

MLE – LS – MAP

CS385 – Machine Learning - Classification

Review previous

CS385 – Machine Learning

Spam detection – Naïve Bayes

$$p(\text{Spam})=3/8$$

Prior

receive a message, not checking it, we believe that with 3/8 probability it is a junk mail

After checking, “secret is secret”, we believe that with 25/26 probability, by using Naïve Bayes, it is a junk mail

posterior

↑

$P(\text{spam} | \text{"secret" is secret})$

$$= \frac{P(\text{"secret"} | \text{spam}) P(\text{"is"} | \text{spam}) \cdot P(\text{"secret"} | \text{spam}) P(\text{spam})}{P(\text{"secret"} | \text{Ham}) \cdot P(\text{"is"} | \text{Ham}) \cdot P(\text{"secret"} | \text{Ham}) P(\text{Ham})}$$

* Bayes Rules

* Conditional Independent.

$C_1 \ C_2 \ \dots \ C_k$

X_1	\dots	X_n	Y (k classes)

$(val_1, val_2, \dots, val_n)$

which class?

we are

$$c_i: \begin{cases} P(C_1 | val_1, val_2, \dots, val_n) \\ P(C_2 | val_1, val_2, \dots, val_n) \\ \vdots \\ P(C_k | val_1, \dots, val_n) \end{cases}$$

$$P(C_i | val_1, \dots, val_n) = \frac{\prod_j P(val_j | C_i)}{P(val_1, \dots, val_n)}$$

Maximum Likelihood Estimate

$\pi = p(\text{Spam}) = 3/8$ how do we learn it?

$p(\text{'is'}|\text{Spam}) = 1/15$ how do we learn it?

$$p(\text{data set} | \pi)$$

$$= p(m_1, m_2, \dots, m_n | \pi)$$

i. i)

$$= p(m_1 | \pi) p(m_2 | \pi) \dots p(m_n | \pi)$$

$$\text{argmax}_{\pi}$$

Laplace Smoothing

CS385 – Machine Learning

Overfitting – Laplace Smoothing

$$P(\text{Spam} | \text{"Today is Secret"}) = 0$$

$$\text{MLE } P(x) = \text{count}(x) / N$$

$$\text{LS } P(x) = (\text{count}(x) + k) / (N + k|x|)$$

$$K=1 \quad p(\text{spam}) = ?$$

$$P(X = \text{Spam})$$

$$P(X = \text{Ham})$$

<u>1 message</u>	<u>1 spam</u>
10 messages	6 spam
100 messages	60 spam

$$2/3$$

$$7/12$$

$$61/102$$

$$= \frac{1+1}{1+1 \times 2}$$

$$\frac{6+1}{10+1 \times 2} = \frac{7}{12}$$

$$P(\text{"Today"} | \text{spam}) \times$$

$$P(\text{"is"} | \text{spam}) \times$$

$$P(\text{"Secret"} | \text{spam}) \times$$

$$P(\text{Spam})$$

Spam Detection – Training Data

Spam

- Offer is secret
- Click secret link
- Secret sports link

Ham

- Play sports today
- Went play sports
- Secret sports event
- Sport is today
- Sport costs money

1 Size of vocabulary? 12

2 $P(\text{Spam}) = ?$ $\frac{3+1}{8+2} = \frac{4}{10} = \frac{2}{5}$

LS Solutions for conditional probability

Spam

- Offer is secret
- Click secret link
- Secret sports link

$p(w = \text{'secret'} | \text{spam})$

$$1 \ P(\text{"Secret"} | \text{Spam}) = \frac{3+1}{9+1 \times 12} = \frac{4}{21}$$

$$2 \ P(\text{"Secret"} | \text{Ham}) = \frac{1+1}{15+1 \times 12} = \frac{2}{27}$$

Ham

- Play sports today
- Went play sports
- Secret sports event
- Sport is today
- Sport costs money

Spam

- Offer is secret
- Click secret link
- Secret sports link

Ham

- Play sports today
- Went play sports
- Secret sports event
- Sport is today
- Sport costs money

