

5. Bayesian decision theory

Chloé-Agathe Azencott

Centre for Computational Biology, Mines ParisTech
chloe-agathe.azencott@mines-paristech.fr



Learning objectives

After this lecture, you should be able to

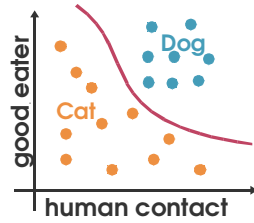
- **Apply Bayes rule** for simple inference and decision problems;
- Explain the connection between **Bayes decision rule**, **empirical risk minimization**, **maximum a priori** and **maximum likelihood**;
- Use a graph to express **conditional independence** among random variables;
- Apply the **Naive Bayes** algorithm.

Probability and inference

- Result of **tossing a coin**: x in {heads, tails}
 - $x = f(z)$ z : **unobserved variables**
 - Replace $f(z)$ (maybe deterministic but unknown) with the **random variable** X in $\{0, 1\}$ drawn from a **probability distribution** $P(X=x)$.
- We do not know P but a **sample** $X = \{x^i\}_{i=1, \dots, n}$
- Goal: **approximate P** (from which X is drawn)
- **Prediction** of next toss:

Classification

- Cat vs. dog
 - Cat = 1 (positive)
 - Dog = 0 (negative)
 - x_1 = human contact
 - x_2 = good eater



- **Prediction:**

$$\hat{y} = \begin{cases} 1 & \text{if} \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{y} = \begin{cases} 1 & \text{if} \\ 0 & \text{otherwise} \end{cases}$$

Bayes rule

Reverend Thomas Bayes

170?-1761



... possibly

Bayes rule

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

Example: rare disease testing

- test is correct 99% of the time
- disease prevalence = 1 out of 10,000

What is the probability that a patient that tested positive actually has the disease?

99% ? 90% ? 10% ? 1% ?

Bayes rule

$$P(y = c|\mathbf{x}) = \frac{\overset{\text{prior}}{P(y = c)}\overset{\text{likelihood}}{p(\mathbf{x}|y = c)}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

$$\begin{aligned} P(y = 0) + P(y = 1) &= 1 & p(\mathbf{x}) &= p(\mathbf{x}|y = 1)P(y = 1) + \\ P(y = 0|\mathbf{x}) + P(y = 1|\mathbf{x}) &= 1 & & p(\mathbf{x}|y = 0)P(y = 0) \end{aligned}$$

Bayes' decision rule:

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|\mathbf{x}) > P(y = 0|\mathbf{x}) \\ 0 & \text{otherwise.} \end{cases}$$

Maximum A Posteriori criterion

• **MAP decision rule:**

- pick the hypothesis that is most probable
- i.e. **maximize the posterior** $P(y = c|\mathbf{x}) = \frac{P(y = c)p(\mathbf{x}|y = c)}{p(\mathbf{x})}$

$$\Lambda_{\text{MAP}}(\mathbf{x}) = \frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})}$$

• **Decision rule:**

Likelihood ratio test (LRT)

$$\Lambda_{\text{MAP}}(\mathbf{x}) = \frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} \quad \Lambda_{\text{MAP}}(\mathbf{x}) >? 1 \quad P(y=c|\mathbf{x}) = \frac{P(y=c)p(\mathbf{x}|y=c)}{p(\mathbf{x})}$$

$$\Lambda_{\text{MAP}}(\mathbf{x}) = \frac{P(y=1)p(\mathbf{x}|y=1)p(\mathbf{x})}{P(y=0)p(\mathbf{x}|y=0)p(\mathbf{x})}$$

$p(\mathbf{x})$ does not affect the decision rule.

- Likelihood ratio test:**

test whether the **likelihood ratio** $\Lambda(\mathbf{x})$ is larger than

$$\Lambda(\mathbf{x}) = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)}$$

Example: LRT decision rule

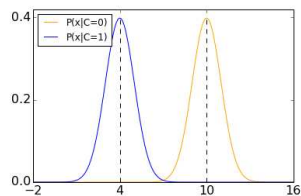
$$\Lambda(x) = \frac{p(x|y=1)}{p(x|y=0)} >? \frac{P(y=0)}{P(y=1)}$$

Assuming the likelihoods below and equal priors, derive a decision rule based on the LRT.

$$p(x|y=1) \sim \mathcal{N}(4, 1) \quad p(x|y=0) \sim \mathcal{N}(10, 1)$$

$Z \sim \mathcal{N}(\mu, \sigma^2)$:

$$f(z) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(z-\mu)^2/(2\sigma^2)}$$



