

# Bayesian Inference and naïve Bayes

Akito Sakurai

## Contents

- Bayes theorem
- MAP and ML
- Bayes optimal classifier and Gibbs algorithm
- Prediction of class or probability?
- Naïve Bayes

## Bayes Theorem

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



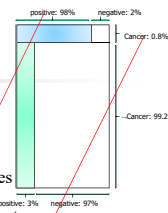
$$P(A, B) = P(A | B) P(B) = P(B | A) P(A)$$

## EX. (Mitchell Chap. 6.2)

Suppose we now observe a new patient for whom the lab test returns a positive result. Should we diagnose the patient as having cancer or not?

The test returns a correct positive result in only 98% of the cases in which the disease is actually present and a correct negative result in only 97% of the cases in which the disease is not present.

Over the entire population of people only .008 have this disease.



$$\begin{aligned}
 P(\text{cancer}) &= .008 & P(\neg \text{cancer}) &= .992 \\
 P(+ | \text{cancer}) &= .98 & P(- | \text{cancer}) &= .02 \\
 P(+ | \neg \text{cancer}) &= .03 & P(- | \neg \text{cancer}) &= .97 \\
 P(+) &= P(+ | \text{c}^r) P(\text{c}^r) + P(+ | \neg \text{c}^r) P(\neg \text{c}^r) = .0376 \\
 P(\text{cancer} | +) &= \frac{P(+ | \text{cancer}) P(\text{cancer})}{P(+)} = .209
 \end{aligned}$$

## EX. (Mitchell Exercise 6.1)

Suppose the doctor decides to order a second laboratory test for the same patient, and suppose the second test returns a positive result as well. What are the posterior probabilities of cancer and  $\neg$ cancer following these two tests? Assume that the two tests are independent.

$$\begin{aligned}
 P(\text{cancer}) &= .008 & P(\neg \text{cancer}) &= .992 \\
 P(+ | \text{cancer}) &= .98 & P(- | \text{cancer}) &= .02 \\
 P(+ | \neg \text{cancer}) &= .03 & P(- | \neg \text{cancer}) &= .97 \\
 P(+_1+_2) &= P(+_1+_2 | \text{c}^r) P(\text{c}^r) + P(+_1+_2 | \neg \text{c}^r) P(\neg \text{c}^r) = .00858 \\
 P(\text{cancer} | +_1+_2) &= \frac{P(+_1+_2 | \text{cancer}) P(\text{cancer})}{P(+_1+_2)} = .896
 \end{aligned}$$

## Basic Probability Formula

Multiplication rule (conditional probability):  
(Product rule)

$$P(A \wedge B) = P(A|B) P(B) = P(B|A) P(A)$$

Addition rule:  
(Sum rule)

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

Law of total probability:

$$P(B) = \sum_{i=1}^n P(B, A_i) = \sum_{i=1}^n P(B | A_i) P(A_i)$$

## Hypothesis selection

$$P(h | D) = \frac{P(D | h) P(h)}{P(D)}$$

$P(h)$  = prior probability of a hypotheses  $h$

$P(D)$  = prior probability that data  $D$  will be observed

$P(h|D)$  = probability that  $h$  holds given that  $D$  is observed

$P(D|h)$  = probability of observing data  $D$  given  $h$

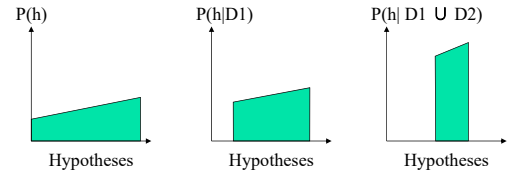
We can estimate the probability that  $h$  holds under the condition that the training data  $D$  is observed.

We can, then, estimate the probability that  $D$  is sampled under  $h$ .

Note: The conditional probability does not necessarily reflect a causal relationship, if any.

Note: Is it possible to think of the "probability that a hypothesis holds"

## Posterior probabilities



## Contents

- Bayes theorem
- MAP and ML
- Bayes optimal classifier and Gibbs algorithm
- Prediction of class or probability?
- Naïve Bayes

## MAP

$$P(h | D) = \frac{P(D | h) P(h)}{P(D)}$$

Finding the most probable hypothesis  $h \in H$  given the observed training data  $D$  should be most interesting.

*Maximum a posteriori hypothesis*  $h_{MAP}$ :

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h | D) \\ &= \arg \max_{h \in H} \frac{P(D | h) P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D | h) P(h) \end{aligned}$$

## ML

Suppose that  $P(h_i) = P(h_j)$  for any  $i, j$ , we get Maximum Likelihood (ML) hypothesis

$$h_{ML} = \arg \max_{h \in H} P(D | h)$$

Compare it with:

$$h_{MAP} = \arg \max_{h \in H} P(D | h) P(h)$$

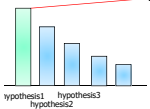
## An interpretation of ML

- In the real world, the prior distribution is thought to be unknown, incomputable, or non-existent.
  - E.g., does a prior distribution exist for words in documents? Doesn't it differ in age groups, social background, and others.
- If the existence of a prior distribution is questionable, likelihood-maximization is a reasonable choice.

## Most probable classification

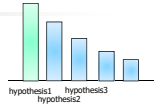
- So far, we have obtained the most probable **hypothesis** given  $D$  ( $h_{MAP}$ ).
- How about most probable **class** of a sample?

- $h_{MAP}(x)$  does not predict the most probable **class**.
  - What is the most probable classification of  $x$  ?
  - 3 hypotheses:  $P(h_1|D)=0.4$ ,  $P(h_2|D)=0.3$ ,  $P(h_3|D)=0.3$
  - Predictions for a sample:  $h_1(x)=+$ ,  $h_2(x)=-$ ,  $h_3(x)=-$



## Bayes optimal classification

$$\arg \max_{c_j \in \{+, -\}} \sum_{h_i \in H} P(c_j | h_i) P(h_i | D)$$



Note: Bayes optimal classifier must not be in  $H$ .

Note: Many papers report that it works well; but when we tried we often found no improvement in accuracy compared to MAP or ML. Why does this happen?

