# Classification using probability models

## A classification problem

You have a bottle of wine whose label is missing.



Which winery is it from, 1, 2, or 3?

Solve this problem using visual and chemical features of the wine.

# The data set

**Training set** obtained from 130 bottles
- Winery 1: 43 bottles
- Winery 2: 51 bottles
- Winery 3: 36 bottles
- For each bottle, 13 features:
  'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash', 'Magnesium',
  'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins',
  'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'
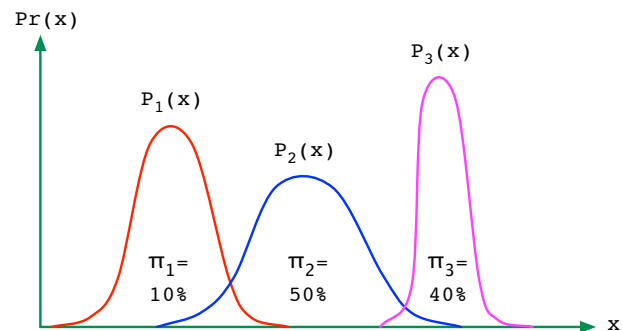
Also, a separate **test set** of 48 labeled points.

# The generative approach to classification

Example:
Data space $\mathcal{X} = \mathbb{R}^d$
Classes/labels $\mathcal{Y} = \{1, 2, 3\}$



The **learning process**:
- Use training data to fit a probability distribution to each class individually.
- For each class $j$, we have:
  - the probability of that class, $\pi_j = \mathrm{Pr}(y = j)$
  - the distribution of data in that class, $P_j(x)$

Overall **joint distribution**: $\mathrm{Pr}(x, y) = \mathrm{Pr}(y)\mathrm{Pr}(x|y) = \pi_y P_y(x)$.
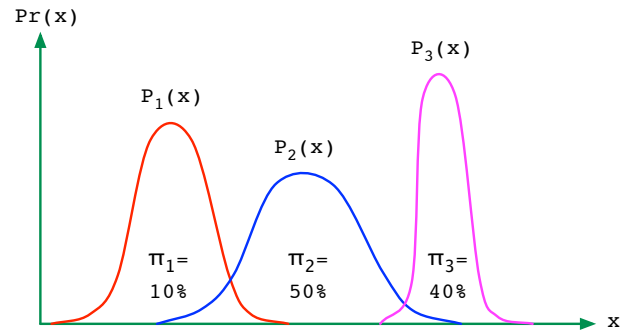
# Classifying a new point

Data space $\mathcal{X} = \mathbb{R}^d$
Classes/labels $\mathcal{Y} = \{1, 2, 3\}$

For each class $j$, we have
- $\pi_j = \Pr(y = j)$
- $P_j(x) = \Pr(x|y = j)$



Now we want to label a new point $x$. For any candidate label $j$,

$$\Pr(y = j|x) = \frac{\Pr(y = j)\Pr(x|y = j)}{\Pr(x)} = \frac{\pi_j P_j(x)}{\Pr(x)}$$

Optimal prediction: the class $j$ with largest $\pi_j P_j(x)$.

# Fitting a generative model

Training set of 130 bottles:
- Winery 1: 43 bottles, winery 2: 51 bottles, winery 3: 36 bottles
- For each bottle, 13 features: 'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash','Magnesium', 'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins', 'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'
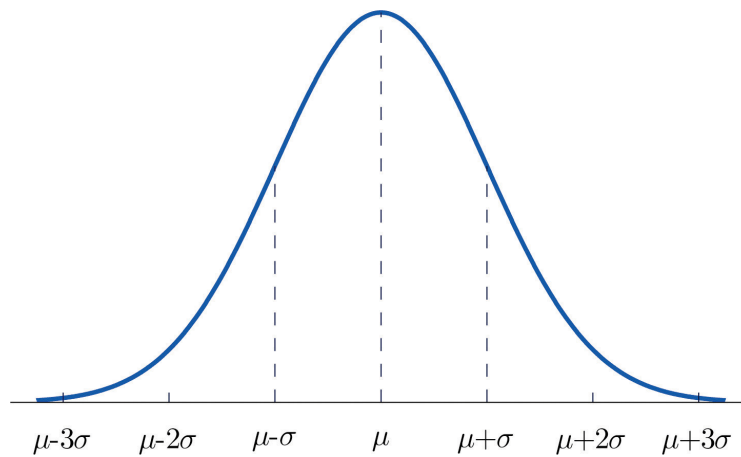
Class weights:
$\pi_1 = 43/130 = 0.33, \quad \pi_2 = 51/130 = 0.39, \quad \pi_3 = 36/130 = 0.28$

