

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'

Each data pt is the description of a bottle:

$$x = (\underset{\substack{\uparrow \\ \text{alcohol}}}{_}, \underset{\substack{\uparrow \\ \text{malic} \\ \text{acid}}}{_}, \underset{\substack{\uparrow \\ \text{ash}}}{_}, _, _, \dots, _, _) \in \mathbb{R}^{13}$$

Label $y = 1 \text{ or } 2 \text{ or } 3$

Three classes or labels

Training set:

130 pairs (x, y)



Use this training set
to learn a classifier
that predicts labels
for new x .

Also, a separate test set of 48 labeled points.

we'll use this to evaluate our classifier

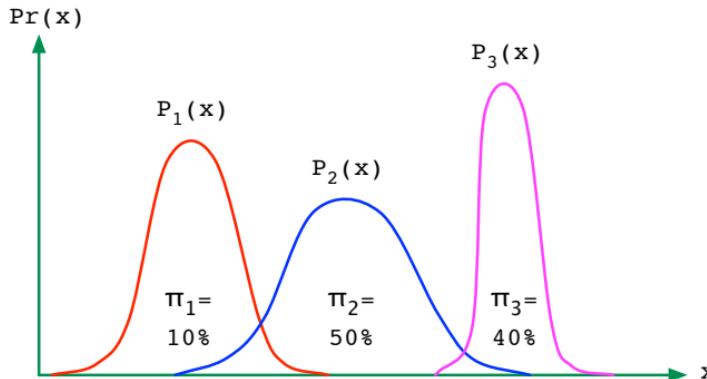
The generative approach to classification

There are many types of classifiers.
We'll look at the generative approach.

Example:

Data space $\mathcal{X} = \mathbb{R}^d$ $d = 13$

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, $\pi_j = \Pr(y = j)$
 - the distribution of data in that class, $P_j(x)$
- What fraction of the data (bottles) are label j ? JUST A NUMBER.
- What is the distribution of x -vectors in class j ? AN ENTIRE DISTRIBUTION.

Overall joint distribution: $\Pr(x, y) = \Pr(y)\Pr(x|y) = \pi_y P_y(x)$.

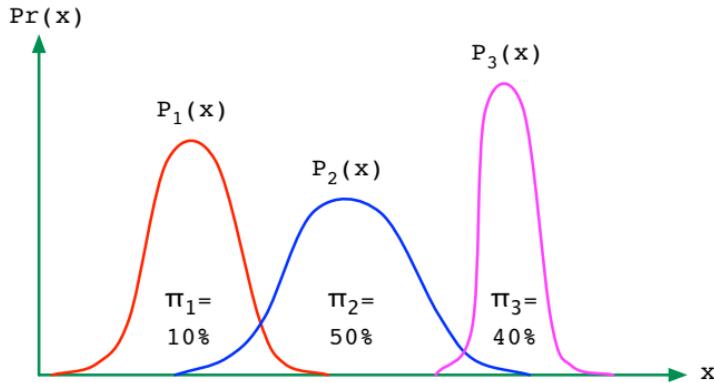
Classifying a new point

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 $\pi_j P_j(x)$. ← very simple rule

Bayes' rule

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'

Fitting a generative model

Proportion of each class : π_1, π_2, π_3

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'

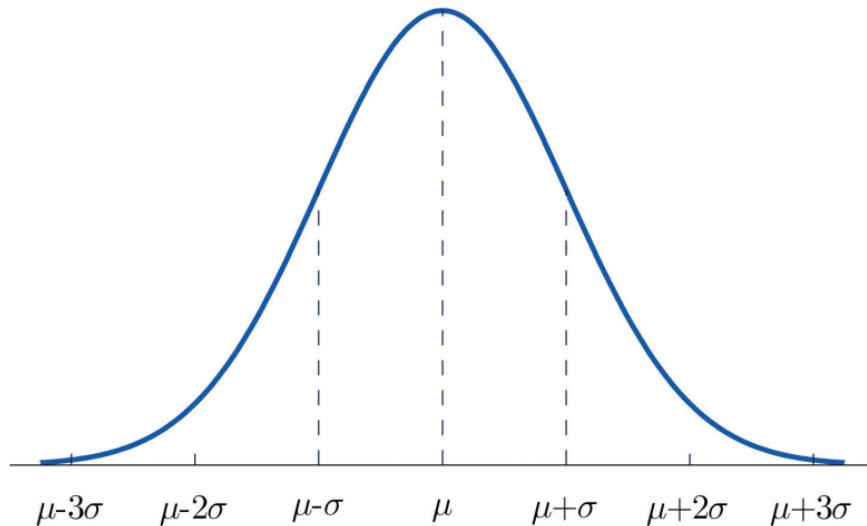
Class weights:

$$\pi_1 = 43/130 = 0.33, \pi_2 = 51/130 = 0.39, \pi_3 = 36/130 = 0.28$$

Need distributions P_1, P_2, P_3 , one per class.
Base these on a single feature: 'Alcohol'.

} Let's start with a simple model that uses only one of the 13 features.

