M]$ where $x_m \in \{1, \dots, K\}$ a word id.

0, 1, 2, ..., K-1

each feature is a count

Generative Story:

for $i \in \{1, \dots, N\}$:

(i)

Bernoulli(ϕ)

for $j \in \{1, \dots, M_i\}$:

$\sim \text{Multinomial}(\theta_{y(i)}, 1)$

newname
10/20/2021 9:04:09 AM

example هنا لكل هتسامل ال y من برنولي ديستريبيوشن .. و هتسامل كل فيتشر عندنا من k-sided .. وده ان عندك Multinomial distribution dice roll .. هتأخذ dye و بعدين ت generate value هتطلع بقيمة و بالطريقة ديه انت بت لل x_j

Model:

newname
10/20/2021 9:02:06 AM

ررقم من واحد ل K ... كذا الدوكيومنت ريزرنتيشن منش هتبقا مجرد انك تقول هل الكلمة ديه موجودة ولا .. انما هتقول how many time does this word appear in that document

الدكتور بيقل ان الموديل هو هو اللي كنا بنعملو قبل كذا . هنضرب البروبابيلتي بتاعت واي في (لكل فيتشر) اضرب احتمالية الفيتشر ديه تحت شرط الواي الفكرة بس ان عندك

K sided dye for each of the possible outcomes of y

$$p_{\phi, \theta}(\mathbf{x}, y) = p_{\phi}(y) \prod_{k=1}^K p_{\theta_k}(x_k | y)$$

$$= (\phi)^y (1 - \phi)^{(1-y)} \prod_{j=1}^{M_i} \theta_{y, x_j}$$

32

Bernoulli

Multi class

Gaussian N.B.

(Multiclass)

Multinomial($\vec{\theta}, 1$)

$$y \sim \text{Bernoulli}(\phi) = p(y|\phi) \leftarrow \text{pdf}$$

$$x_1 \sim \text{Gaussian}(\mu_{y,1}, \sigma_{y,1}^2) = f(x_1 | \dots) \leftarrow \text{pdf}$$

$$x_2 \sim \text{Gaussian}(\mu_{y,2}, \sigma_{y,2}^2) = f(x_2 | \dots) \leftarrow \text{pdf}$$

$$\vdots$$

$$x_M \sim \text{Gaussian}(\mu_{y,M}, \sigma_{y,M}^2) = f(x_M | \dots) \leftarrow \text{pdf}$$

Mth features

in M Ap Likelihood

pdf * pdfs

MAP Estimation and Regularization

MAP

Regularization

$$\hat{\theta}^{MAP} = \underset{\theta}{\operatorname{argmax}} \log p(D|\theta) + \log p(\theta)$$

$$= \underset{\theta}{\operatorname{argmin}} \underbrace{-\log p(D|\theta)}_{J(\theta|D)} + \underbrace{-\log p(\theta)}_{r(\theta)}$$

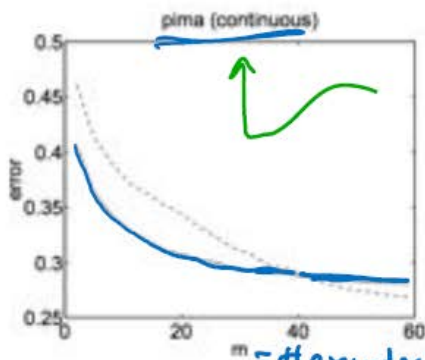
fit the data

keep the model "simple"

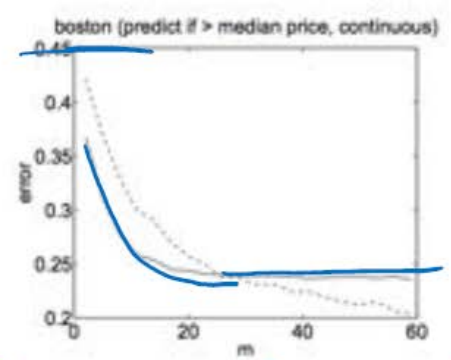
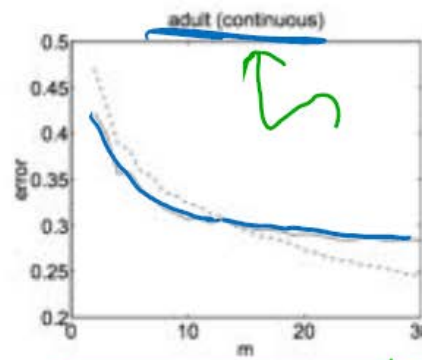
L2 Reg. \equiv Gaussian Prior

Regularizer

show that



m = #examples



solid: NB dashed: LR

Logistic Regression

