

Bayesian classification

CISC 5800

Professor Daniel Leeds

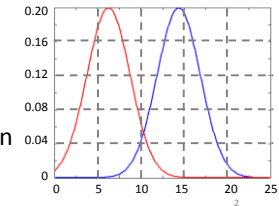
Classifying with probabilities

Example goal: Determine is it cloudy out

- Available data: Light detector: $x \in [0,25]$
- Potential class (atmospheric states): $Y=\{\text{Cloudy}, \text{Non-Cloudy}\}$

Each class (atmospheric state) y has associated probability distribution $P(x)$

Actually each y has a **likelihood** distribution $P(x|\mu_y, \sigma_y)$

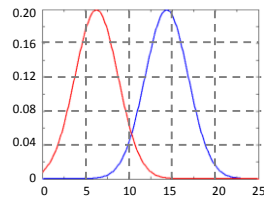


Classifying with probabilities

Example goal: Determine is it cloudy out

- Measure light: x
- Compute $P(x|\mu_y, \sigma_y)$ for $y=\text{Cloudy}$ and $y=\text{Non-Cloudy}$
- Pick y which gives greatest **likelihood** $P(x|\mu_y, \sigma_y)$

$$\operatorname{argmax}_y P(x|\mu_y, \sigma_y)$$



$$x=9$$

$$P(x=9 | \text{Cloudy})=0.12$$

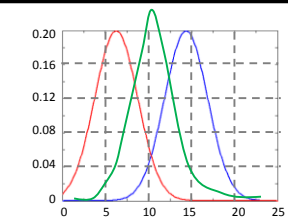
$$P(x=9 | \text{Non-Cloud})=0.02$$

This is **Maximum Likelihood** classification

4

What if there's an eclipse?

- Let's add a third potential class: $Y=\{\text{Cloudy}, \text{Non-Cloudy}, \text{Eclipse}\}$
- What is most likely class if $x=9$?
- Eclipses are low probability!



$x=9$

$$P(x=9 | \text{Cloudy})=0.12$$

$$P(x=9 | \text{Non-Cloud})=0.02$$

$$P(x=9 | \text{Eclipse})=0.16$$

6

Incorporating prior probability

- Define **prior** probabilities for each class $P(y) = P(\mu_y, \sigma_y)$
Probability of class y same as probability of parameters μ_y, σ_y
- **“Posterior probability”** estimated as likelihood \times prior :
 $P(x|\mu_y, \sigma_y) P(\mu_y, \sigma_y)$
- Classify as $\operatorname{argmax}_y P(x|\mu_y, \sigma_y) P(\mu_y, \sigma_y)$
- Terminology: μ_y, σ_y are “parameters.” In general use θ_y
Here: $\theta_y = \{\mu_y, \sigma_y\}$. **“Posterior”** estimate is $P(x|\theta_y) P(\theta_y)$

7

Probability review: Bayes rule

Recall: $P(A|B) = \frac{P(A,B)}{P(B)}$

and: $P(A, B) = P(B|A)P(A)$

so: $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$

Equivalently: $P(y|x) = P(\theta_y|x) = P(\theta_y|D) = \frac{P(D|\theta_y) P(\theta_y)}{P(D)}$

The true posterior



8

The posterior estimate

$$\operatorname{argmax}_{\theta_y} P(\theta_y|\mathbf{D}) \propto P(\mathbf{D}|\theta_y)P(\theta_y)$$

Posterior \propto Likelihood \times Prior \propto - means proportional

We “ignore” the $P(\mathbf{D})$ denominator because \mathbf{D} stays same while comparing different classes (y represented by θ_y)

9

Typical classification approaches

MLE – Maximum Likelihood: Determine parameters/class which maximize probability of the data

$$\operatorname{argmax}_{\theta_y} P(\mathbf{D}|\theta_y)$$

MAP – Maximum A Posteriori: Determine parameters/class that has maximum probability

$$\operatorname{argmax}_{\theta_y} P(\theta_y|\mathbf{D})$$

10

Incorporating a prior

Three classes:

$Y = \{\text{Cloudy}, \text{Non-Cloudy}, \text{Eclipse}\}$

$P(\text{Cloudy}) = 0.4$

$P(\text{Non-Cloudy}) = 0.4$

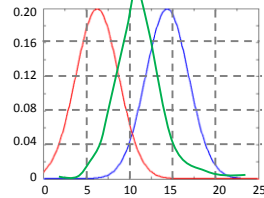
$P(\text{Eclipse}) = 0.2$

$x=9$

$$P(x=9 | \text{Cloudy}) P(\text{Cloudy}) = 0.12 \times 0.4 = .048$$

$$P(x=9 | \text{Non-Cloudy}) P(\text{Non-Cloudy}) = 0.02 \times 0.4 = 0.008$$

$$P(x=9 | \text{Eclipse}) P(\text{Eclipse}) = 0.16 \times 0.2 = .032$$



12

Bernoulli distribution – coin flips

We have three coins with known biases (favoring heads or tails)

How can we determine our current coin?

Flip K times to see which bias it has

Data (\mathbf{D}): {HHTH, TTTH, TTTT} Bias (θ_y): p_y probability of H for coin y

$$P(\mathbf{D} | \theta_y) = p_y^{|\mathbf{H}|} (1 - p_y)^{|\mathbf{T}|} \quad |\mathbf{H}| - \# \text{ heads}, \quad |\mathbf{T}| - \# \text{ tails}$$

14

Bernoulli distribution – reexamined

$$P(\mathbf{D} | \theta_y) = p_y^{|\mathbf{H}|} (1 - p_y)^{|\mathbf{T}|} \quad |\mathbf{H}| - \# \text{ heads}, \quad |\mathbf{T}| - \# \text{ tails}$$

More rigorously: in K trials, $side_k = \begin{cases} 0 & \text{if tails on flip } k \\ 1 & \text{if heads on flip } k \end{cases}$

$$P(\mathbf{D} | \theta_y) = \prod_k p_y^{side_k} (1 - p_y)^{(1 - side_k)}$$

15

Multinomial example

4-sided die – 4 probabilities:

$$p_{side1}, p_{side2}, p_{side3}, p_{side4} \quad (\text{Note: } p_{side4} = 1 - \sum_{k=1}^3 p_{sidek})$$



$$\text{Define: } \delta(x) = \begin{cases} 1 & x = 0 \\ 0 & \text{otherwise} \end{cases}$$

$$P(\mathbf{D} | \theta_y) = \prod_k p_{side1}^{\delta(side_k-1)} p_{side2}^{\delta(side_k-2)} p_{side3}^{\delta(side_k-3)} p_{side4}^{\delta(side_k-4)}$$

16

Optimization: finding the maximum likelihood parameter for a fixed class (fixed coin)

$$\operatorname{argmax}_{\theta} P(\mathbf{D}|\theta_y) = \quad p_y - \text{probability of Head}$$

$$\operatorname{argmax}_p p_y^{|H|} (1 - p_y)^{|T|}$$

Equivalently, maximize $\log P(\mathbf{D}|\theta_y)$

$$\operatorname{argmax}_{p_y} |H| \log p_y + |T| \log(1 - p_y)$$

18

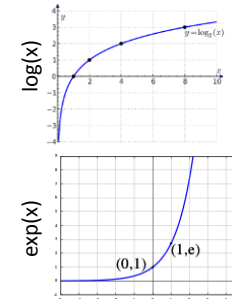
The properties of logarithms

$$e^a = b \leftrightarrow \log b = a$$

$$a < b \leftrightarrow \log a < \log b$$

$$\log ab = \log a + \log b$$

$$\log a^n = n \log a$$



Convenient when dealing with small probabilities

- $0.0000454 \times 0.000912 = 0.0000000414 \rightarrow -10 + -7 = -17$

19

Optimization: finding zero slope

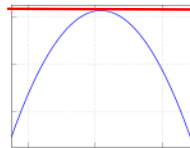
Location of maximum has slope 0 p - probability of Head