Recall: Univariate Gaussian

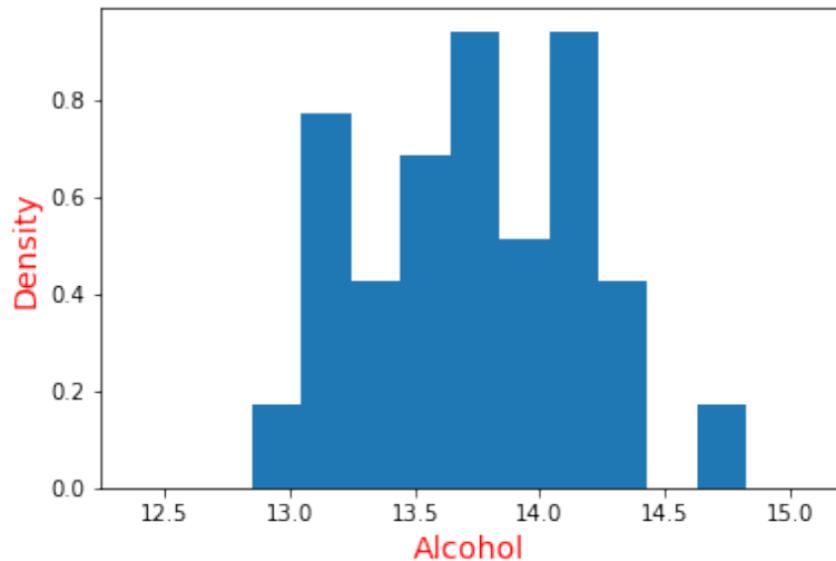


The Gaussian $N(\mu, \sigma^2)$ has mean μ , variance σ^2 , and density function

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

The distribution for winery 1

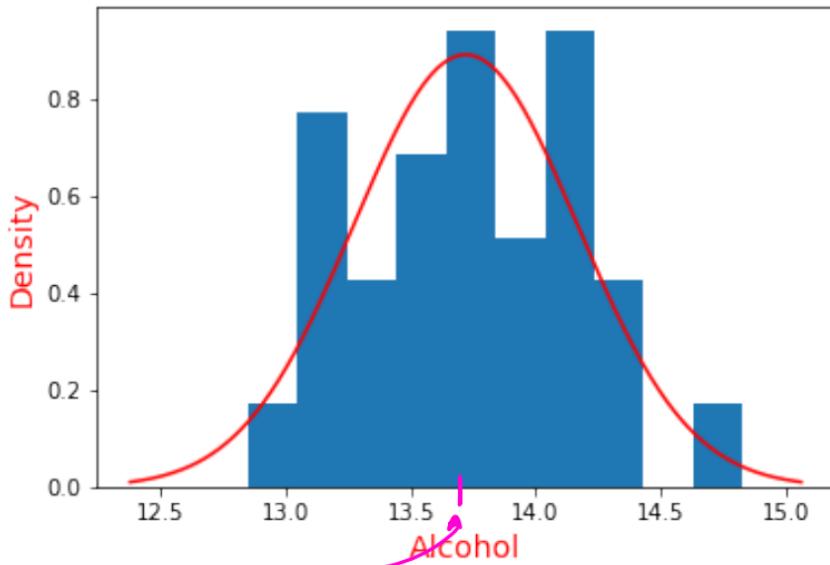
Single feature: 'Alcohol'



The distribution for winery 1

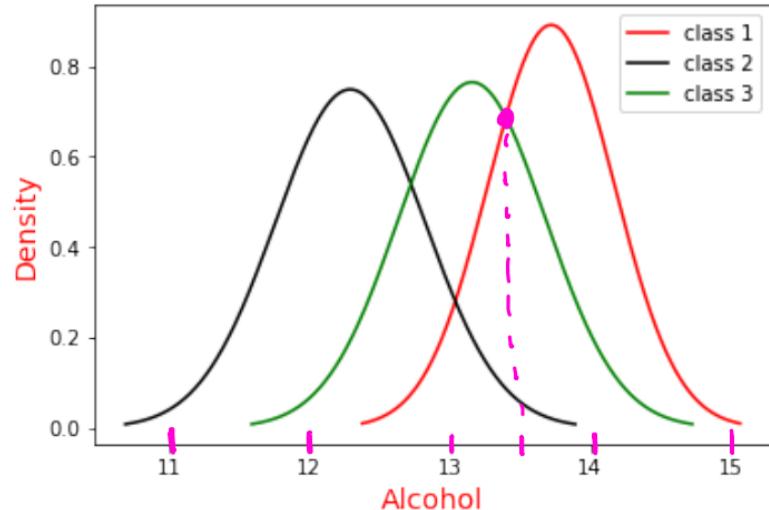
Single feature: 'Alcohol'

The Gaussian smooths out the histogram nicely.



Mean $\mu = 13.72$, Standard deviation $\sigma = 0.44$ (variance 0.20) ↗ two lines of Python

All three wineries



To classify x : compute

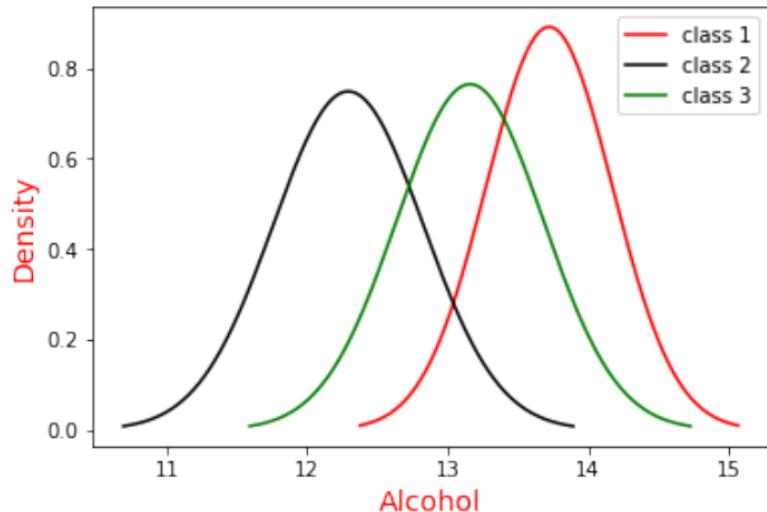
$$\left. \begin{array}{l} \pi_1 P_1(x) \\ \pi_2 P_2(x) \\ \pi_3 P_3(x) \end{array} \right\} \text{pick the largest}$$

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

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

All three wineries

Built model based on training set of 130 bottles.



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

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

evaluate model using test set of 48 bottles.