Maximum likelihood criterion

- Consider **equal priors** $P(y=1) = P(y=0)$

$$\Lambda(\mathbf{x}) = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} >? \frac{P(y=0)}{P(y=1)} \quad \mathbf{1}$$

- Bayes decision rule seeks to maximize $P(\mathbf{x}|y=c)$ and is hence called the **Maximum Likelihood criterion**

$$\Lambda_{\text{ML}}(\mathbf{x}) = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)}$$

- **Decision rule:**

If $\Lambda_{\text{ML}}(\mathbf{x}) > 1$ then choose $y=1$ else choose $y=0$

Bayes rule for $K > 2$

- Bayes rule:**

$$P(y = c_k | \mathbf{x}) = \frac{p(\mathbf{x}|y = c_k)P(y = C_k)}{\sum_{l=1}^K p(\mathbf{x}|y = c_l)P(y = c_l)}$$

- $P(y = c_k) = P(y = c_1) + P(y = c_2) + \dots + P(y = c_K)$

- **What is the decision rule?**

Risk minimization

Losses and risks

- So far we've assumed all errors were **equally costly**.

But misclassifying a cancer sufferer as a healthy patient is much more problematic than the other way around.

- **Action α_k :** assigning class c_k
- **Loss:** quantify the cost λ_{kl} of taking action α_k when the true class is c_l

- **Expected risk:**

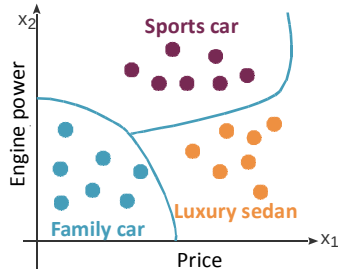
$$R(\alpha_k | \mathbf{x}) = \sum_{l=1}^K \lambda_{lk} P(y = c_l | \mathbf{x})$$

- **Decision (Bayes Classifier):** $\arg \min_k R(\alpha_k | \mathbf{x})$

Discriminant functions

Classification = find K **discriminant functions** f_k s.t. \mathbf{x} is assigned class c_k if $k = \operatorname{argmax}_k f_k(\mathbf{x})$

- Bayes classifier: $f_k(\mathbf{x}) = -R(\alpha_k|\mathbf{x})$
- Defines K **decision regions** $R_k = \{\mathbf{x} : f_k(\mathbf{x}) = \max_l f_l(\mathbf{x})\}$



0/1 Loss

- All misclassifications are **equally costly**.
- $\lambda_{kl} = 0$ if $k=l$ and 1 otherwise

$$\begin{aligned} R(\alpha_k|\mathbf{x}) &= \sum_{l=1}^K \lambda_{lk} P(y = C_l|\mathbf{x}) \\ &= \sum_{l \neq k} P(y = C_l|\mathbf{x}) \\ &= 1 - P(y = C_k|\mathbf{x}) \end{aligned}$$

- **Minimizing the risk:**
 - choose the most probable class (MAP)
 - this is equivalent to the Bayes decision rule.

$$\Lambda(\mathbf{x}) = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} >? \frac{(\lambda_{10} - \lambda_{00})P(y=0)}{(\lambda_{01} - \lambda_{11})P(y=1)}$$

Maximum likelihood criterion

- Consider **equal priors** $P(y=1) = P(y=0)$
- Consider the **0/1 loss function**

$$\Lambda(\mathbf{x}) = \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} >? \frac{(\lambda_{10} - \lambda_{00})P(y=0)}{(\lambda_{01} - \lambda_{11})P(y=1)}$$

Reject

- Add an artificial “reject” class (K+1) for **refusing to take a decision**.

E.g. Zip code detection.

- $\lambda_{kl} = \begin{cases} 0 & \text{if } k = l \\ \lambda & \text{if } k = K+1 \\ 1 & \text{otherwise} \end{cases}$

$$R(\alpha_k | \mathbf{x}) = \sum_{l \neq k} P(y = c_l | \mathbf{x}) = 1 - P(y = c_k | \mathbf{x})$$

$$R(\alpha_{K+1} | \mathbf{x}) = \sum_{l=1}^K \lambda P(y = c_l | \mathbf{x}) = \lambda$$

- **Decision:**

$$\hat{y} = c_k \text{ if } P(y = c_k | \mathbf{x}) > P(y = c_l | \mathbf{x}) \text{ for all } l \neq k \text{ and } P(y = c_k | \mathbf{x}) > 1 - \lambda$$

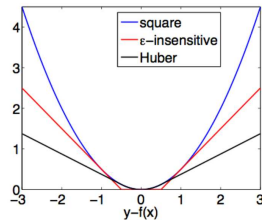
else reject.

