

BAYESIAN DECISION THEORY

SANNE NYGAARD

- Modelling
- Joint probabilities and Bayes' Rule
- Bayesian classification
- Losses, Risk, and Rejection
- Association rules

Modelling

We model data as random process to make up for noise and missing information

EXAMPLE: COIN TOSS

How would we model a coin toss?



Image credit: ICMA Photos

EXAMPLE: COIN TOSS

- There is the *observable variable*, in this case the outcome:
 $Result \in \{Heads, Tails\}$

EXAMPLE: COIN TOSS

- There is the *observable variable*, in this case the outcome:
 $Result \in \{Heads, Tails\}$
- There are *unobservable variables* that we don't have access to, e.g.:
 - ▶ Coin composition
 - ▶ Precise position and orientation
 - ▶ Strength and direction of forces
- If we had access to all info we might calculate the outcome

EXAMPLE: ROULETTE



Image credit: Ralf Roletschek

EXAMPLE: ROULETTE



Image credit: A rather dodgy website

EXAMPLE: COIN TOSS

- Modelled using the Bernoulli distribution:

Result $\in \{Heads, Not\ Heads\} \rightarrow$

Bernoulli random variable: $X \in \{1, 0\}$

$$P\{X = x\} = p^x(1 - p)^{(1-x)}$$

p = probability of '1'

$1 - p$ = probability of '0'

EXAMPLE: COIN TOSS

Estimation

What if we don't know p ? How do we find it?

EXAMPLE: COIN TOSS

We get a sample of observed outcomes:

Sample: $\mathbf{X} = \{x^t\}_{t=1}^N$

Estimation:

$$p = \frac{\#\{Heads\}}{\#\{Tosses\}} = \sum_t \frac{x^t}{N}$$

This illustrates the frequentist approach to probability:

- Probabilities as frequencies of outcomes from repeated experiments
- There is a true value for the parameter p

EXAMPLE: COIN TOSS

Estimation example:

$$\mathbf{X} = \{1, 1, 1, 0, 0, 1, 1, 0, 1\}$$

$$\hat{p} = \frac{\#\{\text{Heads}\}}{\#\{\text{Tosses}\}} = \frac{6}{9}$$



Image source: wikipedia.org

EXAMPLE: COIN TOSS

If instead of nine observations:

$$\mathbf{X} = \{1, 1, 1, 0, 0, 1, 1, 0, 1\}$$

$$\hat{p} = \frac{6}{9}$$

We only had the first three:

$$\mathbf{X} = \{1, 1, 1\}$$

$$\hat{p} = \frac{3}{3} = 1$$

EXAMPLE: COIN TOSS

Predictions

Minimize errors, choose highest probability outcome

Prediction: Heads if $p > \frac{1}{2}$, Tails otherwise
Probability of error: $1 - P(\text{our choice})$

EXAMPLE: COIN TOSS

In summary:

1. Choose a model
2. Use observed data to estimate model parameters
3. Use the trained model to make predictions

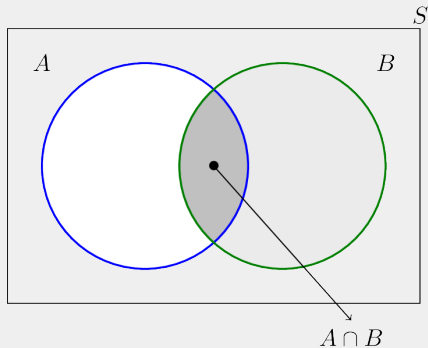
JOINT AND CONDITIONAL PROBABILITIES

Intersection: $P(A \cap B)$

$$\begin{aligned}P(A \cap B) &= P(A|B)P(B) \\ &= P(B|A)P(A)\end{aligned}$$

For independent variables:

$$\begin{aligned}P(A|B) &= P(A) \\ P(A, B) &= P(A)P(B)\end{aligned}$$



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

Image from: www.probabilitycourse.com

BAYES' RULE

This gives us *Bayes' rule*:

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

Example: *The Decaying Fruit Basket*



JOINT PROBABILITY DISTRIBUTIONS

- A not-so-fresh fruit basket with $N = 25$ pieces of fruits: Oranges, Apples, Peaches
- Some of them are good, some have gone bad
- The numbers are:

	Oranges	Apples	Peaches
Good	4	7	3
Bad	3	1	7

JOINT PROBABILITY DISTRIBUTIONS

- Joint distribution over two random variables:
- Fruit type: $X \in \{\text{Oranges}, \text{Apples}, \text{Peaches}\}$
- Freshness: $Y \in \{\text{Good}, \text{Bad}\}$

$$P(X, Y)$$

	Oranges	Apples	Peaches
Good	0.16	0.28	0.12
Bad	0.12	0.04	0.28

JOINT PROBABILITY DISTRIBUTIONS

- The sum rule: marginalizing over one r.v. gives the marginal distribution for the other:

- ▶ Discrete case: $P(X = x) = \sum_j P(X, Y_j)$
- ▶ Continuous case: $p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy$

	Oranges	Apples	Peaches	$P(Y)$
Good	0.16	0.28	0.12	0.56
Bad	0.12	0.04	0.28	0.44
$P(X)$	0.28	0.32	0.40	

JOINT PROBABILITY DISTRIBUTIONS


■ Conditioning

$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

$$P(Y = \text{bad}|X = \text{apple}) = \frac{P(Y = \text{bad}, X = \text{apple})}{P(X = \text{apple})}$$

	Oranges	Apples	Peaches
Good	0.16	0.28	0.12
Bad	0.12	0.04	0.28

$P(Y = \text{bad}, X = \text{apple})$



JOINT PROBABILITY DISTRIBUTIONS

- When picking up a random apple, what is the probability it's bad?

$$P(Y = \text{bad} | X = \text{apple}) = \frac{P(Y=\text{bad}, X=\text{apple})}{P(X=\text{apple})} = \frac{0.04}{0.32} = 0.125$$

- When handed a bad fruit, what is the probability it's a peach?
?

	Oranges	Apples	Peaches	$P(Y)$
Good	0.16	0.28	0.12	0.56
Bad	0.12	0.04	0.28	0.44
$P(X)$	0.28	0.32	0.40	

JOINT PROBABILITY DISTRIBUTIONS

- When picking up a random apple, what is the probability it's bad?

$$P(Y = \text{bad} | X = \text{apple}) = 0.125$$

- When handed a bad fruit, what is the probability it's a peach?
 $P(X = \text{peach} | Y = \text{bad})$

	Oranges	Apples	Peaches
Good	4	7	3
Bad	3	1	7

CLASSIFICATION EXAMPLE: CREDIT SCORING

Credit scoring model: Is a potential customer high-risk or not?
We choose two relevant inputs: Income and savings, represented as two random (observable) variables X_1 and X_2

- Input: $\mathbf{x} = [x_1, x_2]^T$
- Desired output, high-risk or not: $C \in \{1, 0\}$
- Random variable conditioned on X_1 and X_2

CLASSIFICATION

Input: $\mathbf{x} = [x_1, x_2]^T$

Desired output: $C \in \{1, 0\}$

We want to know $P(C|X_1, X_2)$ so we can:

$$\text{Choose} = \begin{cases} C = 1 & \text{if } P(C = 1|x_1, x_2) > P(C = 0|x_1, x_2) \\ C = 0 & \text{otherwise} \end{cases}$$

with probability of error = $1 - \max(P(C = 1|x_1, x_2), P(C = 0|x_1, x_2))$

$$\text{Bayes' rule: } P(C|\mathbf{x}) = \frac{p(\mathbf{x}|C)P(C)}{p(\mathbf{x})}$$

- *Prior* probability $P(C)$: how likely it is to observe a class label C , regardless of \mathbf{x} .
 $p(C) \geq 0$ and $p(C = 1) + p(C = 0) = 1$
- Class *likelihood* $p(\mathbf{x}|C)$: how likely it is that, having observed an example with class label C , the example is at \mathbf{x} .
(The distribution of \mathbf{x} for each class)

$$\text{Bayes' rule: } P(C|\mathbf{x}) = \frac{p(\mathbf{x}|C)P(C)}{p(\mathbf{x})}$$

- *Evidence* $p(\mathbf{x})$: probability of observing an example \mathbf{x} at all (regardless of its class).

$$p(\mathbf{x}) = p(\mathbf{x}|C = 0)p(C = 0) + p(\mathbf{x}|C = 1)p(C = 1)$$

- *Posterior* probability $p(C|\mathbf{x})$: how likely it is that, having observed an example \mathbf{x} , its class label is C .