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

The data set, again

Training set obtained from 130 bottles

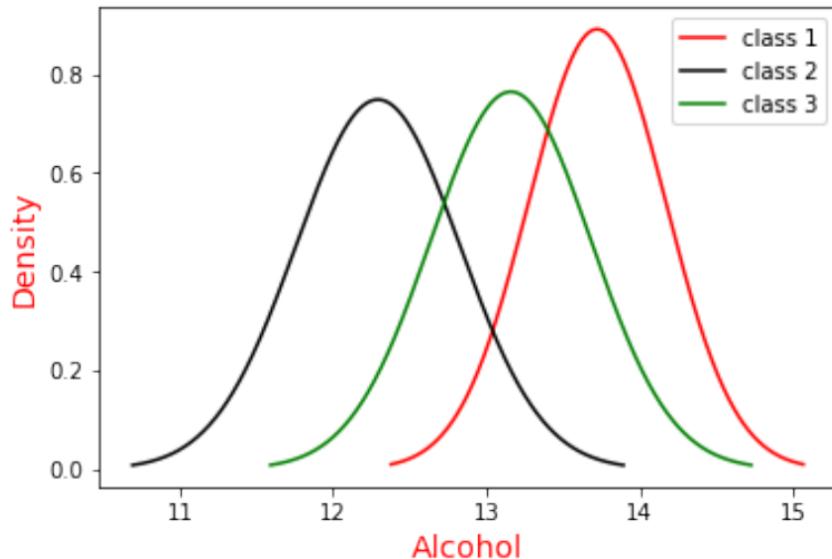
- Winery 1: 43 bottles
- Winery 2: 51 bottles
- Winery 3: 36 bottles
- For each bottle, 13 features:
'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash', 'Magnesium',
'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins',
'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'

Also, a separate test set of 48 labeled points.

This time: 'Alcohol' and 'Flavanoids'.

Why it helps to add features

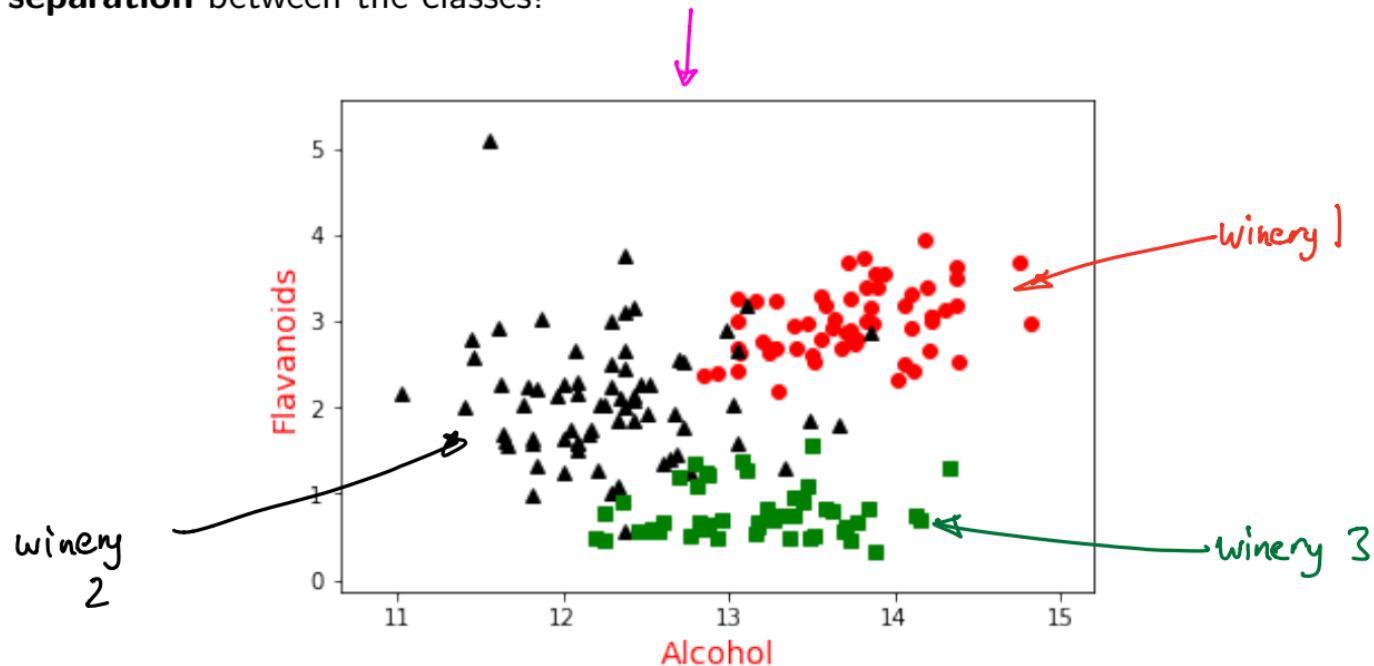
Better **separation** between the classes!



Why it helps to add features

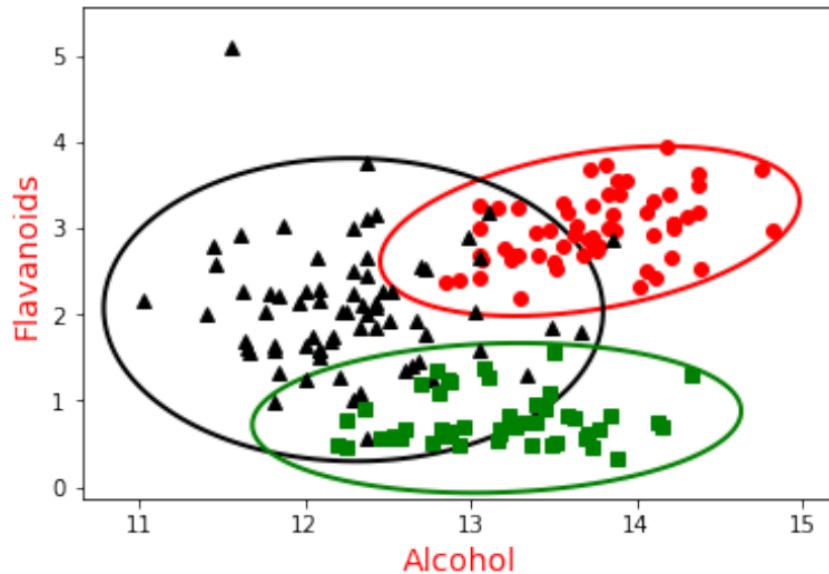
training set of 130 bottles

Better **separation** between the classes!