Note: Feasible? Doesn't it take time for computation?

## Ex. (Mitchell Chap. 6.7)

$$\begin{array}{lll} P(h_1 | D) = .4 & P(- | h_1) = 0 & P(+ | h_1) = 1 \\ P(h_2 | D) = .3 & P(- | h_2) = 1 & P(+ | h_2) = 0 \\ P(h_3 | D) = .3 & P(- | h_3) = 1 & P(+ | h_3) = 0 \end{array}$$

Therefore:

$$\sum_{h_i \in H} P(+ | h_i) P(h_i | D) = .4$$

$$\sum_{h_i \in H} P(- | h_i) P(h_i | D) = .6$$

And:

$$\arg \max_{c_j \in \{+, -\}} \sum_{h_i \in H} P(c_j | h_i) P(h_i | D) = -$$

## Contents

- Bayes theorem
- MAP and ML
- Bayes optimal classifier and Gibbs algorithm
- Prediction of class or probability?
- Naïve Bayes

## Bayes optimal classifier

- Suppose that  $D = \{x_1, \dots, x_n\}$  is observed from a distribution  $P(X; \theta)$  with the parameter  $\theta$ . We want to estimate  $y$  for an unseen  $x$  given  $D$ .
- Method 1: Estimate  $\theta$  and then predict by  $P(X; \theta)$ 
  - MLE (max. likelihood)  $\theta_{MLE} = \arg \max P(D | \theta)$
  - MAP (max. a posteriori)  $\theta_{MAP} = \arg \max P(D | \theta) P(\theta)$
  - Expectation (posterior mean)
 
$$\hat{\theta} = \int \theta P(\theta | D) d\theta = \int \theta P(D | \theta) P(\theta) / P(D) d\theta$$
- Method 2: without estimating the parameter  $\theta$ .
 
$$P(Y, \theta | D) = P(Y, D | \theta) P(\theta) / P(D)$$

$$\rightarrow P(Y | D) = \int P(Y, D | \theta) P(\theta) / P(D) d\theta$$

## Basic ideas of Bayesian inference

- Bayesian view is that we can measure uncertainty, even if there are not a lot of examples
  - What is the probability that a debut team will win the championship league this year?
    - Cannot do this with a frequentist approach
  - What is the probability that a newly minted particular coin will come up as heads?
    - Without much data we utilize our initial belief as the prior
- But as more data comes available we transfer more of our belief to the data (likelihood)
- With all the data, we do not consider the prior at all
- Belief is coded as a probability distribution

## An example: basic ideas

- Assume that we want to infer the mean  $\mu$  of a random variable  $x$  where the variance  $\sigma^2$  is known and we have not yet seen any data
- $P(\mu|D, \sigma^2) = P(D|\mu, \sigma^2)P(\mu)/P(D) \propto P(D|\mu, \sigma^2)P(\mu)$
- A Bayesian would want to represent the prior  $\mu_0$  and the likelihood  $\mu$  as parameterized distributions (e.g. Gaussian, multinomial, uniform, etc.)
- Let's assume a Gaussian just here
- Since the prior is a Gaussian we would like to multiply it by whatever the distribution of the likelihood is in order to get a posterior which is also a parameterized distribution specifically Gaussian

19

## Conjugate Priors

- $P(\mu|D, \sigma^2) = P(D|\mu)P(\mu)/P(D) \propto P(D|\mu)P(\mu)$
- If the posterior is the same distribution as the prior after the multiplication, then we say the prior and posterior are *conjugate* distributions and the prior is a conjugate prior for the likelihood
- In the case of a known variance and a Gaussian prior we can use a Gaussian likelihood and the product (posterior) will also be a Gaussian
- If the likelihood is multinomial then we would need to use a Dirchlet prior and the posterior would be a Dirchlet

20

## Discrete Conjugate Distributions

Likelihood	Model parameters	Conjugate prior distribution	Prior Hyperparameters	Interpretation of Hyperparameters <sup>Wikipedia</sup>	Posterior predictive mean <sup>Wikipedia</sup>
Bernoulli	$\mu$ (probability)	Beta	$\alpha, \beta$	$\alpha - 1$ successes, $\beta - 1$ failures <sup>Wikipedia</sup>	$p(\hat{\mu}) = \frac{\alpha^\alpha \beta^\beta}{(\alpha + \beta)^{\alpha + \beta}}$
Binomial	$\mu$ (probability)	Beta	$\alpha, \beta$	$\alpha - 1$ successes, $\beta - 1$ failures <sup>Wikipedia</sup>	BetaBinomial( $\beta, \alpha, \sigma^2$ ) (beta-binomial)
Negative Binomial with known failure number $r$	$\mu$ (probability)	Beta	$\alpha, \beta$	$\alpha - 1$ total successes, $\beta - 1$ failures <sup>Wikipedia</sup>	$\frac{\alpha^\alpha \beta^\beta}{(\alpha + \beta)^{\alpha + \beta}}$
Poisson	$\lambda$ (rate)	Gamma	$\beta, \theta$	$\beta$ total occurrences = $\beta$ (Poisson)	$\frac{\beta^\beta \theta^\beta}{(\beta + \theta)^{\beta + \theta}}$
Poisson	$\lambda$ (rate)	Gamma	$\alpha, \beta$	$\alpha$ total occurrences = $\beta$ (Poisson)	$\frac{\beta^\beta \theta^\beta}{(\beta + \theta)^{\beta + \theta}}$
Categorical	$\mu$ (probability vector, $k$ number of categories, $k$ size of $\mu$ )	Dirichlet	$\alpha$	$\alpha_i - 1$ occurrences of category $i$ <sup>Wikipedia</sup>	$p(\hat{\mu}) = \frac{\alpha^\alpha}{\sum \alpha_i}$
Multinomial	$\mu$ (probability vector, $k$ number of categories, $k$ size of $\mu$ )	Dirichlet	$\alpha$	$\alpha_i - 1$ occurrences of category $i$ <sup>Wikipedia</sup>	DirMult( $\alpha, \mu$ ) (Dirichlet-multinomial)
Hypergeometric with known total population size $N$	$\mu$ (number of target members)	Beta-binomial <sup>Wikipedia</sup>	$\alpha, \beta$	$\alpha - 1$ successes, $\beta - 1$ failures <sup>Wikipedia</sup>	$\frac{\alpha^\alpha \beta^\beta}{(\alpha + \beta)^{\alpha + \beta}}$
Dirichlet	$\mu$ (probability)	Beta	$\alpha, \beta$	$\alpha - 1$ successes, $\beta - 1$ failures <sup>Wikipedia</sup>	$\frac{\alpha^\alpha \beta^\beta}{(\alpha + \beta)^{\alpha + \beta}}$