Only meaningful if $0 < \lambda < 1$

Losses for regression

- **Square loss:** $L(f(\mathbf{x}), y) = (f(\mathbf{x}) - y)^2$
- **ϵ -insensitive loss:** $L(f(\mathbf{x}), y) = (|f(\mathbf{x}) - y| - \epsilon)_+$
- **Huber loss:** mix of linear and quadratic

$$L_\delta(f(\mathbf{x}), y) = \begin{cases} \frac{1}{2} (y - f(\mathbf{x}))^2 & \text{if } |y - f(\mathbf{x})| \leq \delta \\ \delta |y - f(\mathbf{x})| - \frac{1}{2} \delta^2 & \text{otherwise.} \end{cases}$$



Empirical risk minimization (ERM)

- **Loss:** $L(f(\mathbf{x}), y)$ small when $f(\mathbf{x})$ predicts y well
- **Expected risk:**

$$R = \mathbb{E}[L(f(\mathbf{x}), y)]$$

- **Empirical risk:**

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}^i), y^i)$$

- The **ERM estimator** of the functional class F is the solution, when it exists, of:

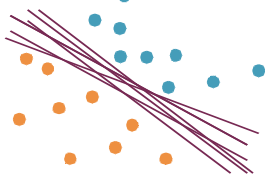
$$\hat{f}_n = \arg \min_{f \in F} R_n(f)$$

Solving ERM

- There can sometimes be an **explicit analytical solution**
- Otherwise: **convex optimization** (if the loss function is convex in f)
- **Limits of ERM:**
 - **ill-posed**
 - **not statistically consistent**
This is particularly true in **high dimension**.

ERM is ill-posed

- **Well-posed problems** (Hadamard):
 - Mathematical models of physical phenomena such that
 - a solution exists;
 - the solution is unique;
 - the solution's behavior changes continuously with the initial conditions.
- It can be that **an infinite number of solutions minimize the empirical risk** to zero.



ERM is not statistically consistent

- **Statistical consistency:** Estimator θ_N of θ that converges in probability towards θ as N increases.

$$\forall \epsilon > 0 \quad \lim_{N \rightarrow \infty} Pr(|\theta_N - \theta| \geq \epsilon) = 0$$

- From the **law of large numbers**,

$$\forall f \in \mathcal{F}, \quad R_N(f) \xrightarrow{N \rightarrow \infty} R(f)$$

but this isn't enough to guarantee that minimizing $R_N(f)$ gives a good estimator of the minimizer of $R(f)$.

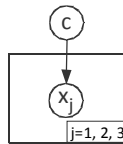
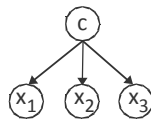
- Vapnik showed that this is only true if the capacity of hypothesis space \mathcal{F} is “not too large”.

Multivariate classification: Naive Bayes

- Multivariate classification: \mathbf{x} is multidimensional
- Assume the variables x_1, x_2, \dots, x_p are **conditionally independent**: $p(x_{j_1}|y = c, x_{j_2}) = p(x_{j_1}|y = c)$

Graphical representation

- We can use a graph to represent **conditional independence**:
 - arc from c to x_j means the distribution of X_j **depends** on c
 - no arc between X_{j_1} and X_{j_2} means that X_{j_1} and X_{j_2} are **independent given C**:
 $p(x_{j_1}|y = c, x_{j_2}) = p(x_{j_1}|y = c)$
- A **plate** represents repeated structure:
 - all X_j inside the same plate follow the same probability distribution.



Naive Bayes

- Multivariate classification: \mathbf{x} is multidimensional
- Assume the variables x_1, x_2, \dots, x_p are **conditionally independent**: $p(x_{j_1}|y = c, x_{j_2}) = p(x_{j_1}|y = c)$

$$P(y = c_k|\mathbf{x}) = \frac{p(\mathbf{x}|y = c_k)P(y = c_k)}{\sum_{l=1}^K p(\mathbf{x}|y = c_l)P(y = c_l)}$$

$$p(x_1, \dots, x_p|y = c) = p(x_1|y = c)p(x_2|y = c) \dots p(x_p|y = c)$$

Hence:

$$P(y = c|x_1, \dots, x_p) = \underbrace{\frac{1}{Z}}_{\text{scaling factor, independent of } c_k} P(y = c)p(x_1|y = c)p(x_2|y = c) \dots p(x_p|y = c)$$

Maximum a posteriori estimation

- **MAP decision rule:** pick the hypothesis that is most probable
- For Naive Bayes:

$$\hat{y} = \arg \max_{k=1, \dots, K} p(y = c_k) \prod_{i=1}^n p(x^i | y = c_k)$$

Naive Bayes spam filtering

- Input: email
- Output: spam / ham
- Naive Bayes assumption: conditional independence

bag of words

$$(x_1, x_2, \dots, x_p) = (0, 1, \dots, 0)$$

rich → 0
CLICK → 1
viagra → 0

Your Mail-Box has exceeded its storage Limit [CLICK=HERE](#) FILL and Click on FINISH for to get more space or you wont be able to send Mail

Dear Dr Azencott,

We obtained your contact information from your excellent papers, and would like to know if our company could serve you. Does your current work require the generation of custom monoclonal antibodies? If so, we would be glad to perform this tedious and time-consuming task on your behalf.

