

Review: 4-14, 4-19

- Types of machine learning problems
- Regularized Logistic Regression
- Naive Bayes Classifier
- Implementing a Gaussian Naives Bayes

- Application of probability, statistics, and prediction for measuring county mortality rates from Twitter.

Gaussian Naive Bayes

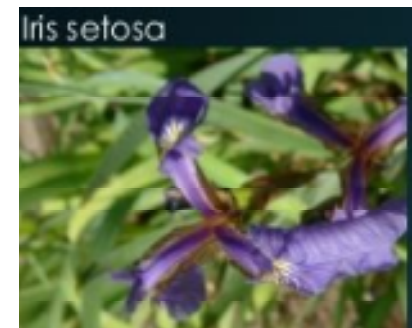
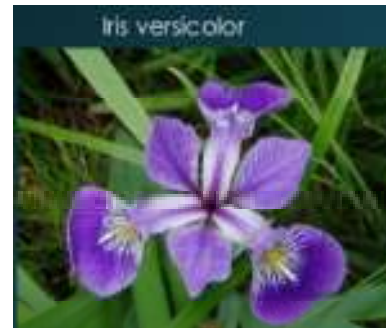
Assume $P(X|Y)$ is *Normal*

Then, training is:

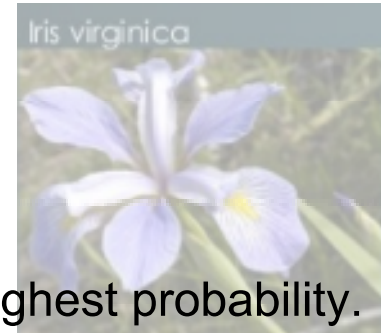
1. Estimate $P(Y = k)$; $\pi_k = \text{count}(Y = k) / \text{Count}(Y = *)$
2. MLE to find parameters (μ, σ) for each class of Y .
(the “class conditional distribution”)

Maximum a Posteriori (MAP): Pick the class with the maximum posterior probability.

$$\hat{y} = \mathit{arg\,max}_y P(y) \prod_{i=1}^m P(X_i|y)$$



Gaussian Naive Bayes



MLE: For which parameters does the observed data have the highest probability.

$$L(\theta) = \prod_{i=1}^n f(X_i; \theta)$$

$$l(\theta) = \log \sum_{i=1}^n f(X_i; \theta)$$



Maximum a Posteriori (MAP): Pick the class with the maximum posterior probability.

Unnormalized Posterior

$$\hat{y} = \arg \max_y P(y) \prod_{i=1}^m P(X_i|y)$$

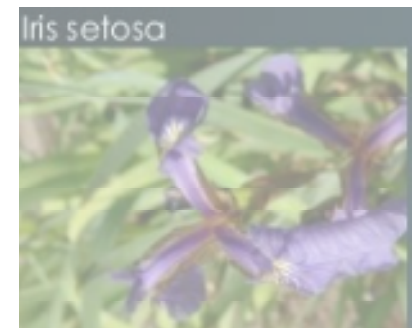
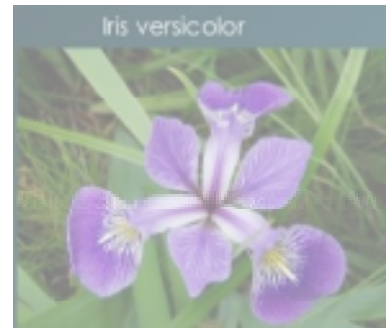
Gaussian Naive Bayes

Assume $P(X|Y)$ is *Normal*

Then, training is:

1. Estimate $P(Y = k)$; $\pi_k = \text{count}(Y = k) / \text{Count}(Y = *)$

Maximum a Posteriori (MAP): Pick the class with the maximum posterior probability.
(the class with the "normal distribution")



Without knowing $P(X)$,
can we turn this into the
(normalized) posterior?

← **Unnormalized Posterior**

$$\hat{y} = \mathop{\text{arg max}}_y P(y) \prod_{i=1}^m P(X_i|y)$$

Gaussian Naive Bayes

Use the **Law of Total Probability**, for all $i = 1 \dots k$, where $A_1 \dots A_k$ partition Ω :

Maximum a Posteriori (MAP): Pick the class with the maximum posterior probability.

Without knowing $P(X)$,
can we turn this into the
(normalized) posterior?