From Wikipedia

## Continuous Conjugate Distribution (1)

Likelihood	Model parameters	Conjugate prior distribution	Prior Hyperparameters	Interpretation of Hyperparameters	Posterior predictive mean <sup>Wikipedia</sup>
Normal with known variance $\sigma^2$	$\mu$ (mean)	Normal	$\mu_0, \nu_0^2$	mean was estimated from observations with total precision (sum of all individual precisions) $1/\nu_0^2$ and with sample mean $\mu_0$	$\mathcal{N}(\hat{\mu} \mu_0, \nu_0^2 + \sigma^2)^{-1}$
Normal with known precision $\tau$	$\mu$ (mean)	Normal	$\mu_0, \tau_0$	mean was estimated from observations with total precision (sum of all individual precisions) $\tau_0$ and with sample mean $\mu_0$	$\mathcal{N}(\hat{\mu} \mu_0, \frac{1}{\tau_0 + \tau})$
Normal with known mean $\mu$	$\sigma^2$ (variance)	Inverse gamma	$\alpha, \beta$	variance was estimated from $2\alpha$ observations with sample variance $\beta/\alpha$ (i.e. with sum of squared deviations $2\beta$ ), where deviations are from known mean $\mu$	$\text{Inv-}\Gamma(\hat{\mu} \mu, \sigma^2 = \beta/\alpha)^{-1}$
Normal with known mean $\mu$	$\sigma^2$ (variance)	Inverse normal chi-squared	$\nu, \nu_0^2$	variance was estimated from $\nu$ observations with sample variance $\nu_0^2$	$\text{Inv-}\chi^2(\hat{\mu} \mu, \nu_0^2)^{-1}$
Normal with known mean $\mu$	$\sigma^2$ (variance)	Gamma	$\alpha, \beta$	precision was estimated from $2\alpha$ observations with sample variance $\beta/\alpha$ (i.e. with sum of squared deviations $2\beta$ ), where deviations are from known mean $\mu$	$\text{Inv-}\chi^2(\hat{\mu} \mu, \sigma^2 = \beta/\alpha)^{-1}$
Normal with $\mu$ and $\sigma^2$ Assuming exchangeability	Normal-inverse gamma	Normal-inverse gamma	$\mu_0, \nu_0, \alpha, \beta$	mean was estimated from $\nu$ observations with sample mean $\mu_0$ , variance was estimated from $2\alpha$ observations with sample mean $\mu_0$ and sum of squared deviations $2\beta$	$\text{Inv-}\chi^2(\hat{\mu} \mu_0, \frac{\sigma^2(\nu_0 + \beta)}{\nu_0 + \beta})$
Normal with $\mu$ and $\sigma^2$ Assuming exchangeability	Normal-inverse gamma	Normal-inverse gamma	$\mu_0, \nu_0, \alpha, \beta$	mean was estimated from $\nu$ observations with sample mean $\mu_0$ , and precision was estimated from $2\alpha$ observations with sample mean $\mu_0$ and sum of squared deviations $2\beta$	$\text{Inv-}\chi^2(\hat{\mu} \mu_0, \frac{\sigma^2(\nu_0 + \beta)}{\nu_0 + \beta})$

Wikipediaより

## Continuous Conjugate Distribution (2)

Multivariate normal with known covariance matrix $\Sigma$	$\mu$ (mean vector)	Multivariate normal	$\mu_0, \Sigma_0$	$(\Sigma_0^{-1} + n\Sigma^{-1})^{-1} (\Sigma_0^{-1}\mu_0 + n\Sigma^{-1}\bar{x})$ , $\Sigma_0^{-1} + n\Sigma^{-1}$	mean was estimated from observations with total precision (sum of all individual precisions) $\Sigma_0^{-1}$ and with sample mean $\bar{x}$	$\mathcal{N}(\hat{\mu} \mu_0, \Sigma_0^{-1} + \Sigma^{-1})^{-1}$
Multivariate normal with known precision matrix $\Psi$	$\mu$ (mean vector)	Multivariate normal	$\mu_0, \Psi_0$	$(\Psi_0 + n\Psi)^{-1} (\Psi_0\mu_0 + n\Psi\bar{x})$ , $(\Psi_0 + n\Psi)$	mean was estimated from observations with total precision (sum of all individual precisions) $\Psi_0$ and with sample mean $\bar{x}$	$\mathcal{N}(\hat{\mu} \mu_0, (\Psi_0^{-1} + \Psi^{-1})^{-1})$
Multivariate normal with known mean $\mu$	$\Sigma$ (covariance matrix)	Inverse-Wishart	$\nu, \Phi$	$\nu + n$ , $\Phi + \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$	covariance matrix was estimated from $\nu$ observations with sum of pairwise deviation products $\Phi$	$\text{Inv-}\Psi(\hat{\mu} \mu, \frac{1}{\nu + n} \Phi)$
Multivariate normal with known mean $\mu$	$\Sigma$ (covariance matrix)	Wishart	$\nu, V$	$\nu + n$ , $(V^{-1} + \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T)^{-1}$	covariance matrix was estimated from $\nu$ observations with sum of pairwise deviation products $V^{-1}$	$\text{Inv-}\Psi(\hat{\mu} \mu, \frac{1}{\nu + n} V)$
Multivariate normal with $\mu$ and $\Sigma$ (covariance matrix)	$\mu$ (mean vector) and $\Sigma$ (covariance matrix)	Normal-inverse-Wishart	$\mu_0, \nu_0, \Psi_0, \Phi_0$	$\frac{\nu_0 \mu_0 + n\bar{x}}{\nu_0 + n}$ , $\nu_0 + n$ , $\Phi_0 + \sum_{i=1}^n (x_i - \mu_0)(x_i - \mu_0)^T$ , $\Psi_0 + C + \sum_{i=1}^n (x_i - \mu_0)(x_i - \mu_0)^T$ , $C = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$	mean was estimated from $\nu_0$ observations with sample mean $\mu_0$ , covariance matrix was estimated from $\nu_0$ observations with sample mean $\mu_0$ and with sum of pairwise deviation products $\Phi_0 = \nu_0 \Sigma_0$	$\text{Inv-}\Psi(\hat{\mu} \mu_0, \frac{\nu_0 + 1}{\nu_0 + n + 1} \Psi_0)$
Multivariate normal with $\mu$ and $\Sigma$ (covariance matrix)	$\mu$ (mean vector) and $\Sigma$ (covariance matrix)	Normal-inverse-Wishart	$\mu_0, \nu_0, \Psi_0, V_0$	$\frac{\nu_0 \mu_0 + n\bar{x}}{\nu_0 + n}$ , $\nu_0 + n$ , $\Phi_0 + \sum_{i=1}^n (x_i - \mu_0)(x_i - \mu_0)^T$ , $V_0^{-1} + C + \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$ , $C = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$	mean was estimated from $\nu_0$ observations with sample mean $\mu_0$ , covariance matrix was estimated from $\nu_0$ observations with sample mean $\mu_0$ and with sum of pairwise deviation products $V_0^{-1}$	$\text{Inv-}\Psi(\hat{\mu} \mu_0, \frac{\nu_0 + 1}{\nu_0 + n + 1} V_0)$