$P(\text{Spam} | \text{"Today is secret"})$

$$P(\text{"is"} | \text{Spam}) = \frac{1+1}{9+12} = \frac{2}{21}$$

$$P(\text{"is"} | \text{Ham}) = \frac{1+1}{15+12} = \frac{2}{27}$$

$$P(\text{"today"} | \text{Spam}) = \frac{0+1}{9+12} = \frac{1}{21}$$

$$P(\text{"today"} | \text{Ham}) = \frac{2+1}{15+12} = \frac{3}{27}$$

LS(k=1) for detection

Spam

- Offer is secret
- Click secret link
- Secret sports link

Ham

- Play sports today
- Went play sports
- Secret sports event
- Sport is today
- Sport costs money

$$P(\text{Spam} | \text{"Today is Secret"}) = 0.4858$$

LS(k=1) for Conditional Probability

	Spam	Spam(LS)	Ham	Ham(LS)
offer	1	1	0	1
is	1	1	1	1
Secret	3	+	1	+
click	1	1	0	1
<u>link</u>	2	1	0	1
Sport	1	1	5	1
Play	0	1	2	1
Today	0	1	2	1
Go	0	1	1	1
Event	0	1	1	1
Cost	0	1	1	1
Money	0	1	1	1
	9	+	15	+
		12		12

8

3 Ham

5 Spam

MAP vs. MLE

CS385 – Machine Learning

Finding π

MLE: choose π that maximizes probability of observed data

$$\hat{\pi} = \arg \max_{\pi} P(data \mid \pi)$$

MAP (Maximum A Posterior): choose π that is most probable given prior probability and the data

$$\begin{aligned}\hat{\pi} &= \arg \max_{\pi} P(\pi \mid data) \\ &= \arg \max_{\pi} \frac{P(data \mid \pi)P(\pi)}{P(data)}\end{aligned}$$

MLE – π , conditional probability θ

$$\hat{\pi} = \arg \max_{\pi} P(\text{data} \mid \pi)$$

$$\pi = P(\text{spam})$$

$$\hat{\pi} = \frac{\# \text{data}\{Y = 1\}}{|\text{data}|} = \frac{3}{8}$$

$$\hat{\theta}_{y=1} = \frac{\# \text{data}\{X_i = x_{ij} \wedge Y = 1\}}{\# \text{data}\{Y = 1\}} = \frac{3}{9}$$

$\theta = P(\text{word} \mid \text{spam})$
↓
secret.

Number of items in
dataset for which $Y=1$

MAP for Binomial (1)