Need distributions $P_1, P_2, P_3$, one per class.
Base these on a single feature: 'Alcohol'.
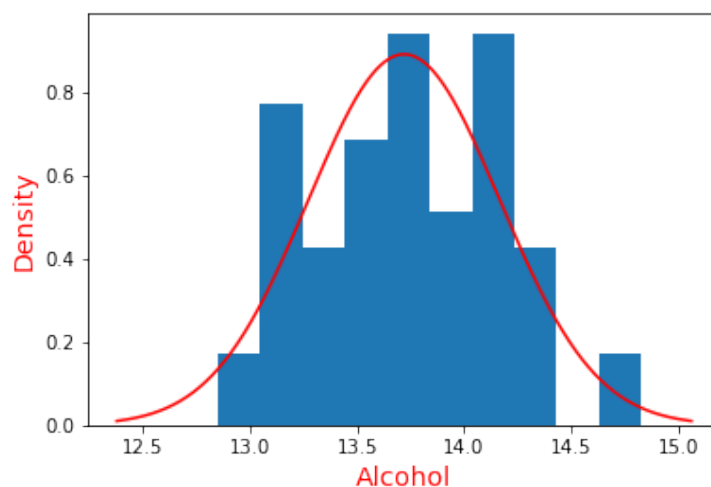
# Recall: Univariate Gaussian



The Gaussian $N(\mu, \sigma^2)$ has mean $\mu$, variance $\sigma^2$, and density function

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$
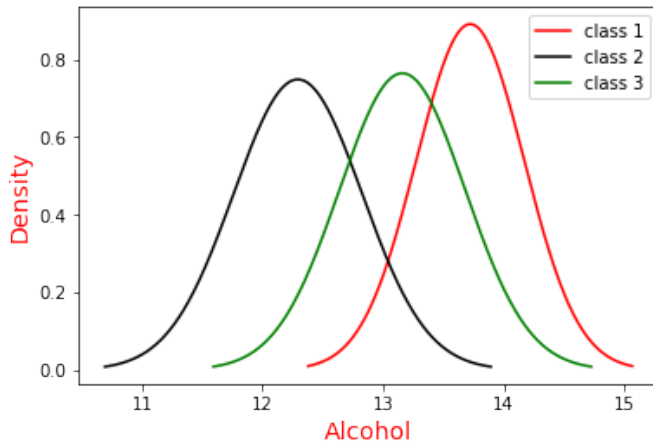
# The distribution for winery 1

Single feature: 'Alcohol'



Mean $\mu = 13.72$, Standard deviation $\sigma = 0.44$ (variance 0.20)

# All three wineries



- $\pi_1 = 0.33$, $P_1 = N(13.7, 0.20)$
- $\pi_2 = 0.39$, $P_2 = N(12.3, 0.28)$
- $\pi_3 = 0.28$, $P_3 = N(13.2, 0.27)$