From Wikipedia

## Continuous Conjugate Distribution (3)

Uniform	$U(a, b)$	Uniform	$a, b$	Base $(x_1, \dots, x_n, \mu_0)$ , $b + n$	$b$ observations with maximum value $b$	$\frac{b^\beta}{\beta!} \Gamma(\beta)$
Poisson with known maximum $\tau_n$	$\lambda$ (rate)	Gamma	$\alpha, \beta$	$\alpha + n$ , $\beta + \sum_{i=1}^n x_i$	$\alpha$ observations with sum $\beta$ of the sizes of magnitude of each observation (i.e. the logarithm of the size of each observation is the maximum $\tau_n$ )	$\frac{\beta^\beta}{\beta!} \Gamma(\beta)$
Dirichlet with known shape $\beta$	$\theta$ (rate)	Inverse gamma <sup>Wikipedia</sup>	$\alpha, \beta$	$\alpha + n$ , $\beta + \sum_{i=1}^n x_i^2$	$\alpha$ observations with sum $\beta$ of the 2th power of each observation	$\frac{\beta^\beta}{\beta!} \Gamma(\beta)$
Log-normal with known precision $\tau$	$\mu$ (mean)	Normal <sup>Wikipedia</sup>	$\mu_0, \tau_0$	$(\nu_0 \mu_0 + \sum_{i=1}^n \ln x_i) / (\nu_0 + n\tau)$ , $\nu_0 + n\tau$	"mean" was estimated from observations with total precision (sum of all individual precisions) $\tau_0$ and with sample mean $\mu_0$	$\frac{\beta^\beta}{\beta!} \Gamma(\beta)$
Log-normal with known mean $\mu$	$\sigma^2$ (variance)	Gamma <sup>Wikipedia</sup>	$\alpha, \beta$	$\alpha + n$ , $\beta + \sum_{i=1}^n (\ln x_i - \mu)^2$	variance was estimated from $2\alpha$ observations with sample variance $\beta/\alpha$ (i.e. with sum of squared log deviations $2\beta$ ), where deviations are from the log of the data points and the "mean"	$\frac{\beta^\beta}{\beta!} \Gamma(\beta)$
Exponential	$\lambda$ (rate)	Gamma	$\alpha, \beta$	$\alpha + n$ , $\beta + \sum_{i=1}^n x_i$	$\alpha$ observations with sum $\beta$	$\frac{\beta^\beta}{\beta!} \Gamma(\beta)$
Gamma with known shape $\alpha$	$\beta$ (rate)	Gamma	$\alpha_0, \beta_0$	$\alpha_0 + n$ , $\beta_0 + \sum_{i=1}^n x_i$	$\alpha_0$ observations with sum $\beta_0$	$\frac{\beta_0^{\beta_0}}{\beta_0!} \Gamma(\beta_0)$
Gamma with known shape $\alpha$	$\beta$ (rate)	Gamma	$\alpha_0, \beta_0$	$\alpha_0 + n$ , $\beta_0 + \sum_{i=1}^n x_i$	$\alpha_0$ observations with sum $\beta_0$	$\frac{\beta_0^{\beta_0}}{\beta_0!} \Gamma(\beta_0)$
Gamma with known rate $\beta$	$\alpha$ (shape)	Gamma	$\alpha_0, \beta_0$	$\alpha_0 + n$ , $\beta_0 + \sum_{i=1}^n x_i$	$\alpha_0$ observations with sum $\beta_0$	$\frac{\beta_0^{\beta_0}}{\beta_0!} \Gamma(\beta_0)$
Gamma with known rate $\beta$	$\alpha$ (shape)	Gamma	$\alpha_0, \beta_0$	$\alpha_0 + n$ , $\beta_0 + \sum_{i=1}^n x_i$	$\alpha_0$ observations with sum $\beta_0$	$\frac{\beta_0^{\beta_0}}{\beta_0!} \Gamma(\beta_0)$
Gamma with known rate $\beta$	$\alpha$ (shape)	Gamma	$\alpha_0, \beta_0$	$\alpha_0 + n$ , $\beta_0 + \sum_{i=1}^n x_i$	$\alpha_0$ observations with sum $\beta_0$	$\frac{\beta_0^{\beta_0}}{\beta_0!} \Gamma(\beta_0)$