maximize $\log P(\mathbf{D}|\theta)$

$$\operatorname{argmax}_p |H| \log p + |T| \log(1 - p) :$$

$$\frac{d}{dp} |H| \log p + |T| \log(1 - p) = 0$$

$$\frac{|H|}{p} - \frac{|T|}{1-p} = 0$$



21

Intuition of the MLE result

$$p_y = \frac{|H|}{|H| + |T|}$$

- Probability of getting heads is # heads divided by # total flips

22

Finding the maximum a posteriori

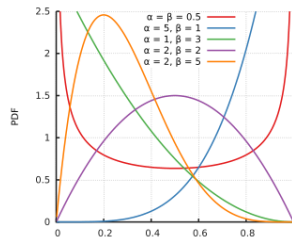
$$P(\theta_y | \mathbf{D}) \propto P(\mathbf{D} | \theta_y) P(\theta_y)$$

- Incorporating the Beta prior:

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

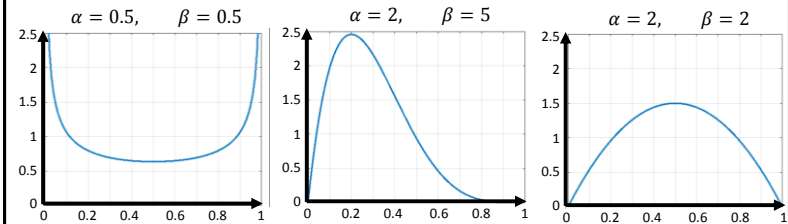
$$\operatorname{argmax}_{\theta} P(\mathbf{D} | \theta_y) P(\theta_y) =$$

$$\operatorname{argmax}_{\theta} \log P(\mathbf{D} | \theta_y) + \log P(\theta_y)$$



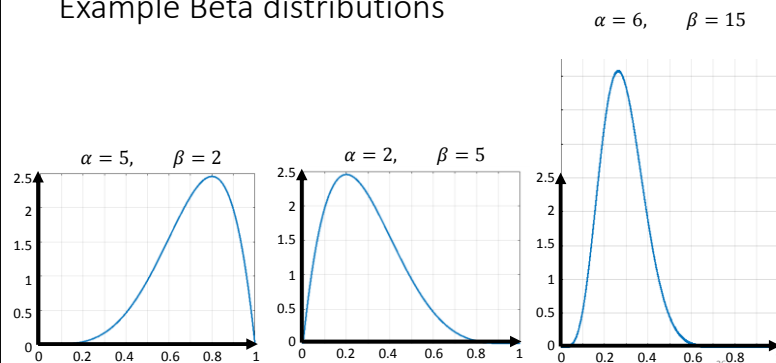
23

Example Beta distributions



25

Example Beta distributions



26

MAP: estimating θ (estimating p)

$$\operatorname{argmax}_{\theta} \log P(\mathbf{D} | \theta) + \log P(\theta)$$

$$\operatorname{argmax}_p |H| \log p + |T| \log(1-p) +$$

$$(\alpha - 1) \log p + (\beta - 1) \log(1-p) - \log(B(\alpha, \beta))$$

↓ *Set derivative to 0*

$$\frac{|H|}{p} - \frac{|T|}{1-p} + \frac{(\alpha - 1)}{p} - \frac{(\beta - 1)}{1-p} = 0$$

$$(1-p)|H| - p|T| + (1-p)(\alpha - 1) - p(\beta - 1) = 0$$

$$|H| + (\alpha - 1) = (|H| + |T| + (\alpha - 1) + (\beta - 1))p$$

27

Intuition of the MAP result

$$p_y = \frac{|H| + (\alpha - 1)}{|H| + (\alpha - 1) + |T| + (\beta - 1)}$$

- Prior has strong influence when $|H|$ and $|T|$ small
- Prior has weak influence when $|H|$ and $|T|$ large

- $\alpha > \beta$ means expect to find coins biased to heads
- $\beta > \alpha$ means expect to find coins biased to tails