To classify $x$: Pick the $j$ with highest $\pi_j P_j(x)$

**Test error: $14/48 = 29\%$**

# The data set, again

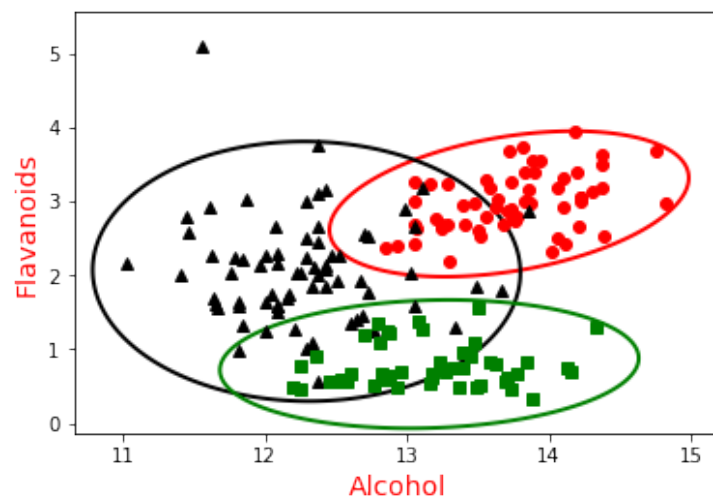Training set obtained from 130 bottles
- Winery 1: 43 bottles
- Winery 2: 51 bottles
- Winery 3: 36 bottles
- For each bottle, 13 features:
  'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash','Magnesium',
  'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins',
  'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'

Also, a separate test set of 48 labeled points.
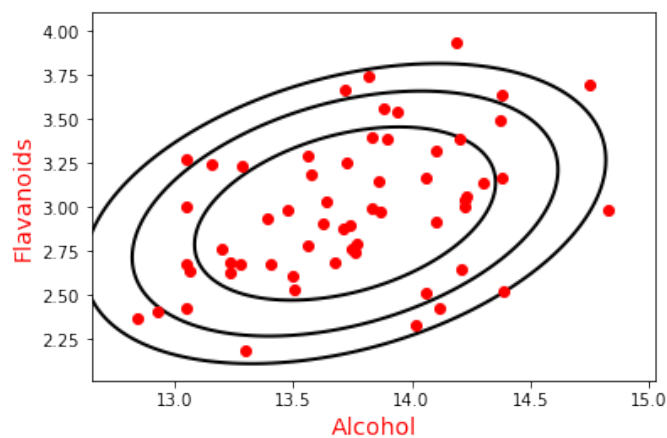
This time: 'Alcohol' and 'Flavanoids'.

# Why it helps to add features

Better **separation** between the classes!



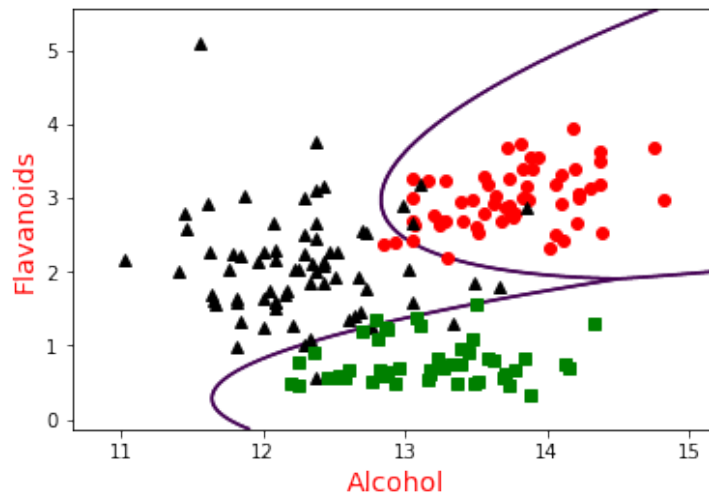Error rate drops from 29% to 8%.

# Bivariate Gaussian



Model class 1 by a bivariate Gaussian, parametrized by:

$$\text{mean } \mu = \begin{pmatrix} 13.7 \\ 3.0 \end{pmatrix} \text{ and covariance matrix } \Sigma = \begin{pmatrix} 0.20 & 0.06 \\ 0.06 & 0.12 \end{pmatrix}$$
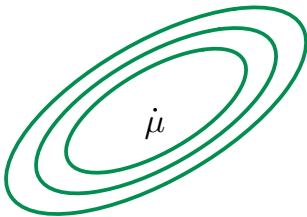
# The decision boundary

Go from 1 to 2 features: error rate goes from 29% to 8%.



What kind of function is this? And, can we use more features?

# Recall: Multivariate Gaussian



$N(\mu, \Sigma)$: Gaussian in $\mathbb{R}^d$
- mean: $\mu \in \mathbb{R}^d$
- covariance: $d \times d$ matrix $\Sigma$

Generates points $X = (X_1, X_2, \ldots, X_d)$.