From Wikipedia

## An illustrative example of Bayes inference

- Prior dist.:  $P(\mu) = \mathcal{N}(\mu | \mu_0, \sigma_0^2)$
- Posterior dis.:  $P(\mu | D) = \mathcal{N}(\mu | \mu_N, \sigma_N^2)$

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML}$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

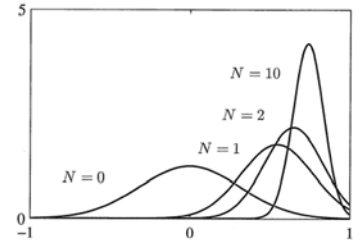
$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad \sigma^2 = \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

- What we believe moves from the prior distribution to data

25

## An illustrative example of Bayes inference

An illustration of Bayesian Inference for the location parameter  $\mu$  of a Gaussian distribution assuming the variance is given. The curves show the prior distribution of  $\mu$  (the curve labelled  $N = 0$ ) (this, too, is a Gaussian), along with the posterior distribution when increasing  $N$ . The data points are generated from a Gaussian of location and variance parameters 0.8 and 0.1 respectively, and the prior is set to have location parameter 0. Also in the prior distribution and the likelihood function, the true variance is known.



26

## An example of Bayes inference

- In this example, if the mean is known and the variance is unknown, the conjugate prior is the inverse-Gamma.
  - If precision (inverse of variance) is used, the conjugate prior is gamma distribution.
- If mean and variance is unknown, the conjugate prior is normal-inverse-gamma (a combination of normal and inverse-gamma distribution).
- A generalization of this for multivariate case is distribution for multiple dimensions is the normal-inverse-Wishart distribution.

27

## An example of Bayes inference

- $P(\mu, \sigma^2 | D) = P(D | \mu, \sigma^2) P(\mu, \sigma^2) / P(D)$   
 $\propto P(D | \mu, \sigma^2) P(\mu | \sigma^2) P(\sigma^2)$
  - prior:  $P(\mu | \sigma^2) = \mathcal{N}(\mu | \mu_0, \sigma^2 / k_0)$ ,  
 $P(\sigma^2) = \text{IG}(\sigma^2 | r_0/2, s_0/2)$   
 $N\_IG(\mu, \sigma^2 | \mu_0, k_0, r_0, s_0)$
  - posterior:  $P(\mu | \sigma^2, D) = \mathcal{N}(\mu | \mu_N, \sigma^2 / k_N)$ ,  
 $P(\sigma^2) = \text{IG}(\sigma^2 | r_N/2, s_N/2)$   
 $N\_IG(\mu, \sigma^2 | \mu_N, k_N, r_N, s_N)$
- $$\mu_N = \frac{k_0}{k_0 + N} \mu_0 + \frac{N}{k_0 + N} \mu_{ML} \quad r_N = r_0 + N$$
- $$k_N = k_0 + N \quad s_N = r_0 + (N - 1)$$

28

## Gibbs classifier

1. Select a hypothesis randomly according to  $P(h|D)$
2. Classify a new example following the h

Good news: If a hypothesis is randomly sampled from  $P(h)$ ,

$$E[\text{error}_{\text{Gibbs}}] \leq 2E[\text{error}_{\text{BayesOptimal}}]$$

(See "Mitchell Machine Learning Chap. 6.8")

Effective when there are so many hypothesis that a Bayes optimal is hard to calculate and we repeat the inferences

## Contents

- Bayes theorem
- MAP and ML
- Bayes optimal classifier and Gibbs algorithm
- Prediction of class or probability?
- Naïve Bayes

## Two types of target values

- **Categorical val.:** classification problem

- To divide (explanatory variable) space

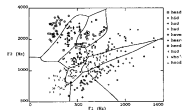
- Boundaries are to be obtained

- Category ↔ value range

- **Continuous val.:** regression

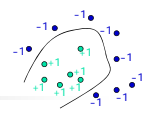
- Discrete values, too

When continuous func. is used, these are the same.  
 (1) e.g. 0 value-set, 1 value-set etc. are boundaries  
 (2) e.g. non-integer value area represents a category



Note: discrete function is not easily approximated by continuous one, so that regression framework is not straightforwardly applied to it

## Note



- For categorical values (assuming searching for continuous functions that assumes 0 on boundaries)

- Close to 0 = close to boundaries = not sure

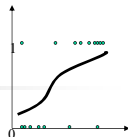
- Suppose that confidence level is represented by a real value between 0 and 1, it is, in the sub-area of a category,

- in the middle = confident = close to 1,
    - close to boundary = not sure = close to 0,

then the framework is of regression

- The values are ids of category

## Another viewpoint



- **Probability and the number of samples**

- If we use probability, not definite value, to represent the level of belongingness to a class, we could suppose that the frequency of samples reflects the probability

- **Confidence level and the # of samples**

- Confidence level of being in a subarea for a point is considered to be proportional to the number of samples around the point.

## In summary

Target value	Categorical value	Continuous value
Concept and method	Find boundary to minimize the # of errors	Regression of target values (minimize average errors)
	Regression of Categorical id values	
	Estimate distributions considering inputs as samples from a population (density estimation)	

## How to understand NN outputs

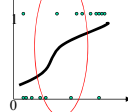
### Category id for classification

### Value for regression

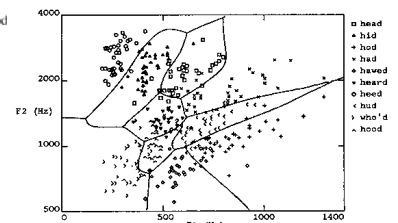
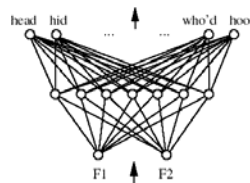
Although classification class is taught, output values is considered to be probability.

Although classification class is taught, output value is rounded to be considered as category id.

When a standard sigmoid function or the softmax is used as the final activation function, outputs are between 0 and 1.