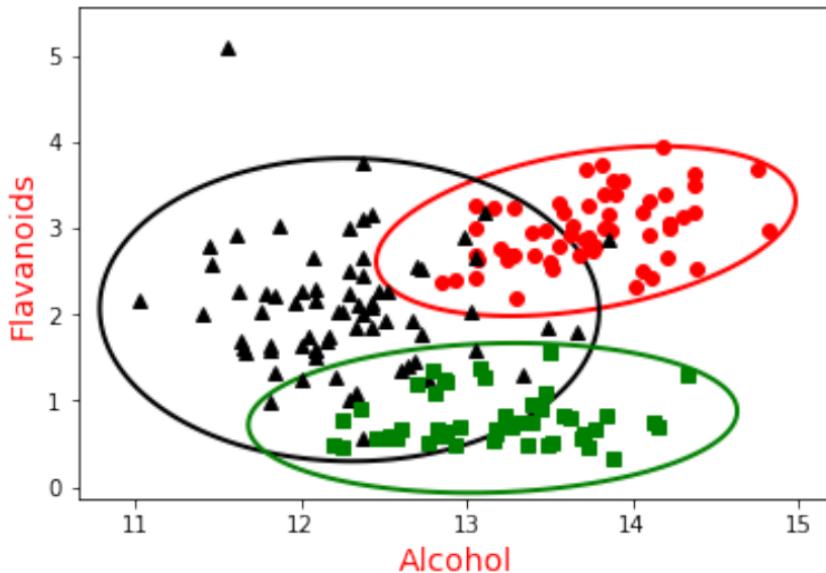
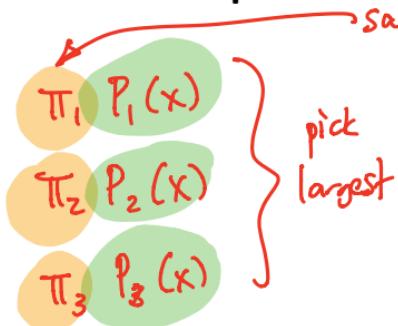
Why it helps to add features

Better **separation** between the classes!



Why it helps to add features

Better **separation** between the classes!



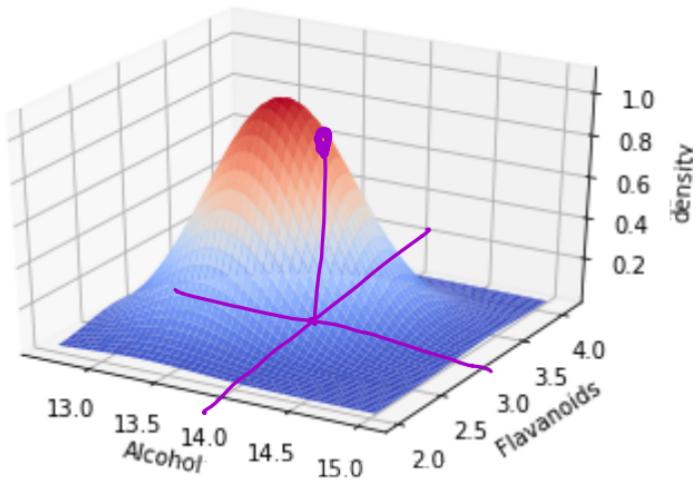
Error rate drops from 29% to 8%.

Test set has 48 points.
How many did we get wrong? 4.

$$\frac{8}{100} \times 48$$

Random guessing	67%
One feature	29%
Two features	8%

Bivariate Gaussian



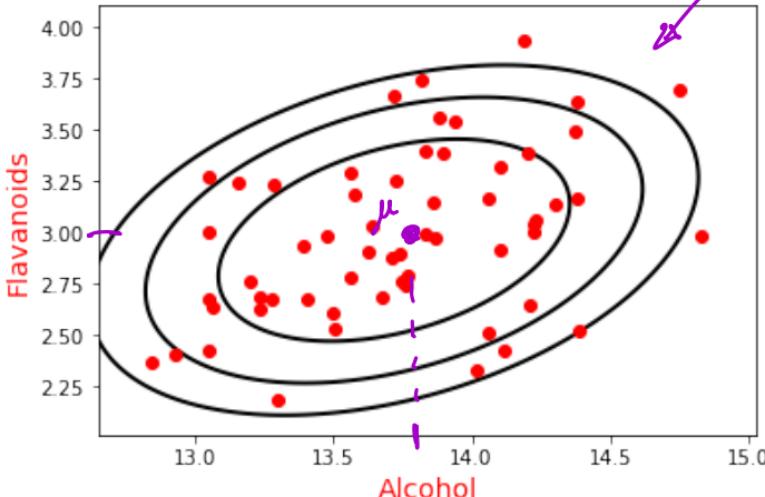
Bivariate Gaussian

What is the correlation btwn these two features?

$$\frac{0.06}{\sqrt{0.2 \times 0.12}} = 0.39$$

weak correlation

Density for winery #1



$P_1(x)$ density for winery 1

$$P_1 : N(\mu_1, \Sigma_1)$$
$$N(\mu^{(1)}, \Sigma^{(1)})$$

Model class 1 by a bivariate Gaussian, parametrized by:

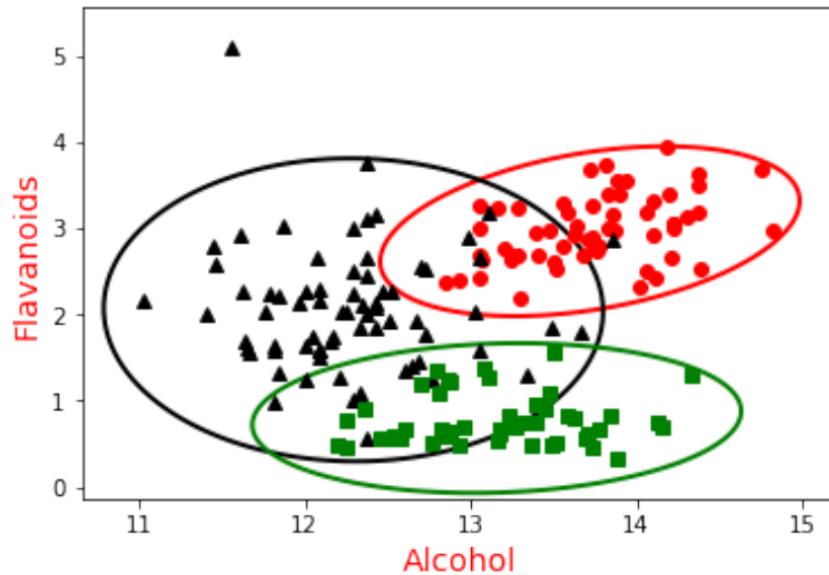
$$\text{mean } \mu = \begin{pmatrix} 13.7 \\ 3.0 \end{pmatrix} \text{ and covariance matrix } \Sigma = \begin{pmatrix} 0.20 & 0.06 \\ 0.06 & 0.12 \end{pmatrix}$$