- $\mu$ is the vector of coordinate-wise means:

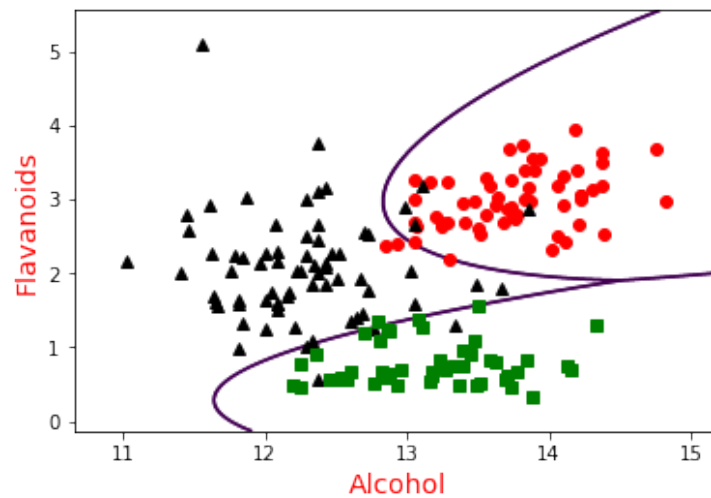$$\mu_1 = \mathbb{E}X_1, \ \mu_2 = \mathbb{E}X_2, \ldots, \ \mu_d = \mathbb{E}X_d.$$

- $\Sigma$ is a matrix containing all pairwise covariances:

$$\Sigma_{ij} = \Sigma_{ji} = \text{cov}(X_i, X_j) \quad \text{if } i \neq j$$
$$\Sigma_{ii} = \text{var}(X_i)$$

Density $\quad p(x) = \dfrac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\dfrac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$

# Back to the winery data

Go from 1 to 2 features: test error goes from 29% to 8%.



**With all 13 features: test error rate goes to zero.**

# Binary classification with Gaussian generative model

- Estimate class probabilities $\pi_1, \pi_2$
- Fit a Gaussian to each class: $P_1 = N(\mu_1, \Sigma_1), \ P_2 = N(\mu_2, \Sigma_2)$

Given a new point $x$, predict class 1 if

$$\pi_1 P_1(x) > \pi_2 P_2(x) \quad \Leftrightarrow \quad x^T M x + 2w^T x \geq \theta,$$

where:

$$M = \frac{1}{2}(\Sigma_2^{-1} - \Sigma_1^{-1})$$
$$w = \Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2$$

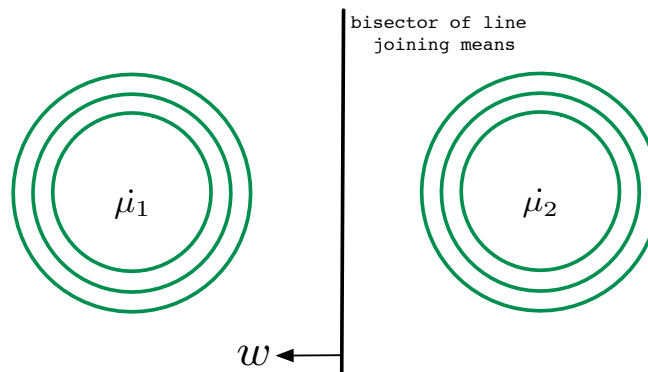and $\theta$ is a threshold depending on the various parameters.

**Linear** or **quadratic** decision boundary.

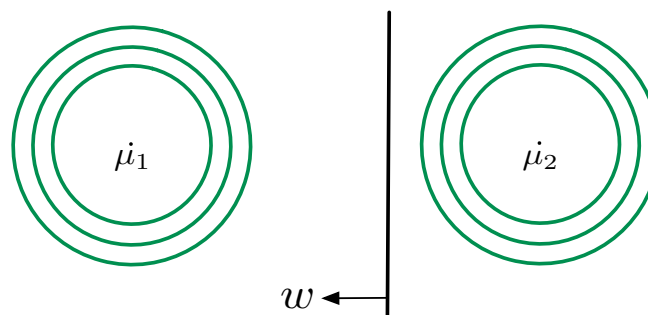# Common covariance: $\Sigma_1 = \Sigma_2 = \Sigma$