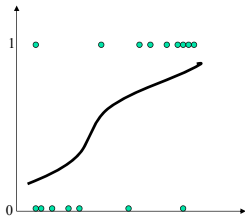


## Class areas and boundaries



## Learn to predict probability

Why do NNs learn probability, although just classes are taught?

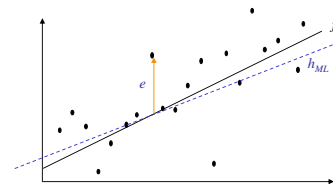


In online learning, frequency is learning by counting.

Is it similar to logistic regression?

## Regression to learn reals

By the way, what is regression?



## Statistical interpretation

Training samples:  $\langle x_i, d_i \rangle$  where

$$d_i = f(x_i) + e_i$$

$e_i$  is noise = prob. var. distributed to a normal distribution (iid) of mean=0 and finite deviation

iid=independent, identically distributed

Then:

$$h_{ML} = \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

## Statistical interpretation

$$\begin{aligned} h_{ML} &= \arg \max_{h \in H} \ln p(D | h) \\ &= \arg \max_{h \in H} \ln \prod_{i=1}^m e^{-\frac{1}{2} \left( \frac{d_i - h(x_i)}{\sigma} \right)^2} \\ &= \arg \max_{h \in H} \sum_{i=1}^m -\frac{1}{2} \left( \frac{d_i - h(x_i)}{\sigma} \right)^2 \\ &= \arg \max_{h \in H} \sum_{i=1}^m -(d_i - h(x_i))^2 \\ &= \arg \min_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2 \end{aligned}$$

## Squared error is not appropriate to predict probability

### Ex. Learning survival rate from data

$$\begin{aligned} h_{ML} &= \arg \max_{h \in H} \ln p(D | h) \\ &= \arg \max_{h \in H} \ln \prod_{i=1}^m P(d_i | h, x_i) P(x_i) && d_i \text{ is 0 or 1 (or probability belonging to a class)} \\ &= \arg \max_{h \in H} \sum_{i=1}^m \ln [P(d_i | h, x_i) P(x_i)] \\ &= \arg \max_{h \in H} \sum_{i=1}^m \ln (h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} P(x_i)) \\ &= \arg \max_{h \in H} \sum_{i=1}^m d_i \ln h(x_i) + (1 - d_i) \ln (1 - h(x_i)) \end{aligned}$$

note: cross entropy  $H(p, q) = -\sum_x p(x) \log q(x) = H(p) + D_{KL}(p \| q)$

## Contents

- Bayes theorem
- MAP and ML
- Bayes optimal classifier and Gibbs algorithm
- Prediction of class or probability?
- Naïve Bayes

## 目次

- Bayes 定理
- MAP と ML
- Bayes 最適分類器, Gibbs アルゴリズム
- クラスの推定か確率の推定か
- Naïve Bayes

知的情報処理の復習

## Naïve Bayes classifier

- Since (although?) simple, it is wellknown
  - More accurate than expected, although simple
  - Fast as is expected, since simple
- Bayes Theorem + Assumption *conditional independence*
  - The assumption hardly holds in the real world
  - In the real world, though, it works well
- Successful applications:
  - Text classification,
  - Diagnosis, and many others

Naïve Bayes is not a Bayesian

## Difficulty in Bayes classifier

- Recall that for a set of attributes  $\langle a_1, \dots, a_n \rangle$  of  $x$ , to infer the class that  $x$  belongs

$$\begin{aligned}c_{MAP} &= \arg \max_{c_j \in C} P(c_j | a_1, a_2, \dots, a_n) \\ &= \arg \max_{c_j \in C} \frac{P(a_1, a_2, \dots, a_n | c_j) P(c_j)}{P(a_1, a_2, \dots, a_n)} \\ &= \arg \max_{c_j \in C} P(a_1, a_2, \dots, a_n | c_j) P(c_j)\end{aligned}$$

- Difficulty: Huge dataset is required to infer  $P(a_1 \dots a_n | c_j)$ , since there are huge number of parameters ( $\prod |A_i|$ ) (for two value attributes,  $2^n$  parameters for  $n$  attributes)

## Naïve Bayes classifier

- **Naïve Bayes assumption:** attributes are mutually independent when the class is given
  - $P(a_1, \dots, a_n | c_j) = P(a_1 | c_j) P(a_2 | c_j) \dots P(a_n | c_j)$
  - *conditional independence* (given the class)
  - Reduces the number of parameters to infer:  
 $\prod |A_i| (=O(2^n)) \rightarrow \sum |A_i| (=O(n))$
- Under this assumption,  $c_{MAP}$  becomes

$$c_{NB} = \arg \max_{c_j \in C} P(c_j) \prod_i P(a_i | c_j)$$

## Naïve Bayes: an Algorithm

### Training (for a set of instances)

Estimate the probability that an instance  $x$  belongs to a class  $c_j$

$$\hat{P}(c_j) = P(c_j) \text{ 's estimator}$$

Estimate the probability that the  $i$ -th attribute value of an instance  $x$  belonging to the class  $c_j$  is  $a_i$ .

$$\hat{P}(a_i | c_j) = P(a_i | c_j) \text{ 's estimator}$$

Class( $x$ )

$$c_{NB} = \arg \max_{c_j \in C} \hat{P}(c_j) \prod_i \hat{P}(a_i | c_j)$$

## Naïve Bayes: Estimation

- How can we estimate  $P(c_j)$  and  $P(a_i | c_j)$  ?
  - A standard method that statistics tells us
    - Use frequency of the samples
    - $P(c)$  is estimated by  $\text{count}(c) / N$
    - $P(A|B)$  is estimated by  $\text{count}(A \wedge B) / \text{count}(B)$
  - Example: 100 samples. 70 + and 30 -
    - $P(+)=0.7$  and  $P(-)=0.3$
    - Among 70 positives, in 35  $a_1=\text{SUNNY}$
    - $P(a_1=\text{SUNNY} | +)=0.5$