28

Multinomial distribution **Classification**

- What is mood of person in current minute? $M = \{\text{Happy, Sad}\}$
- Measure his/her actions every ten seconds: $A = \{\text{Cry, Jump, Laugh, Yell}\}$

Data (\mathbf{D}): {LLJLCY, JJLYJL, CCLLLJ, JJJJJJ}Bias (θ_y): Probability table

	Happy	Sad
Cry	0.1	0.5
Jump	0.3	0.2
Laugh	0.5	0.1
Yell	0.1	0.2

$$P(\mathbf{D}|\theta_y) = (p_y^{\text{Cry}})^{|\text{Cry}|} (p_y^{\text{Jump}})^{|\text{Jump}|} (p_y^{\text{Laugh}})^{|\text{Laugh}|} (p_y^{\text{Yell}})^{|\text{Yell}|}$$

29

Multinomial distribution – reexamined

$$P(\mathbf{D}|\theta_y) = (p_y^{\text{Cry}})^{|\text{Cry}|} (p_y^{\text{Jump}})^{|\text{Jump}|} (p_y^{\text{Laugh}})^{|\text{Laugh}|} (p_y^{\text{Yell}})^{|\text{Yell}|}$$

More rigorously: in K measures,

$$\delta(\text{trial}_k = \text{Action}) = \begin{cases} 0 & \text{if } \text{trial}_k \neq \text{Action} \\ 1 & \text{if } \text{trial}_k = \text{Action} \end{cases}$$

$$P(\mathbf{D}|\theta_y) = \prod_k \prod_i (p_y^{\text{Action}_i})^{\delta(\text{trial}_k = \text{Action}_i)}$$

Classification: Given known likelihoods for each action, find mood that maximizes likelihood of observed sequence of actions (assuming each action is independent in the sequence)

30

Learning parameters

$$\text{MLE: } P(A = a_i | M = m_j) = p_j^i = \frac{\#D\{A=a_i \wedge M=m_j\}}{\#D\{M=m_j\}}$$

$$\text{MAP: } P(A = a_i | M = m_j) = \frac{\#D\{A=a_i \wedge M=m_j\} + (\gamma_i - 1)}{\#D\{M=m_j\} + \sum_k (\gamma_k - 1)}$$

$$P(Y = y_j) = \frac{\#D\{M=m_j\} + (\beta_j - 1)}{|D| + \sum_m (\beta_m - 1)}$$

γ_k is prior probability of each action class a_k

β_k is prior probability of each mood class m_k

32

Multiple multi-variate probabilities

Mood based on Action, Tunes, Weather

	Happy	Sad
Cry, Jazz, Sun	0.003	0.102
Cry, Jazz, Rain	0.024	0.025
	⋮	
Cry, Rap, Snow	0.011	0.115
	⋮	
Laugh, Rap, Rain	0.042	0.007
	⋮	
Yell, Opera, Wind	0.105	0.052

$\operatorname{argmax}_{\theta_y} P(A, T, W | \theta_y)$

How many entries in probability table?

params = $|M| \times (|A| \times |T| \times |W| - 1)$

36

Naïve bayes:

Assuming independence of input features

	Happy	Sad
Jazz	0.05	0.4
Rap	0.5	0.3
Opera	0.45	0.3

$\operatorname{argmax}_{\theta_y} P(A, T, W | \theta_y) =$

	Happy	Sad
Sun	0.6	0.2
Rain	0.05	0.3
Snow	0.3	0.3
Wind	0.05	0.2

$\operatorname{argmax}_{\theta_y} P(A | \theta_y) P(T | \theta_y) P(W | \theta_y)$

How many entries in probability tables?