Linear decision boundary: choose class 1 if

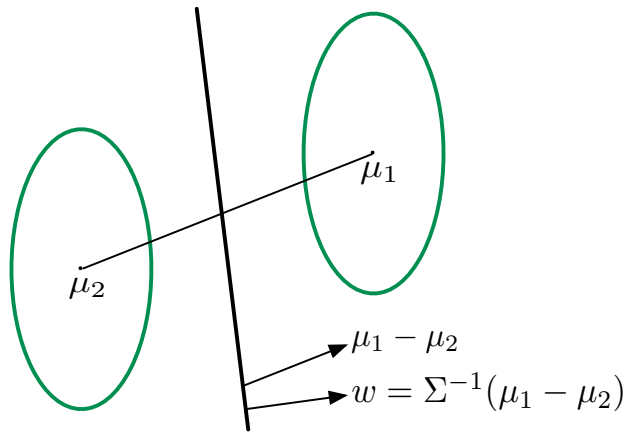$$x \cdot \underbrace{\Sigma^{-1}(\mu_1 - \mu_2)}_{w} \geq \theta.$$

Example 1: Spherical Gaussians with $\Sigma = I_d$ and $\pi_1 = \pi_2$.



Example 2: Again spherical, but now $\pi_1 > \pi_2$.

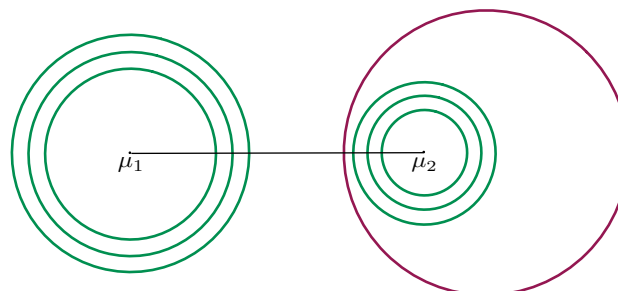Classification rule: $w \cdot x \geq \theta$

- Choose $w$ as above
- Common practice: fit $\theta$ to minimize training or validation error
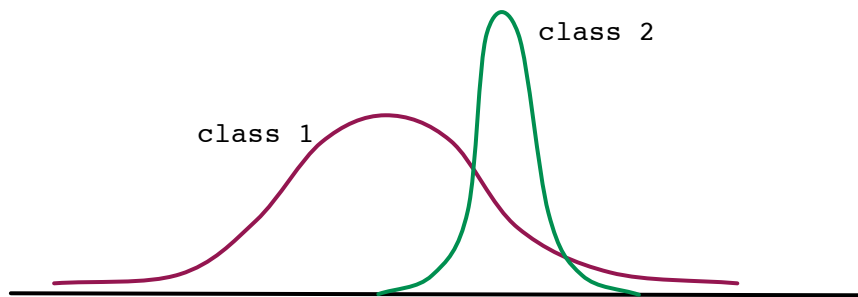
# Different covariances: $\Sigma_1 \neq \Sigma_2$

Quadratic boundary: choose class 1 if $x^T M x + 2 w^T x \geq \theta$, where:

$$M = \frac{1}{2}(\Sigma_2^{-1} - \Sigma_1^{-1})$$
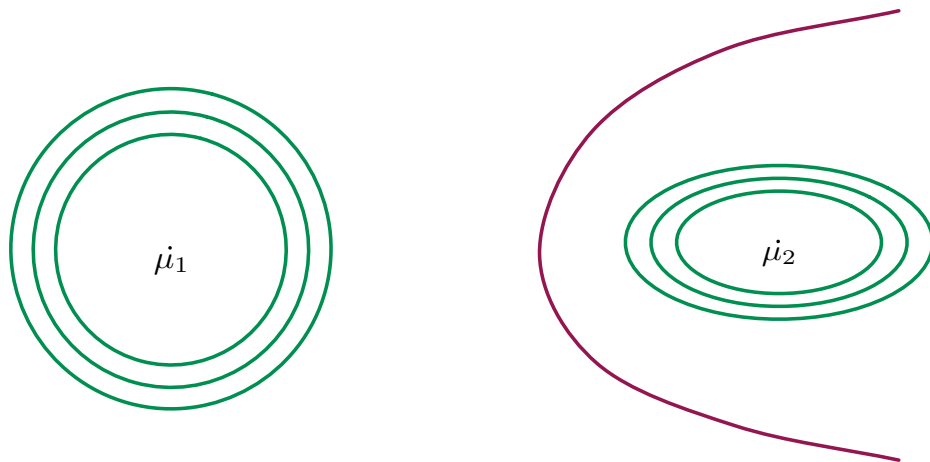$$w = \Sigma_1^{-1}\mu_1 - \Sigma_2^{-1}\mu_2$$