## Example: Play Tennis

Outlook	Temp.	Humidity	Windy	Play
Sunny	Hot	High	No	No
Sunny	Hot	High	Yes	No
Overcast	Hot	High	No	Yes
Rainy	Mild	High	No	Yes
Rainy	Cool	Normal	No	Yes
Rainy	Cool	Normal	Yes	No
Overcast	Cool	Normal	Yes	Yes
Sunny	Mild	High	No	No
Sunny	Cool	Normal	No	Yes
Rainy	Mild	Normal	No	Yes
Sunny	Mild	Normal	Yes	Yes
Overcast	Mild	High	Yes	Yes
Overcast	Hot	Normal	No	Yes
Rainy	Mild	High	Yes	No

There are two classes: to play tennis (Play=Yes) and not to play tennis (Play=No)

Please infer if on the following day they played tennis or not

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

From Tom Mitchell's Machine Learning

49

## A solution

- For the *PlayTennis*, and a new instance <Outlook=sunny, Temp=cool, Humid=high, Windy=true>
- We want to calculate:

$$c_{NB} = \arg \max_{c_j \in C} \hat{P}(c_j) \prod_i \hat{P}(a_i | c_j)$$

- $\hat{P}(Y) \hat{P}(\text{sunny} | Y) \hat{P}(\text{cool} | Y) \hat{P}(\text{high} | Y) \hat{P}(\text{true} | Y) = 0.0053$
  - $\hat{P}(N) \hat{P}(\text{sunny} | N) \hat{P}(\text{cool} | N) \hat{P}(\text{high} | N) \hat{P}(\text{true} | N) = 0.0206$
- $$\Rightarrow c_{NB} = N$$

## Naïve Bayes: Conditional independence

- Is it necessary?
- What happens if the assumption does not hold?
  - i.e. if  $P(a_1, \dots, a_n | c_j) \neq P(a_1 | c_j) P(a_2 | c_j) \dots P(a_n | c_j)$
- If the following (weak) condition holds, the prediction is the same as Bayes classifier:

$$\arg \max_{c_j \in C} P(a_1 | c_j) P(a_2 | c_j) \dots P(a_n | c_j) P(c_j)$$

$$= \arg \max_{c_j \in C} P(a_1, a_2, \dots, a_n | c_j) P(c_j)$$

- But, the probability obtained in the prediction happens to be unrealistically close to 0 or 1

## Naïve Bayes: a Problem

- What happens if an attribute value  $a_i$  is not observed for a class  $c_j$ ?
  - Estimator of  $P(a_i | c_j) = 0$  because  $\text{count}(a_i \wedge c_j) = 0$
  - Big impacts: if this is 0, any products are 0!
- A solution: use Laplace correction.
  - $\hat{P}(a_i | c_j) = \frac{n_c + mp}{n + m}$
  - $n$  : # of training samples for  $c = c_j$
  - $n_c$  : # of training samples for  $c = c_j$  and  $a = a_i$
  - $p$  : prior probability (estimator)  $P^*(a_i | c_j)$  (uniform distribution is common)
  - $m$  : pseudo-count (commonly the number of attribute values)

$m=1$  is another choice which works better in many cases, too

## Note: Laplace correction

- (in parameter estimations from frequency) supposing a prior distribution for the parameter, obtain a MAP estimator.
- Beta distribution is the prior:
 
$$f(x; \alpha, \beta) = x^{\alpha-1} (1-x)^{\beta-1} / B(\alpha, \beta)$$
- The posterior mean of the parameter is the Laplace correction. If the likelihood is a result of a Bernoulli trial:  $\hat{\theta} = (n_0 + \alpha) / ((n_0 + n_1) + \alpha + \beta)$

## Note: smoothing

- In an estimation of statistical model, assigning a small probability to events that did not occur is called smoothing
  - [http://www.jaist.ac.jp/project/NLP\\_Portal/doc/glossary/index.html](http://www.jaist.ac.jp/project/NLP_Portal/doc/glossary/index.html)
- In natural language processing, frequencies of a word or a sequence of  $n$  words ( $n$ -gram) are often used. When  $n$  grows,  $n$ -gram becomes scarce, i.e., many  $n$ -grams do not occur. To solve the problem many techniques were invented.
  - Laplace smoothing (additive smoothing)
  - Linear interpolation
  - Good-Turing smoothing
  - Katz smoothing
  - Church-Gale smoothing
  - Witten-Bell smoothing
  - Kneser-Ney smoothing
  - ....
  - Hierarchical Pitman-Yor language model

54

## Text classification



- Texts classification:
  - Classifying documents (mail, news, web pages, etc. or a paragraph, a sentence, etc.)
  - Classifying e-mails into spam or not.
  - Classifying blogs into splog or not
  - Classifying news into interesting or not (to a person)
  - Classifying reviews of a product into groups of good reputation or not
  - Classifying reviews into trustable or not
  - Classifying open ended questions for questionnaire surveys
  - Classifying Q and A's at a call-center.
- Naïve Bayes works well
  - How to apply Naïve Bayes ?
  - Point: How to represent a sample (i.e. document), attributes?

## Document representation

### Bag-of-words

- Document as a vector of frequency of words in it
  - "Bag" implies discarding positions where the word occurs, and disregarding the sequences (contexts) of a word
  - i.e. if keio, gijuku, and university are words, there would be no difference between keio gijuku university, keio university giju, and gijuku keio university
- "what are words" is important, which should not differ among documents.
  - In English, "dog" and "dogs" should be treated as the same
- Ignore words not relevant to classification
  - In Japanese, particles such as ha, ga, mo, ya, etc are the ones
  - In English, prepositions
  - The words that have syntactic function but have no meaning are called functional words.
- Ignore words that are close to noise
  - Very low frequent words such as appearing just once.

