

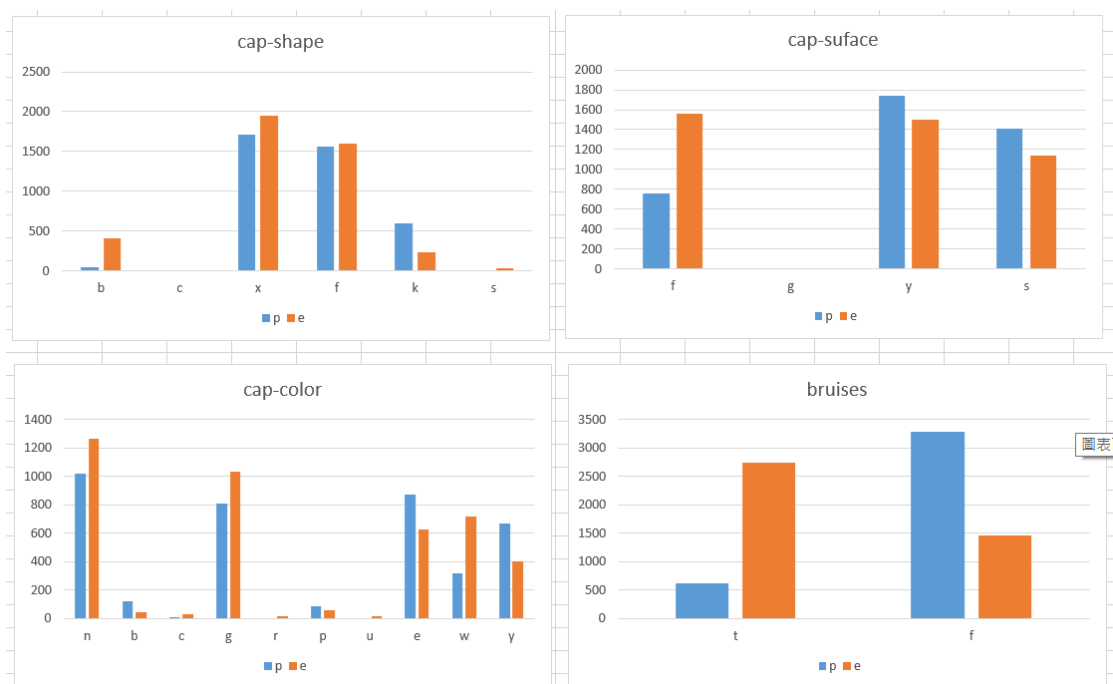
Mushroom

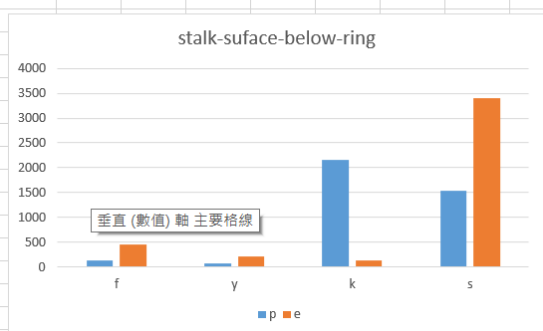
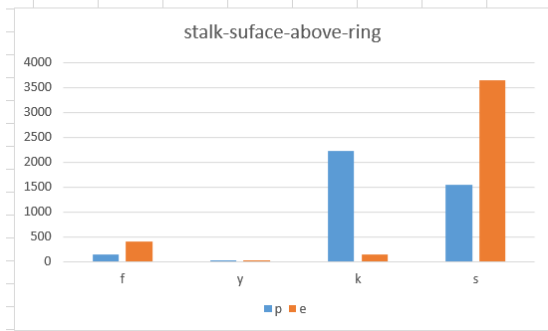
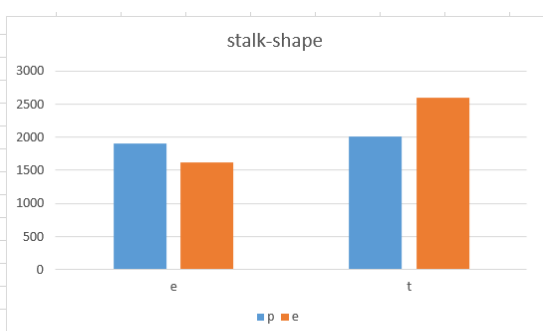
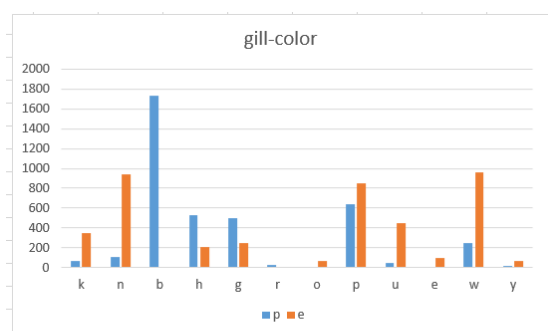