$P(\pi | data)$

$$\hat{\pi} = \arg \max_{\pi}$$

$$\frac{P(data | \pi)P(\pi)}{P(data)}$$

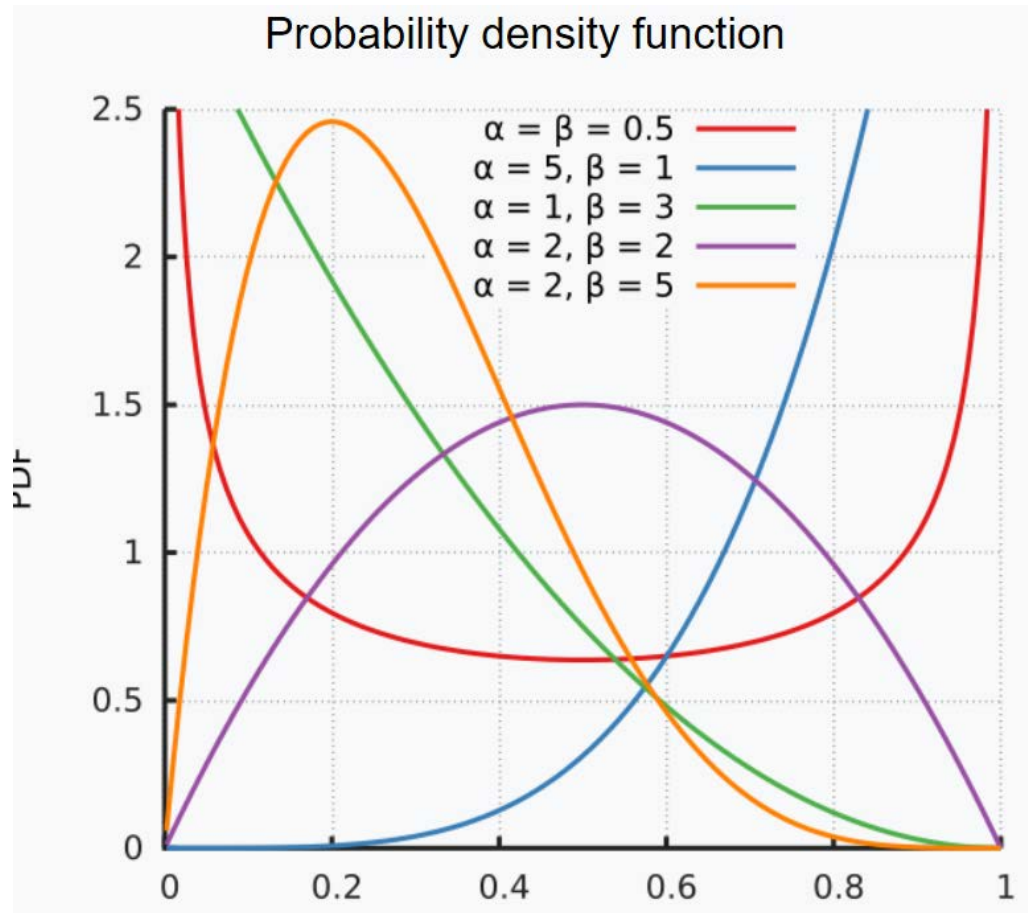
Conjugate priors: $P(\pi)$ and $P(\pi | data)$ have the same form, say, Beta distribution for Binomial

Likelihood $P(data | \pi)$ is Binomial

$$P(data | \pi) = \pi^{count(y_i=1)} \cdot (1 - \pi)^{count(y_i=0)}$$

$$= \pi^3 \cdot (1 - \pi)^5$$

Beta Distribution



Notation	Beta(α, β)
Parameters	$\alpha > 0$ shape (real) $\beta > 0$ shape (real)
Support	$x \in [0, 1]$ or $x \in (0, 1)$
PDF	$\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$ <p>where $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$</p>

$$f(x; \alpha, \beta) = \text{constant} \cdot x^{\alpha-1}(1-x)^{\beta-1}$$

$$= \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\int_0^1 u^{\alpha-1}(1-u)^{\beta-1} du}$$

$$= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$

$$= \frac{1}{B(\alpha, \beta)} x^{\alpha-1}(1-x)^{\beta-1}$$

MAP for Binomial (2)

Prior is Beta distribution

$$P(\pi) = \frac{\pi^{\alpha_1-1}(1-\pi)^{\alpha_0-1}}{B(\alpha_1, \alpha_0)} \sim \text{Beta}(\alpha_1, \alpha_0)$$

Posterior is Beta distribution

$$\begin{aligned} P(\pi \mid \text{data}) &= \frac{\pi^{\text{count}(y_i=1)} \cdot (1-\pi)^{\text{count}(y_i=0)} \cdot \frac{\pi^{\alpha_1-1}(1-\pi)^{\alpha_0-1}}{B(\alpha_1, \alpha_0)}}{\int_0^1 \pi^{\text{count}(y_i=1)} \cdot (1-\pi)^{\text{count}(y_i=0)} \cdot \frac{\pi^{\alpha_1-1}(1-\pi)^{\alpha_0-1}}{B(\alpha_1, \alpha_0)} d\pi} \\ &= \frac{\pi^{\text{count}(y_i=1)+\alpha_1-1} \cdot (1-\pi)^{\text{count}(y_i=0)+\alpha_0-1}}{\int_0^1 \pi^{\text{count}(y_i=1)+\alpha_1-1} \cdot (1-\pi)^{\text{count}(y_i=0)+\alpha_0-1}} \\ &= \frac{\pi^{\text{count}(y_i=1)+\alpha_1-1} \cdot (1-\pi)^{\text{count}(y_i=0)+\alpha_0-1}}{B(\text{count}(y_i=1) + \alpha_1, \text{count}(y_i=0) + \alpha_0)} \\ &\sim \text{Beta}(\text{count}(y_i=1) + \alpha_1, \text{count}(y_i=0) + \alpha_0) \end{aligned}$$

MAP for Binomial (3)

MAP

$$\hat{\pi} = \frac{\# \text{data}\{Y = 1\} + \alpha_1}{\underbrace{|\text{data}|}_{\text{Q}} + \alpha_1 + \alpha_0}$$

Handwritten red notes: "3 + α 1" above the numerator and "Q" below the denominator.