Example 1: $\Sigma_1 = \sigma_1^2 I_d$ and $\Sigma_2 = \sigma_2^2 I_d$ with $\sigma_1 > \sigma_2$

Example 2: Same thing in 1-d. $\mathcal{X} = \mathbb{R}$.

class 2

class 1

Example 3: A parabolic boundary.

$\mu_1$

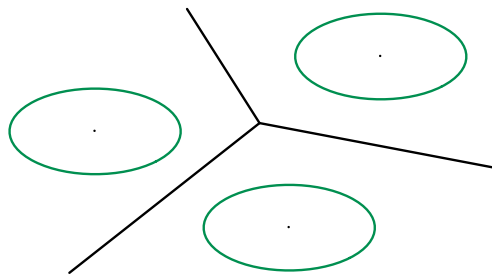$\mu_2$

# Multiclass discriminant analysis

$k$ classes: weights $\pi_j$, class-conditional densities $P_j = N(\mu_j, \Sigma_j)$.

Each class has an associated **quadratic** function

$$f_j(x) = \log\left(\pi_j\, P_j(x)\right)$$

To classify point $x$, pick $\arg\max_j f_j(x)$.

If $\Sigma_1 = \cdots = \Sigma_k$, the boundaries are **linear**.
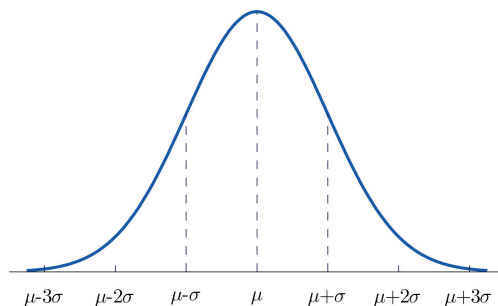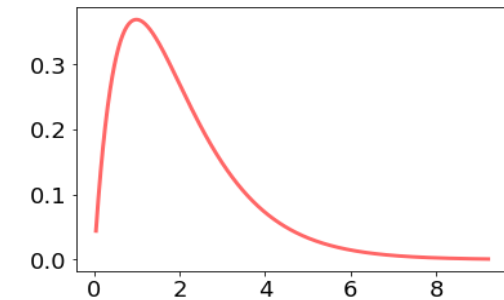


# Beyond Gaussians

The generative methodology:
- Fit a **distribution** to each class separately
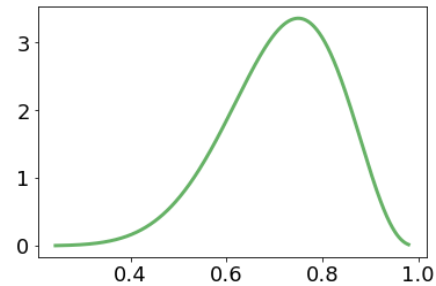- Use Bayes' rule to classify new data

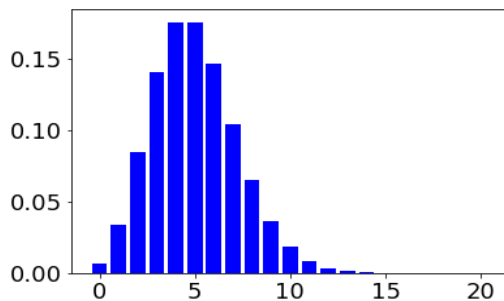What distribution to use? Are Gaussians enough?

# Exponential families of distributions


GAMMA


BETA


POISSON


CATEGORICAL

# Multivariate distributions

We've described a variety of distributions for **one-dimensional** data.
What about higher dimensions?