$$p(C = 1|\mathbf{x}) + p(C = 0|\mathbf{x}) = 1$$

(This is what we want to know)

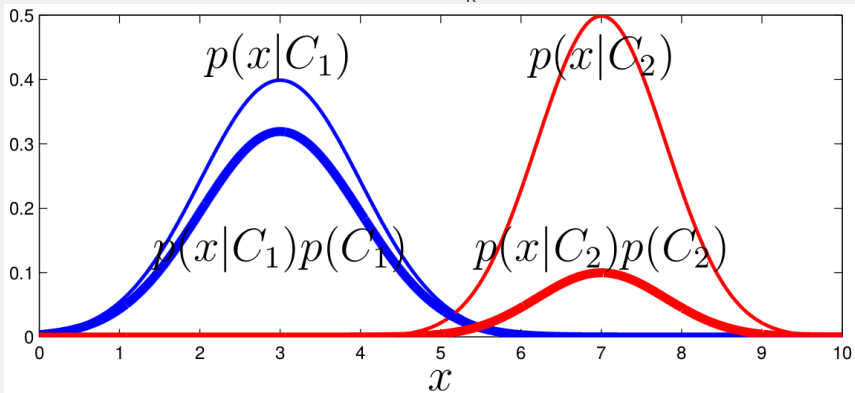
Bayes' rule:

$$P(C|\mathbf{x}) = \frac{p(\mathbf{x}|C)P(C)}{p(\mathbf{x})}$$

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior}}{\textit{evidence}}$$

EXAMPLE, LIKELIHOODS AND PRIORS

Class likelihoods $p(x|C_k) \sim \mathcal{N}(\mu_k, \sigma_k^2)$



with different prior probabilities $p(C_k)$

<http://faculty.ucmerced.edu/mcarreira-perpinan>

EXAMPLE: DISEASE TEST RESULTS (EX. 1)

- Rare disease (d) with occurrence 10^{-6}
- Test (t) for disease:
 - ▶ 99% chance of finding the disease, if present
 - ▶ 10^{-3} chance of wrongly finding the disease, when not present
- A random person is tested, and the result comes out positive:
What is the probability that this person has the disease?

EXAMPLE: DISEASE TEST RESULTS

using 0, 1 encoding for *True, False*

$$\begin{aligned}P(d = 1|t = 1) &= \frac{P(t = 1|d = 1)P(d = 1)}{P(t = 1)} \\&= \frac{P(t = 1|d = 1)P(d = 1)}{P(t = 1|d = 1)P(d = 1) + P(t = 1|d = 0)P(d = 0)} \\&= \frac{0.99 \cdot 10^{-6}}{0.99 \cdot 10^{-6} + 10^{-3} \cdot (1 - 10^{-6})} = 0.00098902\end{aligned}$$

So still more likely that the person does not have the disease

BAYES' RULE $K > 2$ CLASSES

Generalization to more than two classes:

$$P(C_i|\mathbf{x}) = \frac{P(\mathbf{x}|C_i)P(C_i)}{P(\mathbf{x})} = \frac{p(\mathbf{x}|C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x}|C_k)P(C_k)}$$

where:

$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

Choose the class with $\operatorname{argmax}_{k=1,\dots,K} p(C_k|\mathbf{x})$

LOSSES AND RISKS

- So far, we've classified according to highest posterior probability
- This minimizes the expected classification error
- But what if some mistakes are more costly than others?
E.g. cancer diagnosis, earthquake prediction



LOSSES AND RISKS

- So far, we've classified according to highest posterior probability
- This minimizes the expected classification error
- But what if some mistakes are more costly than others?
E.g. cancer diagnosis, earthquake prediction



Choose the class with lowest risk instead

LOSSES AND RISKS

Define:

- Actions: α_j
- Loss (cost) of α_j when real class is C_k : λ_{jk}

Example: Loss matrix:

	cancer	normal
cancer	0	1000
normal	1	0

Image from: Bishop 2006

- Actions: α_i
- Loss (cost) of α_i when real class is C_k : λ_{ik}
- Expected risk:

$$R(\alpha_i|\mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k|\mathbf{x})$$

- choose α_i if $R(\alpha_i|\mathbf{x}) = \min_k R(\alpha_i|\mathbf{x})$

A special case is:

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

This gives us:

$$\begin{aligned} R(\alpha_i|\mathbf{x}) &= \sum_{k=1}^K \lambda_{ik} P(C_k|\mathbf{x}) \\ &= \sum_{k \neq i} P(C_k|\mathbf{x}) \\ &= 1 - P(C_i|\mathbf{x}) \end{aligned}$$

THE REJECT OPTION

- Sometimes it is better to discard uncertain classifications (and leave the final decision to human review or follow-up systems).
- to do that we add 'reject' as an extra action: α_{K+1}

THE REJECT OPTION

For the 0/1 loss:

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = K + 1, \ 0 < \lambda < 1 \\ 1 & \text{otherwise} \end{cases}$$

$$R(\alpha_{K+1}|\mathbf{x}) = \sum_{k=1}^K \lambda P(C_k|\mathbf{x}) = \lambda$$

$$R(\alpha_i|\mathbf{x}) = \sum_{k \neq i} P(C_k|\mathbf{x}) = 1 - P(C_i|\mathbf{x})$$

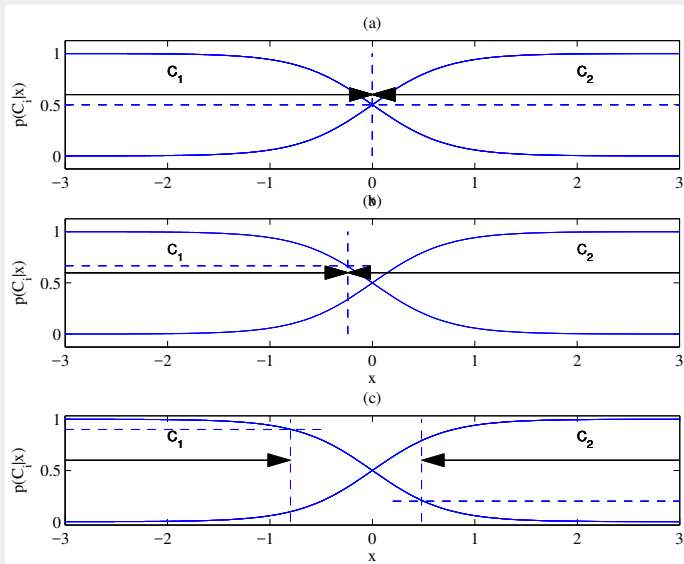
choose $\begin{cases} C_i & \text{if } P(C_i|\mathbf{x}) > P(C_k|\mathbf{x}) \ \forall k \neq i \text{ and } P(C_i|\mathbf{x}) > 1 - \lambda \\ \text{else reject} \end{cases}$

THE REJECT OPTION

In extreme cases of λ :

- $\lambda = 0$: always reject (rejecting is less costly than a correct classification).
- $\lambda = 1$: never reject (rejecting is costlier than any misclassification)

EFFECT OF LOSS AND REJECTION

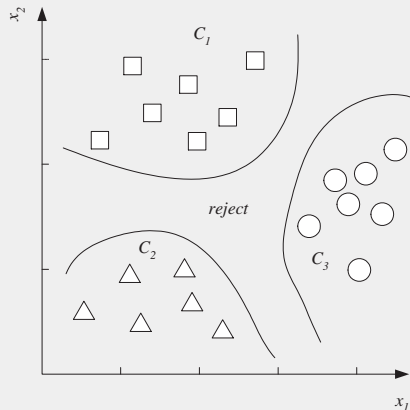


DISCRIMINANT FUNCTIONS

- Classification can also be seen as a set of discriminant functions: $g_i(\mathbf{x})$, $i = 1, \dots, K$
- Choose C_i if $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$
- With 0/1 Loss: $g_i(\mathbf{x}) = P(C_i|\mathbf{x})$
Ignoring $p(x)$: $g_i(\mathbf{x}) = P(\mathbf{x}|C_i)P(C_i)$

DISCRIMINANT FUNCTIONS

- Discriminant functions divide the feature space into K decision regions R_1, \dots, R_K
- where R_i is the region $\{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$
- The regions are separated by decision boundaries where there are ties



ASSOCIATION RULES

- Association rule: $X \rightarrow Y$ (X : antecedent, Y : consequent)
- E.g.: People who {buy/click/visit} X are also likely to {buy/click/visit} Y
- Used for basket analysis and product recommendations
- $X \rightarrow Y$ implies association, not necessarily causation