← **Unnormalized Posterior**

$$\hat{y} = \mathop{\text{arg max}}_y P(y) \prod_{i=1}^m P(X_i|y)$$

Gaussian Naive Bayes

Use the **Law of Total Probability**, for all $i = 1 \dots k$, where $A_1 \dots A_k$ partition Ω :

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

Maximum a Posteriori (MAP): Pick the class with the maximum posterior probability.

Without knowing $P(X)$,
can we turn this into the
(normalized) posterior?

← **Unnormalized Posterior**

$$\hat{y} = \arg \max_y P(y) \prod_{i=1}^m P(X_i|y)$$

Gaussian Naive Bayesian Inference

Use the **Law of Total Probability**, for all $i = 1 \dots k$, where $A_1 \dots A_k$ **partition** Ω :

$$P(A_i|B) = \frac{P(B, A_i)}{P(B)} = \frac{P(B|A_i)P(A_i)}{\sum_{i=1}^k P(B|A_i)P(A_i)} \quad \text{discrete}$$

$$P(A|B) = \frac{P(B|A)P(A)}{\int P(B|A)P(A)dA} \quad \begin{array}{l} \text{continuous} \\ A \text{ is} \\ \text{"marginalized"} \\ \text{out} \end{array}$$

Without knowing $P(X)$,
can we turn this into the
(normalized) posterior?

← **Unnormalized Posterior**

$$\hat{y} = \mathop{\text{arg max}}_y P(y) \prod_{i=1}^m P(X_i|y)$$

Gaussian Naive Bayesian Inference

Q: What distinguishes Bayesian inference? **A:** Assume a

$P(\theta)$ – prior

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

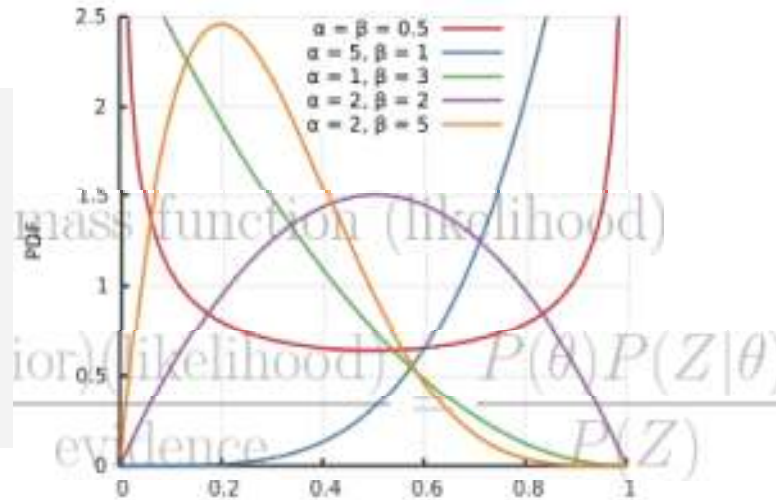
Types of priors:

- Uninformative (Improper: not a probability (e.g. constant))
- Belief-based
- **Conjugate** to a likelihood: if the posterior is in the same family as the prior.

Bayesian Inference

Example: Beta(α, β) is conjugate to a Bernoulli likelihood.

https://en.wikipedia.org/wiki/Conjugate_prior#Table_of_conjugate_distributions



Types of priors:

- Uninformative (Improper: not a probability (e.g. constant))
- Belief-based
- **Conjugate** to a likelihood: if the posterior is in the same family as the prior.

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

$$P(\theta|Z) = \frac{P(\theta)P(Z|\theta)}{\int P(\theta)P(Z|\theta)d\theta}$$

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

$$P(\theta|Z) = \frac{P(\theta)P(Z|\theta)}{\int P(\theta)P(Z|\theta)d\theta}$$

$$P(z^{new}|Z) = \int P(z^{new}|\theta)P(\theta|Z)d\theta \text{ -- predictive distribution}$$

Bayesian Inference

$$Z = X_{training}$$

Given:

$P(Z|\theta)$ – probability density or mass function (likelihood)

$P(\theta)$ – prior

Goal: Compute the posterior = $\frac{(\text{prior})(\text{likelihood})}{\text{evidence}} = \frac{P(\theta)P(Z|\theta)}{P(Z)}$

$$P(\theta|Z) = \frac{P(\theta)P(Z|\theta)}{\int P(\theta)P(Z|\theta)d\theta}$$

Like a posterior-weighted average of $P(z^{new}|\theta)$

$$P(z^{new}|Z) = \int P(z^{new}|\theta)P(\theta|Z)d\theta \quad \text{-- predictive distribution}$$