mean alcohol
mean flavanoid

alcohol variance
covariance btwn alcohol, flavanoid
flavanoid variance

The decision boundary

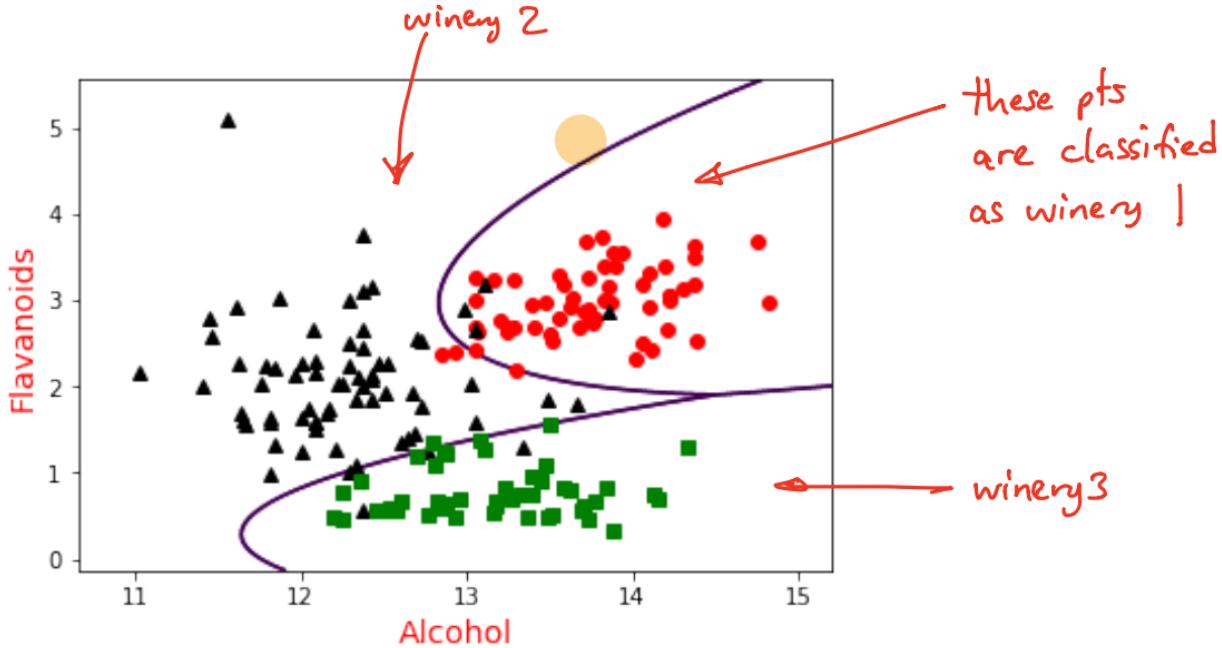
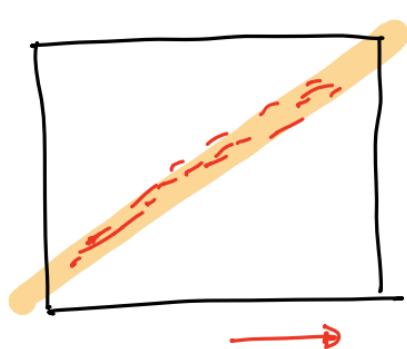
Go from 1 to 2 features: error rate goes from 29% to 8%.



The decision boundary

$$\Pr(y=j|x) = \frac{\pi_j p_j(x)}{\pi_1 p_1(x) + \pi_2 p_2(x) + \pi_3 p_3(x)}$$

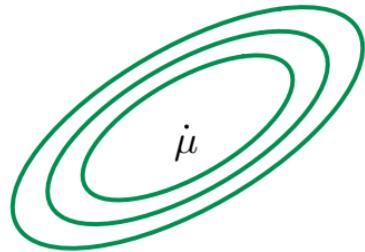
Go from 1 to 2 features: error rate goes from 29% to 8%.



Quadratic

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

Recall: Multivariate Gaussian

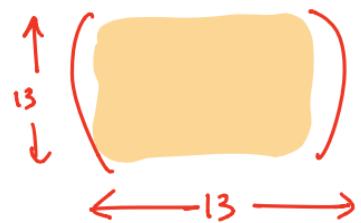


$N(\mu, \Sigma)$: Gaussian in \mathbb{R}^d

- mean: $\mu \in \mathbb{R}^d$
- covariance: $d \times d$ matrix Σ

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

13 feature:
for each winery



- μ is the vector of coordinate-wise means:

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

- Σ is a matrix containing all pairwise covariances:

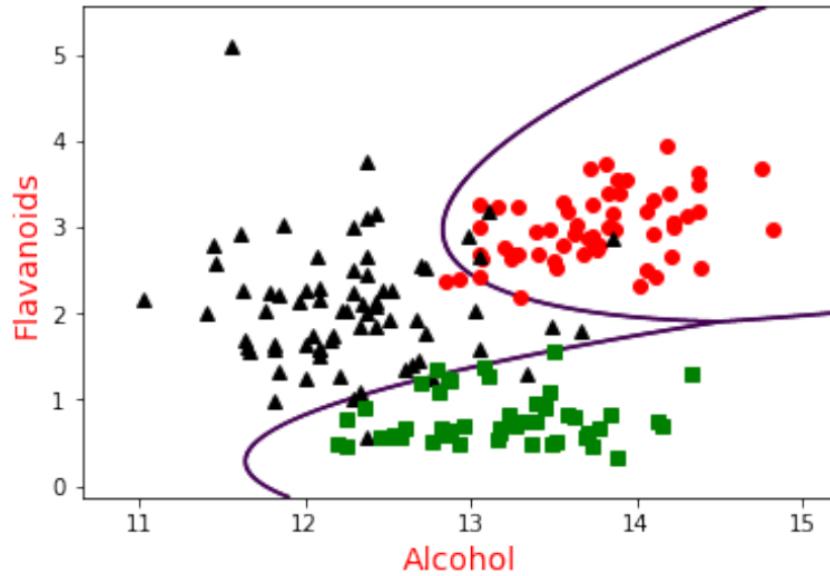
$$\Sigma_{ij} = \Sigma_{ji} = \text{cov}(X_i, X_j) \quad \text{if } i \neq j$$