MLE

$$\hat{\pi} = \frac{\# \text{data}\{Y = 1\}}{|\text{data}|}$$

$\alpha_0 \alpha_1$: "Imaginary"
Examples

LS

$$\hat{\pi} = \frac{\# \text{data}\{Y = 1\} + k}{|\text{data}| + 2k}$$

$P(\text{word}|\text{Spam})$ Binomial Multinomial ?

- $P(\text{spam})$ $P(\text{ham})$
- An email is either a spam or a ham
- $P(\text{word}|\text{spam})$
- A word can be “secret”, “today”, “sport”,
- A word can be one of the words in the dictionary
- So it follows a multinomial distribution

MAP for multinomial (1)

$\eta \rightarrow p(\text{spam}) \quad \theta_1$
 $1-\eta \rightarrow p(\text{ham}) \quad \theta_2$
 $p(\text{word}) \dots \dots$
 $\{\theta_1, \theta_2, \theta_3\}$

Dirichlet distribution

Likelihood is $\sim \text{Multinomial}(\theta = \{\theta_1, \theta_2, \dots, \theta_k\})$

$$P(\mathcal{D} | \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$\theta_1^{\beta_1-1} \cdot \theta_2^{\beta_2-1} \cdot \theta_3^{\beta_3-1}$

$$P(\theta) = \frac{\prod_{i=1}^k \theta_i^{\beta_i-1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta | \mathcal{D}) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \dots, \beta_k + \alpha_k)$$

MAP for multinomial (2)

scored spam!

$$\text{MAP} \quad \hat{\theta}_{y=1} = \frac{\# \text{data}\{X_i = x_{ij} \wedge Y = 1\} + \beta_{(y=1, x_i=x_{ij})}}{\# \text{data}\{Y = 1\} + \sum_m \beta_{(y=1, x_i=x_{im})}}$$

$3 + \beta_i$
 $9 + \beta_1 + \beta_2 + \dots + \beta_{12}$

MLE

$$\hat{\theta}_{y=1} = \frac{\# \text{data}\{X_i = x_{ij} \wedge Y = 1\}}{\# \text{data}\{Y = 1\}}$$

LS

$$\hat{\theta}_{y=1} = \frac{\# \text{data}\{X_i = x_{ij} \wedge Y = 1\} + k}{\# \text{data}\{Y = 1\} + |\text{vocabulary}| \cdot k}$$

$\beta_1 = \beta_2 = \dots = \beta_{12} = 1$

Advanced Spam Filters

- Known spamming IP?
- Have you emailed person before?
- Have 1000 other people received same message?
- Email header consistent
- All caps?
- Are you addressed by name?

Gaussian Naïve Bayes

CS385 – Machine Learning

Continuous value

Own Property	Marital Status	Annual Revenue(k)	Delinquency
Y	Single	125	N
N	Married	100	N
N	Single	70	N
Y	Married	120	N
N	Divorced	95	Y
N	Married	60	N
Y	Divorced	220	N
N	Single	85	Y
N	Married	75	N
N	Single	90	Y

Continuous X_i

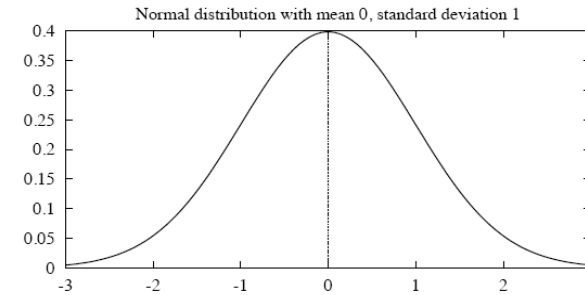
- Y_i is discrete, but X_i is real (continuous)