## Document representation (cntd.)

- Representation itself is like naïve Bayes
  - Because representation is not inference, it is not naïve Bayes, but it really looks like naïve Bayes.
  - Probability of the occurrence of a document is formulated in naïve Bayes fashion.
  - Suppose that for each class of documents, the probability that a specific word occurs in a document is known as  $P(w_1 | c_j)$ ,  $P(w_2 | c_j)$ , ...,  $P(w_n | c_j)$ . If  $w_1, w_2, \dots, w_n$  are the words that occur in a document, then the probability that the document occurs is  $P(\text{doc} | c_j) = P(w_1 | c_j)^{TF(w_1)} P(w_2 | c_j)^{TF(w_2)} \dots P(w_n | c_j)^{TF(w_n)}$  where  $TF(w)$  is the term frequency of a word  $w$  in a document  $\text{doc}$

出現確率をこう書けば naïve Bayes といえよう

## Document classification by Naïve Bayes

- For a document  $\text{doc}$ ,

$$c_{NB} = \arg \max_{c_j \in C} P(c_j) \prod_{w_k \in \text{Voc}} P(w_k | c_j)^{TF(w_k, \text{doc})}$$

where  $TF(w_k, \text{doc})$  = frequency of  $w_k$  in  $\text{doc}$  and  $\text{Voc}$  is a set of all the words that we consider

- To represent word frequencies in a document, we need Laplace correction. The following estimator is used: where  $n_j$  = the number of words in a class  $c_j$ ,  $n_{k,j}$  = the number of occurrences of word  $w_k$  in class  $c_j$ .

$$P(w_k | c_j) = \frac{n_{k,j} + 1}{n_j + |\text{Voc}|}$$

## Twenty News Groups (Joachims 1996)

- 1000 training documents in each group
- Assign new documents to one of newsgroups

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x & rec.sport.hockey	rec.sport.hockey

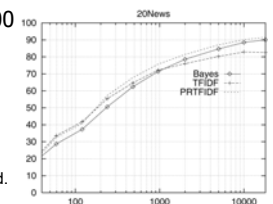
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

T. Joachims. A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization. In Proceedings of the 14th International Conference on Machine Learning, Nashville, TN, 1997, pp.143-151.

## Twenty News Groups (Joachims 1996)

- Naïve Bayes: 89% accuracy of classification
  - Highly frequent 100 word (the and of ...) are deleted
    - The words such as functional words, words relatively useless for classification are categorized as stop words and are deleted
  - The words occurring less than 3 times are deleted
  - The words remained: 38,500

Note: the accuracy is overly high.. In every text in 20 Newsgroups has a "subject" field which is very helpful for classification. Although the subject field is now deleted, in the previous works the field might be utilized.



Accuracy vs. # of training data (1/3 is reserved for test)

## 20 Newsgroups in R

- Naive Bayes package in R is not appropriate to large dataset.
  - Because the data matrix with naive implementation becomes huge (in the previous R program, the matrices  $xy$ ,  $xy$ ,  $tt$ ) (2000 rows for documents and 40,000 columns for words).
  - But non-zero entries are small in number, sparse matrix representation is useful.
  - You have to pay for computational overhead.
  - Then let us write an efficient program by ourselves!?

## 20 Newsgroups: データ

- In the "20 Newsgroups" site:
  - <http://people.csail.mit.edu/jrennie/20Newsgroups/>
  - Redirected to <http://qwone.com/~jason/20Newsgroups/>
- There is a preprocessed version:
  - 20news-bydate-matlab.tgz
- We use `train.data`, `train.label`, `test.data`, and `test.label`
- A program is uploaded to the class web page.
- Only the confusion matrix is in the next slide.
- Accuracy is around 78.2%.

```
> cm
      correct
predicted 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
1  237  3  3  0  0  0  0  1  0  4  2  0  2 10  3  7  2 12  7 47
2  0 299 33  8  8 42  9  1  1  1  0  5 18  7  8  2  0  1  1  3
3  0  7 208 15 10  8  4  0  0  0  0  1  0  1  0  1  0  0  0  0  0
4  0 12  58 306 38 10 49  2  0  1  0 1 28  3  0  0  0  0  0  0  0
5  0  7 11 21 275  2 21  0  0  1  0  2  8  0  0  1  1  0  0  0  0
6  1 21 30  2  3 306  1  1  0  2  0  1  3  0  2  2  0  0  1  0
7  0  1  0  4  4 1227  5  1  3  1  1  1  1  0  0  2  0  0  0  0
8  0  3  2  6  4  0 32 356 25  3  1  0  9  3  0  0  2  2  1  0
9  0  1  2  0  1  2  5  4 353  1  0  0  2  0  1  0  1  1  0  1
10 0  0  2  0  1  1  0  2  2 345  4  0  0  2  0  0  1  1  0  0
11 1  0  1  1  0  0  1  0  0 16 381  0  0  0  1  0  0  1  0  0
12 1 16 17  5  5 10  3  1  1  2 1 361 45  0  3  1  3  4  3  1
13 1  4  1 23 16  0 11  4  1  2  0  3 260  3  4  0  0  0  0  0
14 2  3  4  0  7  0  2  0  1  0  2  2  6 324  4  1  1  0  3  3
15 3  6  4  1  2  3  3  2  0  0  1  0  3  3 333  0  2  0  7  5
16 43  4  5  0  0  1  3  0  1  3  2  2  6 16  5 377  3  7  2  69
17 3  0  0  0  3  1  1  5  4  1  0  7  0  3  1  2 324  3  95 19
18  9  0  0  0  0  1  3  1  2  2  1  0  2  6  2  2  2 323  5  5
19  7  2  9  0  6  2  6  9  5  9  3  8  0 10 24  1 16 21 184  8
20 10  0  1  0  0  0  1  1  0  1  0  1  0  1  1  1  4  0  1  90
```

## Bayes inference and NB

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- Overview of the learning algorithm:
  - ML: maximize  $P(D|h)$
  - MAP: maximize  $P(h|D) \propto P(D|h) P(h)$
  - Posterior mean:
    - Bayes optimal classifier:  $P(c|D) = \int P(c|h)P(h|D) dh$
    - Hypotheses distribute!
- Regression under Gaussian noise:
  - $\Leftrightarrow$  minimization of mean squared error
- Learning of probability of binary events
  - $\Leftrightarrow$  minimization of cross-entropy
- Naive Bayes: rough assumption but practical
  - Ex. Document classification