$$\Sigma_{ii} = \text{var}(X_i)$$

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

Back to the winery data

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



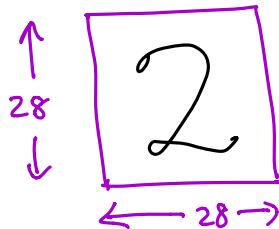
With all 13 features: test error rate goes to zero.

Is this a perfect classifier?

Classification of handwritten digits

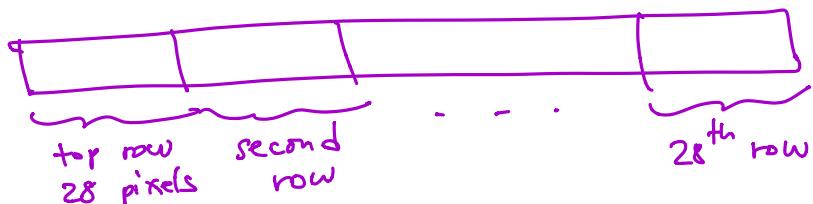
MNIST → 60,000 training images of digits 0-9
→ 10,000 test images

Each image



grayscale
each pixel is 0-255
(0 = white, 255 = black)

Vector of length 784:



Data pts: $x \in \mathbb{R}^{784}$
 $y \in \{0, 1, \dots, 9\}$

10 digits

Training set size 60,000

Roughly 6,000 of each

- For each class:
Covariance matrix $784 \times 784 \leftarrow$ lots of parameters

Need to smooth these matrices.

say Σ_j = empirical covariance for class j
(numpy.cov)

Instead, use $\Sigma_j + cI$ to "smooth" it.

↑ how to choose this?
should only use the training set

Instead of
 $\underset{j}{\operatorname{argmax}} \pi_j P_j(x)$

Can use
 $\underset{j}{\operatorname{argmax}} \log \pi_j + \log P_j(x)$

- In high dimensional spaces, probabilities are tiny.
To avoid problems, often useful to work with log-probabilities.

Binary classification with Gaussian generative model

- Estimate class probabilities π_1, π_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

$$\pi_1 P_1(x) > \pi_2 P_2(x)$$

Binary classification with Gaussian generative model

- Estimate class probabilities π_1, π_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

$$\pi_1 P_1(x) > \pi_2 P_2(x)$$

Binary classification with Gaussian generative model

- Estimate class probabilities π_1, π_2
- Fit a Gaussian to each class: $P_1 = N(\mu_1, \Sigma_1)$, $P_2 = N(\mu_2, \Sigma_2)$

$$\pi_1 \frac{1}{(2\pi)^{d/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_1)^\top \Sigma_1^{-1} (x - \mu_1)\right) > \pi_2 \frac{1}{(2\pi)^{d/2} |\Sigma_2|^{1/2}} \exp\left(-\dots\right)$$

Given a new point x , predict class 1 if

$$\boxed{\pi_1 P_1(x) > \pi_2 P_2(x)} \Leftrightarrow \underbrace{x^\top M x}_{\text{quadratic term}} + \underbrace{2w^\top x}_{\text{linear term}} \geq \theta,$$

where:

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

Q: When is $M = 0$?

A: When $\Sigma_1 = \Sigma_2$.

and θ is a threshold depending on the various parameters.

Linear or quadratic decision boundary.

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

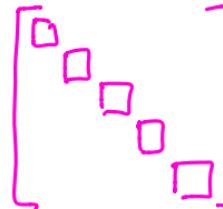
Linear decision boundary: choose class 1 if

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

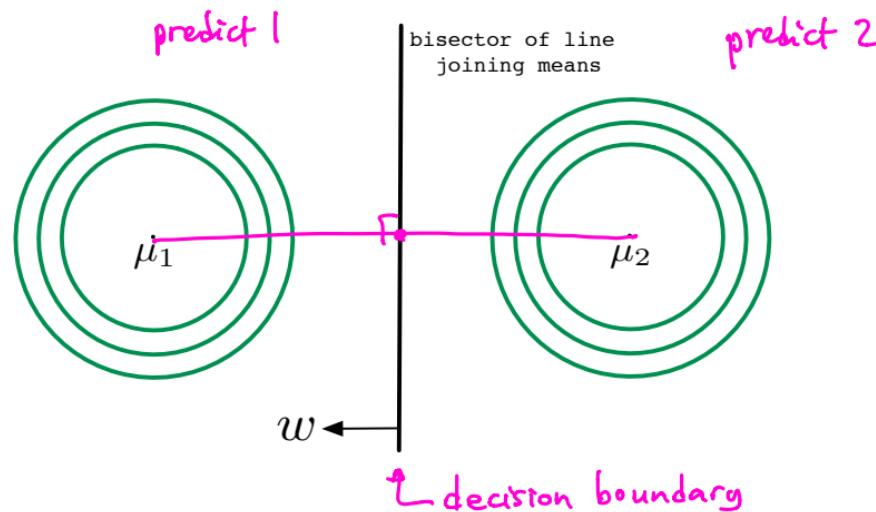
Common covariance: $\Sigma_1 = \Sigma_2 = \Sigma$

Linear decision boundary: choose class 1 if

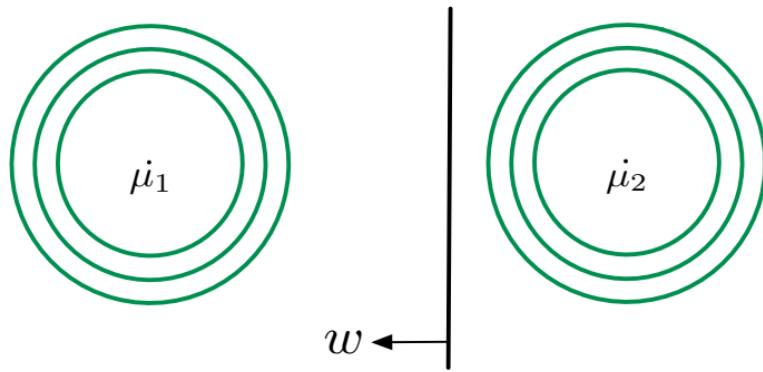
$$x \cdot \underbrace{\Sigma^{-1}(\mu_1 - \mu_2)}_w \geq \theta.$$



Example 1: Spherical Gaussians with $\Sigma = I_d$ and $\pi_1 = \pi_2 = \frac{1}{2}$.