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

- Common approach
 - Assume $P(X_i | Y = y_k)$ follows a Normal (Gaussian) distribution

Normal (Gaussian) distribution

- $p(x)$ is a probability density function whose integral is 1



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

The probability that X will fall into the interval (a, b) is given by

$$\int_a^b p(x) dx$$

- Expected, or mean value of X , $E[X]$, is

$$E[X] = \mu$$

- Variance of X is

$$Var(X) = \sigma^2$$

- Standard deviation of X , σ_X , is

$$\sigma_X = \sigma$$

GNB - Gaussian Naïve Bayes

- GNB assume
$$p(X_i = x|Y = y_k) = \frac{1}{\sqrt{2\pi\sigma_{ik}^2}} e^{-\frac{1}{2}\left(\frac{x-\mu_{ik}}{\sigma_{ik}}\right)^2}$$
- Train Naïve Bayes
 - For each y_k , estimate $\pi_k = p(Y = y_k)$ [how many parameters? $n-1$, n : number of classes]
 - For each attribute x_i , estimate $p(x_i|Y=y_k)$, class conditional mean μ_{ik} , variance σ_{ik} [how many parameters? $2dn$, d : number of attributes]

▪ Classify(X^{new})

$$Y^{\text{new}} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{\text{new}}|Y = y_k)$$
$$Y^{\text{new}} \leftarrow \arg \max_{y_k} \pi_k \prod_i \mathcal{N}(X_i^{\text{new}}; \mu_{ik}, \sigma_{ik})$$

MLE: Y discrete, X_i continuous

- $\pi_k = p(Y = y_k)$
 - Same as Naïve Bayes Discrete
- $p(x_i | Y = y_k)$

$$\hat{\mu}_{ik} = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_i^j \delta(Y^j = y_k)$$

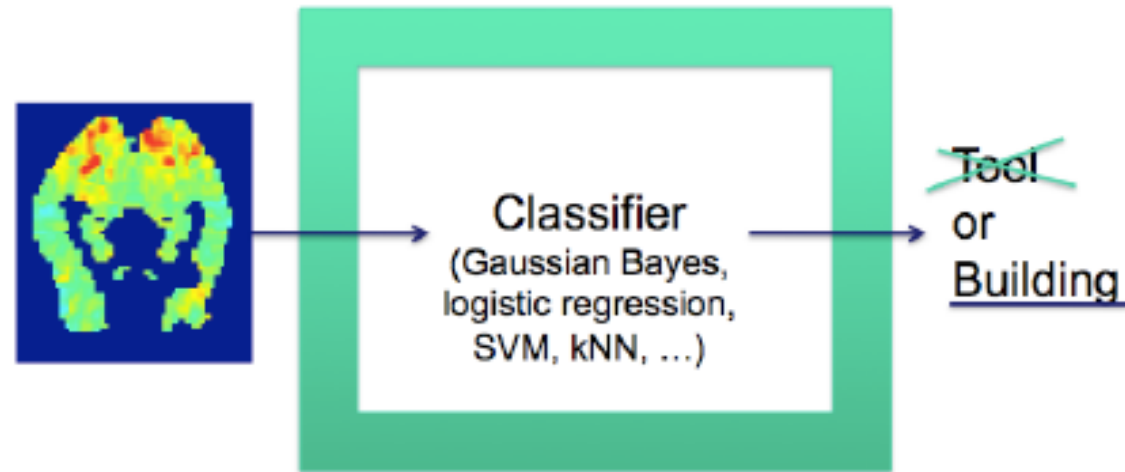
Diagram illustrating the formula for $\hat{\mu}_{ik}$ (mean of feature i for class k):

- $\hat{\mu}_{ik}$: i th feature, k th class
- X_i^j : j th training example
- $\delta(Y^j = y_k)$: $\delta() = 1$ if $(Y^j = y_k)$ else 0

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_i^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

GNB Example

- Classify a person's cognitive state based on brain image
- Reading a word describing a “tool” or “building”?



Naïve Bayes – Conclusion

- Designing classifiers based on Bayes rule
- Conditional independence
- How to train Naïve Bayes classifiers
 - MLE and MAP estimates
 - with discrete and/or continuous inputs X_i