❶ **Naive Bayes**: Treat coordinates as independent.
For $x = (x_1, \ldots, x_d)$, fit separate models $\mathrm{Pr}_i$ to each $x_i$, and assume

$$\mathrm{Pr}(x_1, \ldots, x_d) = \mathrm{Pr}_1(x_1)\mathrm{Pr}_2(x_2) \cdots \mathrm{Pr}_d(x_d).$$

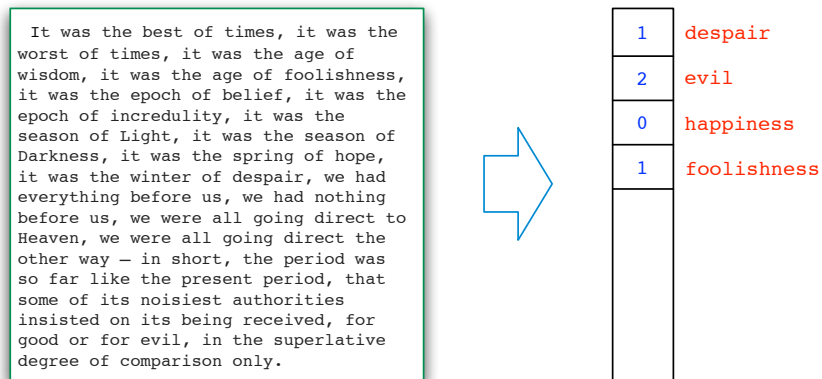This assumption is typically inaccurate.

❷ **Multivariate Gaussian**.
Model correlations between features: we've seen this in detail.

❸ **Graphical models**.
Arbitrary dependencies between coordinates.

# Handling text data

Bag-of-words: vectorial representation of text documents.



```
 It was the best of times, it was the
worst of times, it was the age of
wisdom, it was the age of foolishness,
it was the epoch of belief, it was the
epoch of incredulity, it was the
season of Light, it was the season of
Darkness, it was the spring of hope,
it was the winter of despair, we had
everything before us, we had nothing
before us, we were all going direct to
Heaven, we were all going direct the
other way — in short, the period was
so far like the present period, that
some of its noisiest authorities
insisted on its being received, for
good or for evil, in the superlative
degree of comparison only.
```

| 1 | despair |
| 2 | evil |
| 0 | happiness |
| 1 | foolishness |
| | |

- Fix $V$ = some vocabulary.
- Treat each document as a vector of length $|V|$:

$$x = (x_1, x_2, \ldots, x_{|V|}),$$

where $x_i = \#$ of times the $i$th word appears in the document.

A standard distribution over such document-vectors $x$: the **multinomial**.

# Multinomial naive Bayes

**Multinomial** distribution over a vocabulary $V$:

$$p = (p_1, \ldots, p_{|V|}), \quad \text{such that} \quad p_i \geq 0 \text{ and } \sum_i p_i = 1$$

Document $x = (x_1, \ldots, x_{|V|})$ has probability $\propto p_1^{x_1} p_2^{x_2} \cdots p_{|V|}^{x_{|V|}}$.

For naive Bayes: one multinomial distribution per class.

- Class probabilities $\pi_1, \ldots, \pi_k$
- Multinomials $p^{(1)} = (p_{11}, \ldots, p_{1|V|}), \ldots, p^{(k)} = (p_{k1}, \ldots, p_{k|V|})$

Classify document $x$ as

$$\arg\max_j \quad \pi_j \prod_{i=1}^{|V|} p_{ji}^{x_i}.$$

(As always, take log to avoid underflow: linear classifier.)

# Improving performance of multinomial naive Bayes

A variety of heuristics that are standard in text retrieval, such as:

**1** Compensating for burstiness.
Problem: Once a word has appeared in a document, it has a much higher chance of appearing again.

Solution: Instead of the number of occurrences $f$ of a word, use $\log(1 + f)$.

**2** Downweighting common words.
Problem: Common words can have a unduly large influence on classification.

Solution: Weight each word $w$ by **inverse document frequency**:

$$\log \frac{\# \text{ docs}}{\#(\text{docs containing } w)}$$