Dear Dr Azencotte,

Thank you very much for your review of manuscript CHIN-D-15-00031. We greatly appreciate your assistance.

Best wishes,
Samuel Winthrop
Journal of Cheminformatics

SPAM

NOT SPAM

- **P(spam | (x₁, x₂, ... x_p))**
= 1/Z p(spam) p(x₁ | spam) p(x₂ | spam) ... p(x_p | spam)
- **P(ham | (x₁, x₂, ... x_p))**
= 1/Z p(ham) p(x₁ | ham) p(x₂ | ham) ... p(x_p | ham)
- **Decision:**
If P(spam | (x₁, x₂, ..., x_p)) > P(ham | (x₁, x₂, ..., x_p)) then spam else ham
- **Inference:** we need to determine
p(spam), p(ham), p(x_j | spam), p(x_j | ham)

- **Bernoulli Naive Bayes:**

- Each email is the outcome of p Bernoulli trials
- **Naive assumption:** the trials are independent
word co-occurrences in a category aren't independent
still, independence assumptions can give good results

$$p(x_j | \text{spam}) = p_j^{x_j} (1 - p_j)^{(1-x_j)}$$

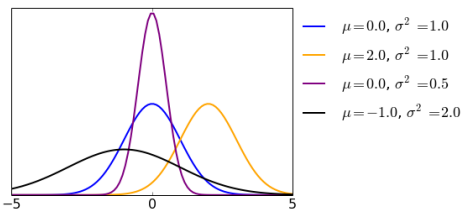
- **Direct estimate of p_j :** $p_j = S_j / S$
 - S = # spams in train set
 - S_j = # spams containing word j in train set

Gaussian naive Bayes

- Assume

$p(x_j | y=c_k)$ **univariate Gaussian**

$$p(x_j | y = c_k) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x_j - \mu)^2 / (2\sigma^2)}$$



Bayesian model selection

- **Priors on model:** $p(\text{model})$

$$p(\text{model} | \text{data}) = \frac{p(\text{data} | \text{model}) p(\text{model})}{p(\text{data})}$$

- Regularization \equiv prior that favors simpler models.

- Take the log

$$\log p(\text{model} | \text{data}) = \underbrace{\log p(\text{data} | \text{model})}_{\equiv \text{training error}} + \underbrace{\log p(\text{model})}_{\equiv \text{model complexity}} - c$$

- MAP similar to minimizing

$$E' = \text{empirical error} + \lambda \text{ model complexity}$$

Summary

$$P(y = c|\mathbf{x}) = \frac{\overset{\text{prior}}{P(y = c)} \overset{\text{likelihood}}{p(\mathbf{x}|y = c)}}{\underset{\text{evidence}}{p(\mathbf{x})}}$$

posterior

- **Bayes decision rule** \equiv **likelihood ratio test**
choose the most probable class, given evidence (data) and prior belief.
- Equivalent to **minimizing Bayes risk**
usually achieved approximately through **empirical risk minimization** (not equivalent!!)
- For the 0/1 loss, equivalent to **maximizing the posterior**.
- For the 0/1 loss and equal priors (uniform prior), equivalent to **maximizing the likelihood**.

References

- *A Course in Machine Learning*.
http://ciml.info/dl/v0_99/ciml-v0_99-all.pdf
 - **Bayes classifier**: Chap 2.1
 - **LRT**: Chap 9.4
 - **Naive Bayes**: Chap 9.3
- *The Elements of Statistical Learning*.
<http://web.stanford.edu/~hastie/ElemStatLearn/>
 - **Bayes classifier**: Chap 2.4
 - **Maximum Likelihood**: Chap 2.6.3, Chap 8.3
- **Probabilistic machine learning**
<https://www.repository.cam.ac.uk/bitstream/handle/1810/248538/Ghahramani%202015%20Nature>
- **Spam detection**: <http://www.paulgraham.com/spam.html>
- **Naive Bayes**:
<https://nlp.stanford.edu/IR-book/pdf/13bayes.pdf>