Frequently bought together



Total price: **£71.84**

[Add both to basket](#)

1. One of these items is dispatched sooner than the other. [Show details](#)

- ✓ **This item:** Introduction to Machine Learning [Adaptive Computation and Machine Learning Series] by [Ethem Alpaydm](#)
- ✓ **The Hundred-Page Machine Learning Book** by [Arndt Borkov](#) Paperback **£27.85**

Customers who viewed this item also viewed



Machine Learning: The New AI (The MIT Press Essential Knowledge...)
• [Ethem Alpaydm](#)
★★★★☆ 3



Pattern Recognition and Machine Learning (Information Science...)
• [Christopher M. Bishop](#)
★★★★☆ 32



Deep Learning (Adaptive Computation and Machine Learning Series)
• [Ian Goodfellow](#)
★★★★☆ 32

BASKET ANALYSIS - MEASURES

Simple but popular measures:

- $\text{Support}(X, Y)$
- $\text{Confidence}(X \rightarrow Y)$
- $\text{Lift}(X \rightarrow Y)$



Image source: dlpng.com

$$\text{Support}(X, Y) \equiv P(X, Y) = \frac{\text{\#purchases of } X \text{ and } Y}{\text{\#purchases}}$$

For the rule to be significant, the support should be large. High support means items X and Y are frequently bought together

$$\text{Confidence}(X \rightarrow Y) \equiv P(Y|X) = \frac{P(X, Y)}{P(X)} = \frac{\text{\#purchases of } X \text{ and } Y}{\text{\#purchases of } X}$$

For the rule to hold with enough confidence, should be $\gg p(Y)$ and close to 1

$$\begin{aligned}\text{Lift}(X \rightarrow Y) &= \frac{P(X, Y)}{P(X)P(Y)} = \frac{P(Y|X)}{P(Y)} \\ &= \frac{\text{\#purchases of X and Y}}{(\text{\#purchases of X}) * (\text{\#purchases of Y})}\end{aligned}$$

This is the ratio of the observed joint probability to that expected under independence

- Lift < 1: X makes Y less likely
- Lift = 1: X and Y are independent
- Lift > 1: X makes Y more likely

BASKET ANALYSIS EXAMPLE (EX. 3.7)

Given the following data of transactions at a shop, calculate the support and confidence values of milk \rightarrow bananas, bananas \rightarrow milk, milk \rightarrow chocolate, and chocolate \rightarrow milk.

Transaction	Items in basket
1	milk, bananas, chocolate
2	milk, chocolate
3	milk, bananas
4	chocolate
5	chocolate
6	milk, chocolate

THE APRIORI ALGORITHM

Agrawal et al., 1996

- For (X, Y, Z) , a 3-item set, to be frequent (have enough support), (X, Y) , (X, Z) , and (Y, Z) should be frequent.
- If (X, Y) is not frequent, none of its supersets can be frequent.
- Once we find the frequent k -item sets, we convert them to rules:
 - ▶ $X, Y \rightarrow Z, \dots$
 - ▶ $X \rightarrow Y, Z, \dots$

THE APRIORI ALGORITHM

- Initially, scan DB once to get frequent 1-itemset
- Repeat
 - ▶ Generate length $(k+1)$ candidate itemsets from length k frequent itemsets
 - ▶ Test the candidates against DB to find frequent $(k+1)$ itemsets
 - ▶ Set $k = k+1$
- Terminate when no frequent or candidate set can be generated
- Return all the frequent itemsets derived

Convert the found itemsets into rules with enough confidence, by splitting the itemset into antecedent and consequent

For $X \rightarrow Y, Z$ to have enough confidence, $X, Y \rightarrow Z$ must have enough confidence

- Start by considering a single consequent and test the confidence for all possible single consequent
- Then consider two consequents for the added rules; etc.

SOURCES AND RESOURCES

- 'Pattern Recognition and Machine Learning', Bishop, 2006
- probabilitycourse.com
- <http://www.cs.cornell.edu/courses/cs4780/2018fa/lectures/index.html>
- <http://faculty.ucmerced.edu/mcarreira-perpinan/teaching/CSE176/lecturenotes.pdf>
- <https://slideplayer.com/slide/4778004/> and [/6275779/](https://slideplayer.com/slide/6275779/)