	Happy	Sad
Cry	0.1	0.5
Jump	0.3	0.2
Laugh	0.5	0.1
Yell	0.1	0.3

• # params = $|M| \times ((|A|-1) + (|T|-1) + (|W|-1)) = 2 \times (3+2+3) = 16$

Benefits of Naïve Bayes

Very fast learning and classifying:

- **For multinomial problem:**
 - **Naïve independence:** learn $|Y| \times \sum_i (|X_i| - 1)$ parameters
 - **Non-naïve:** learn $|Y| \times (\prod_i |X_i| - 1)$ parameters

$|Y|$ is number of possible classes

$|X_i|$ is number of possible values for i^{th} feature

Often works even if features are NOT independent

39

Typical Naïve Bayes classification

$\operatorname{argmax}_{\theta_y} P(\theta_y | \mathbf{D}) \rightarrow \operatorname{argmax}_{\theta_y} P(\mathbf{D} | \theta_y) P(\theta_y)$ $P(\theta_y)$ prior class probability

$P(\mathbf{D} | \theta_y) = \prod_i P(X^i | \theta_y)$ where $\mathbf{D} = \begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix}$ is a list of feature values

e.g., $x^1 = \text{Action}$, $x^2 = \text{Tunes}$

NB (Naïve Bayes): Find class y with θ_y to maximize $P(\theta_y | \mathbf{D})$

Assuming variables independent

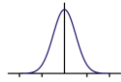
40

Multi-dimensional probability functions

- Multiple features as vector: $\mathbf{x} = \begin{bmatrix} \text{temperature} \\ \text{windSpeed} \\ \text{musicVolume} \end{bmatrix}$

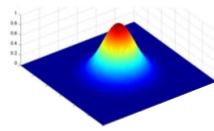
- In 1D: likelihood $P(\text{temperature} \mid \text{mood})$

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



- In 2D: likelihood $P\left(\begin{bmatrix} \text{temp} \\ \text{wind} \end{bmatrix} \mid \text{mood}\right)$

$$L = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^2 |\boldsymbol{\Sigma}|}}$$



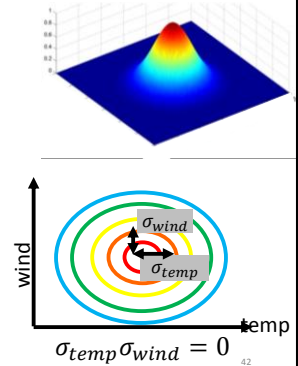
Multi-dimensional Gaussian – Naïve

- In 2D: likelihood $P\left(\begin{bmatrix} \text{temp} \\ \text{wind} \end{bmatrix} \mid \text{mood}\right)$

$$L = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)}{\sqrt{(2\pi)^2 |\boldsymbol{\Sigma}|}}$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{\text{temp}}^2 & \sigma_{\text{temp}}\sigma_{\text{wind}} \\ \sigma_{\text{temp}}\sigma_{\text{wind}} & \sigma_{\text{wind}}^2 \end{bmatrix}$$

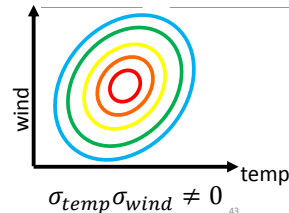
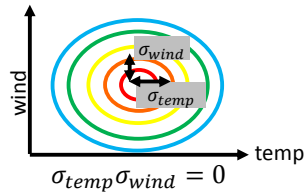
$$\boldsymbol{\mu} = [\mu_{\text{temp}} \quad \mu_{\text{wind}}]$$



Multi-dimensional Gaussian – Non-naïve

- In 2D: likelihood $P\left(\begin{bmatrix} \text{temp} \\ \text{wind} \end{bmatrix} \mid \text{mood}\right)$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{\text{temp}}^2 & \sigma_{\text{temp}}\sigma_{\text{wind}} \\ \sigma_{\text{temp}}\sigma_{\text{wind}} & \sigma_{\text{wind}}^2 \end{bmatrix}$$



Gaussian parameter counts

For k dimensions

- Naïve: $k + k \approx 2k$ parameters
- Non-naïve: $k + \frac{k(k-1)}{2} \approx \frac{k^2}{2}$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \cdots & \sigma_1\sigma_k \\ \vdots & \ddots & \vdots \\ \sigma_1\sigma_k & \cdots & \sigma_k^2 \end{bmatrix}$$