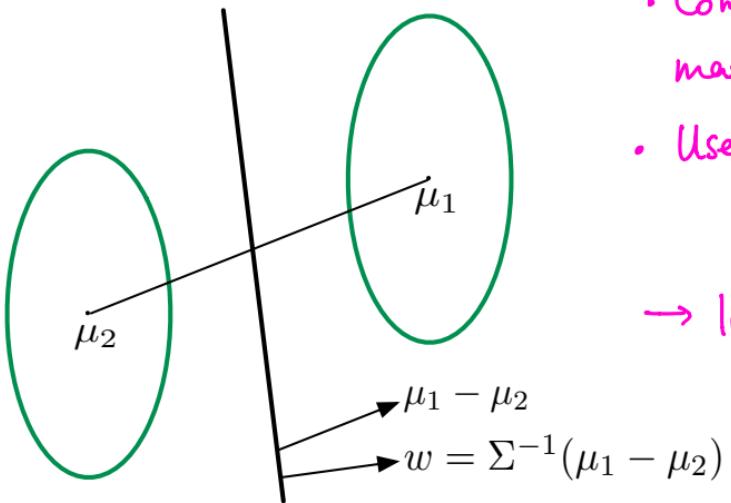


Example 2: Again spherical, but now $\pi_1 > \pi_2$.



Example 3: Non-spherical.

$$\Sigma_1 = \Sigma_2$$



Common case:

- Compute empirical covariance matrices $\hat{\Sigma}_1$ and $\hat{\Sigma}_2$
- Use a common covariance $\Sigma_1 = \Sigma_2 = \frac{1}{2} (\hat{\Sigma}_1 + \hat{\Sigma}_2)$
→ linear decision boundary

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

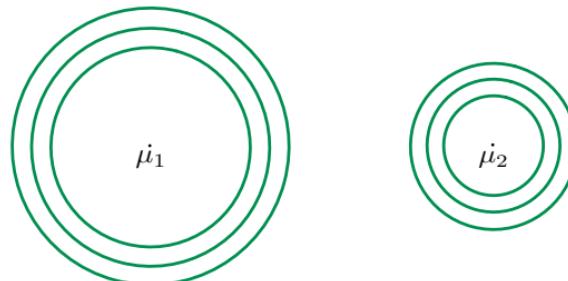
- Choose w as above
- Common practice: fit θ to minimize training or validation error

Different covariances: $\Sigma_1 \neq \Sigma_2$

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

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

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

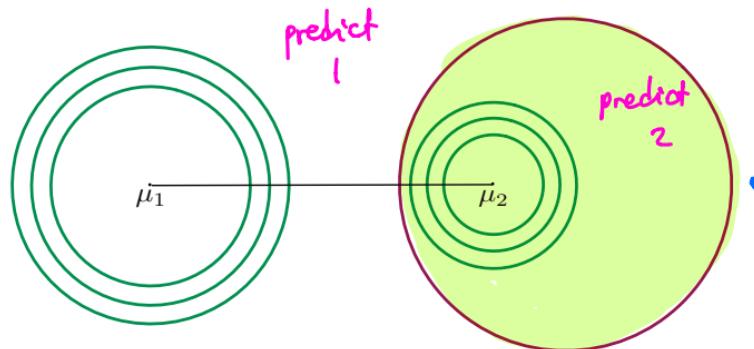


Different covariances: $\Sigma_1 \neq \Sigma_2$

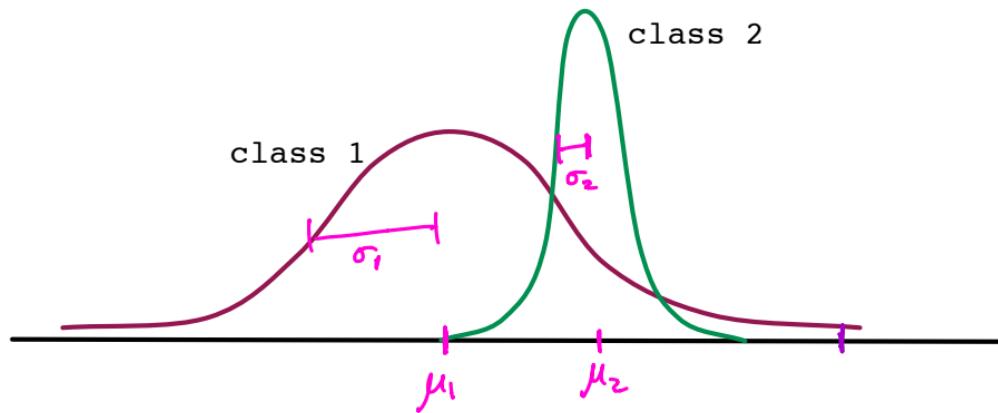
Quadratic boundary: choose class 1 if $x^T Mx + 2w^T x \geq \theta$, where:

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

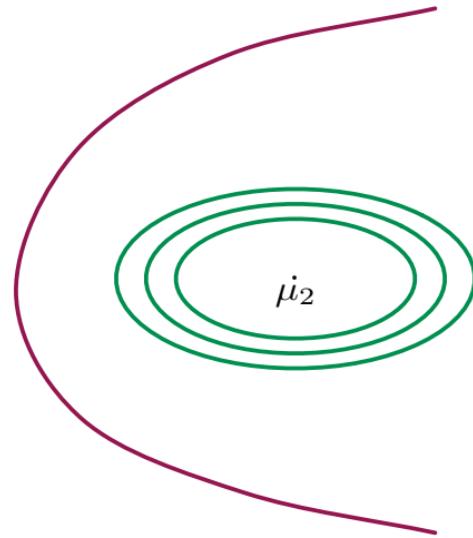
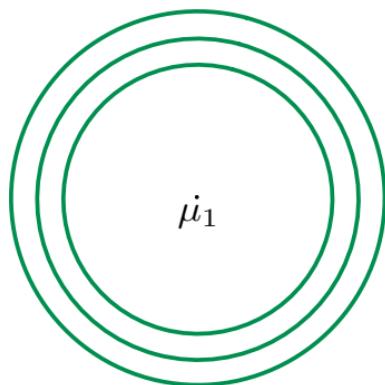
Example 1: $\Sigma_1 = \sigma_1^2 I_d$ and $\Sigma_2 = \sigma_2^2 I_d$ with $\sigma_1 > \sigma_2$



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



Example 3: A parabolic boundary.



Multiclass discriminant analysis

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

Each class has an associated **quadratic** function

$$f_j(x) = \log(\pi_j P_j(x))$$

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

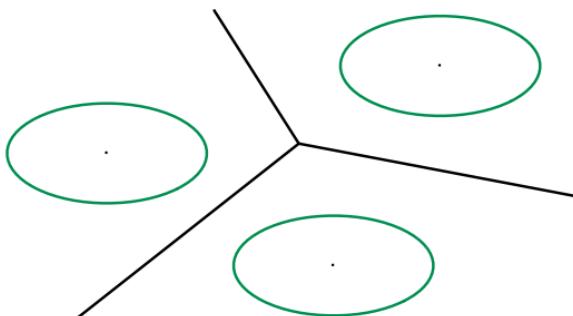
If $\Sigma_1 = \dots = \Sigma_k$, the boundaries are **linear**.

To predict at x , return

$$\operatorname{argmax}_j \pi_j P_j(x)$$

$$\equiv \operatorname{argmax}_j \underbrace{\log(\pi_j P_j(x))}_{f_j(x)}$$

quadratic function

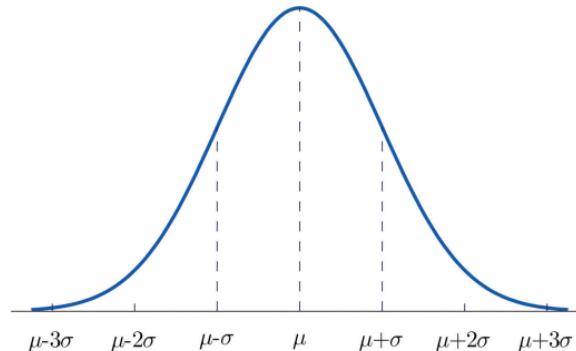


Beyond Gaussians

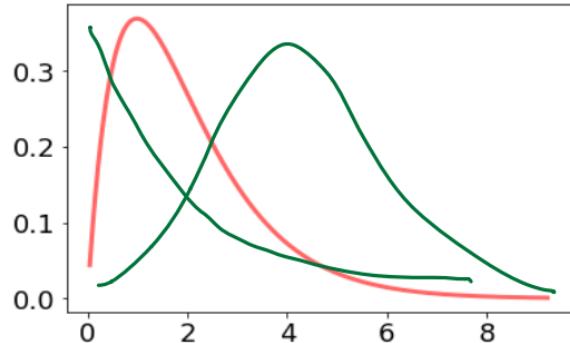
The generative methodology:

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

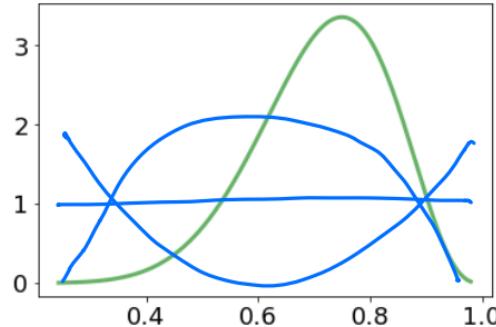
What distribution to use? Are Gaussians enough?



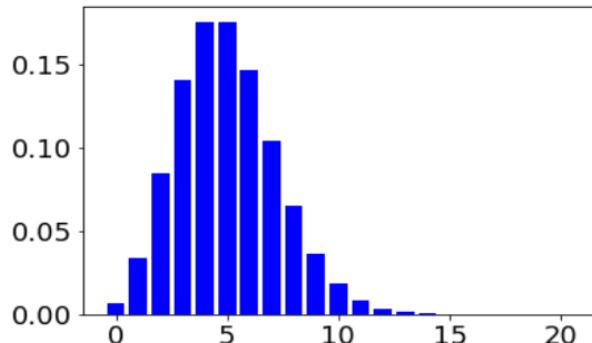
Exponential families of distributions



GAMMA



BETA



POISSON

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.



A legend for a categorical distribution. It shows four categories: 'despair' (blue), 'evil' (red), 'happiness' (green), and 'foolishness' (orange). The legend is organized into four rows: 1 (despair), 2 (evil), 0 (happiness), and 1 (foolishness).

1	despair
2	evil
0	happiness
1	foolishness

CATEGORICAL

Multivariate distributions

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

- ① **Naive Bayes:** Treat coordinates as independent.

For $x = (x_1, \dots, x_d)$, fit separate models \Pr_i to each x_i , and assume

$$\Pr(x_1, \dots, x_d) = \Pr_1(x_1)\Pr_2(x_2) \cdots \Pr_d(x_d).$$

This assumption is typically inaccurate.

- ② **Multivariate Gaussian.**

Model correlations between features: we've seen this in detail.

- ③ **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 $|V|$:

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

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 :

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

Document $x = (x_1, \dots, x_{|V|})$ 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 π_1, \dots, π_k
- Multinomials $p^{(1)} = (p_{11}, \dots, p_{1|V|}), \dots, p^{(k)} = (p_{k1}, \dots, p_{k|V|})$

Classify document x as

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

(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.

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 $\log(1 + f)$.

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.

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**:

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