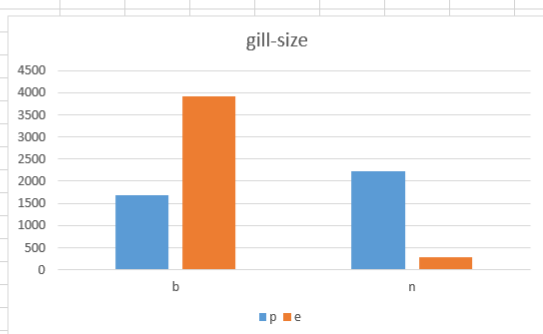
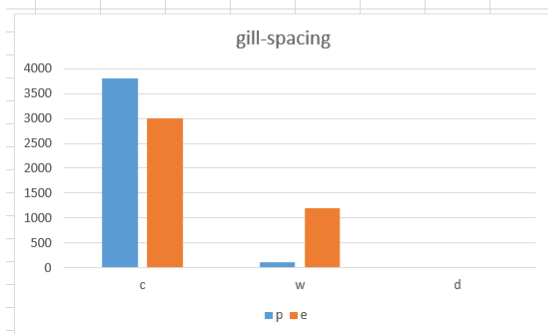
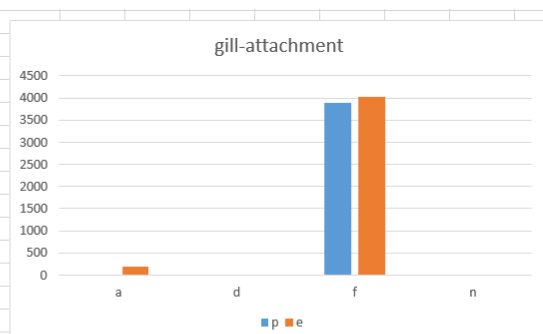
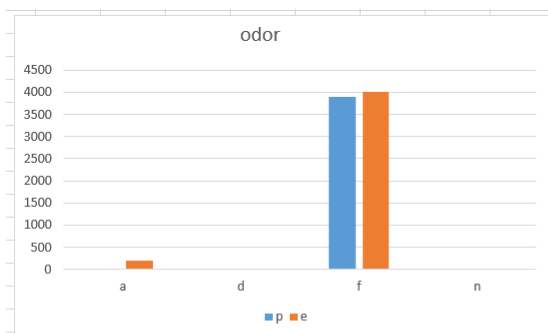
Data input:

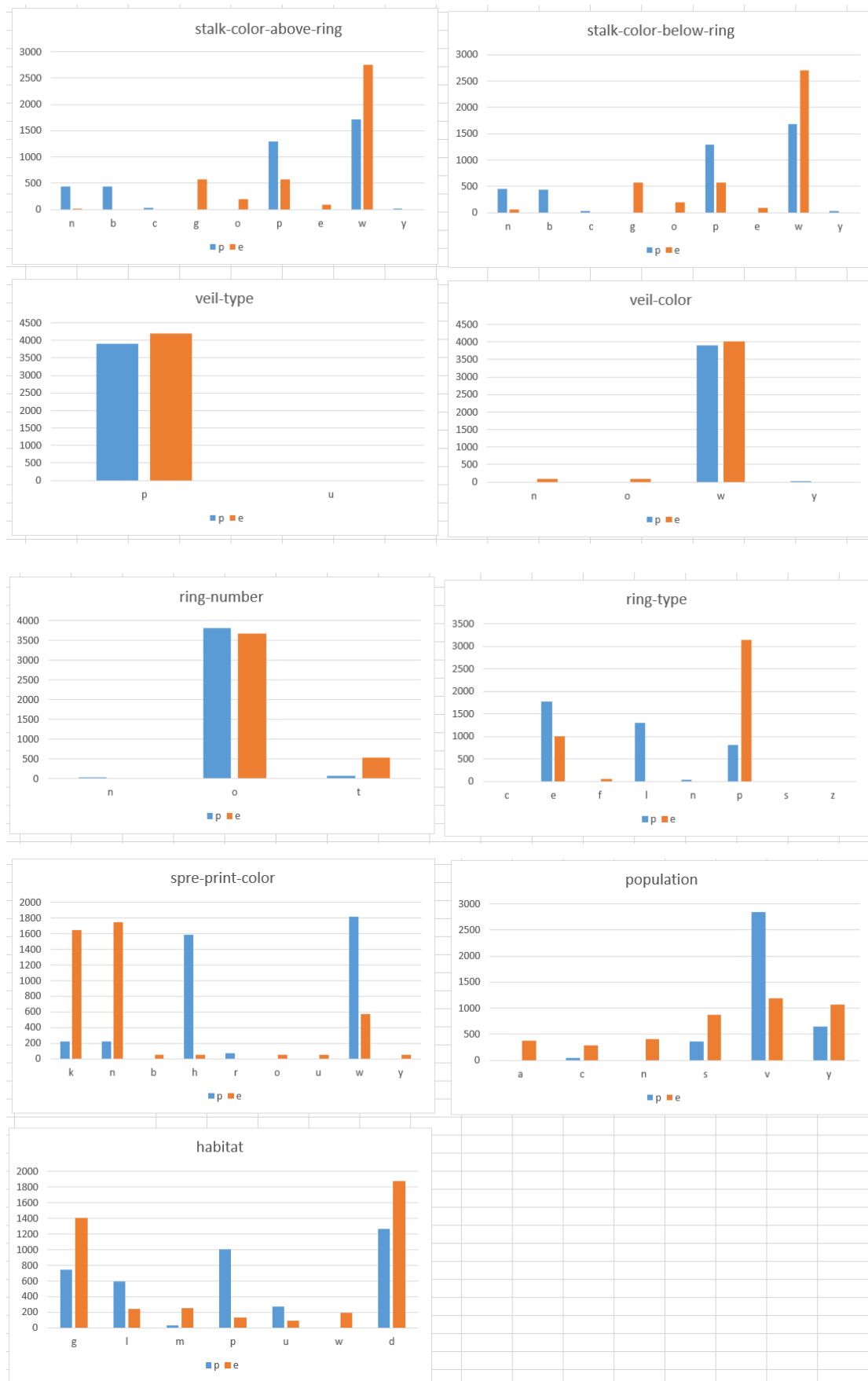
用 pandas 的 `read_csv()` 讀取資料 `agaricus-lepiota.data`

Data Visualization:

使用 Excel:







Data preprocessing:

1. 將有 missing 的 feature 移除，作法是先將 "?" 替換成 nan，再用 dropna()把那一欄刪除
2. 打亂順序

Model construction:

無 Laplace smoothing 的 model: 寫在 naïve_bayes 函式中，主要原理為:

$$(q) = \operatorname{argmax}_{Y \in \mathbb{T}} [\log P(Y) + \sum_{mi=1} \log P(X_i|Y)]$$

$$P(X_i|Y) = N(X_i|Y) / N(Y)$$

有做 Laplace smoothing 的 model: 寫在 naïve_bayes_smoothing 函式中，原理為:

$$(q) = \operatorname{argmax}_{Y \in \mathbb{T}} [\log P(Y) + \sum_{mi=1} \log P(X_i|Y)]$$

$$P(X_i|Y) = (N(X_i|Y) + k) / (N(Y) + k\tau), \text{ 取 } k=1$$

Train-and-spilt:

分為 train : test = 7 : 3

```
train = df.sample(frac=0.7, random_state=200)
```

```
test = df.drop(train.index)
```

Results:

```
Result without Laplace smoothing:
Actual edible      Predicted edible      Predicted poisonous
Actual edible      1158                      123
Actual poisonous   255                      901
Accuracy: 0.844891
Sensitivity: 0.903981
Precision: 0.819533

Result with Laplace smoothing(k=1):
Actual edible      Predicted edible      Predicted poisonous
Actual edible      1277                      4
Actual poisonous   97                      1059
Accuracy: 0.958556
Sensitivity: 0.996877
Precision: 0.929403
```

Without Laplace smoothing:

Confusion Matrix:

	Predicted edible	Predicted poisonous
Actual edible	1158	123
Actual poisonous	255	901

Accuracy: 0.884891

Sensitivity: 0.903981

Precision: 0.819533

With Laplace smoothing:

Confusion Matrix

	Predicted edible	Predicted poisonous
Actual edible	1277	4
Actual poisonous	97	1059

Accuracy: 0.958556

Sensitivity: 0.996877

Precision: 0.929403

Comparison & Conclusion:

因為每次都是取隨機的樣本來訓練 model，所以結果每次都不一樣，但是有做 Laplace smoothing 的 model 都會有較高的 Accuracy, Sensitivity, Precision，所以顯示做 Laplace smoothing 可使 model 準確度提高。

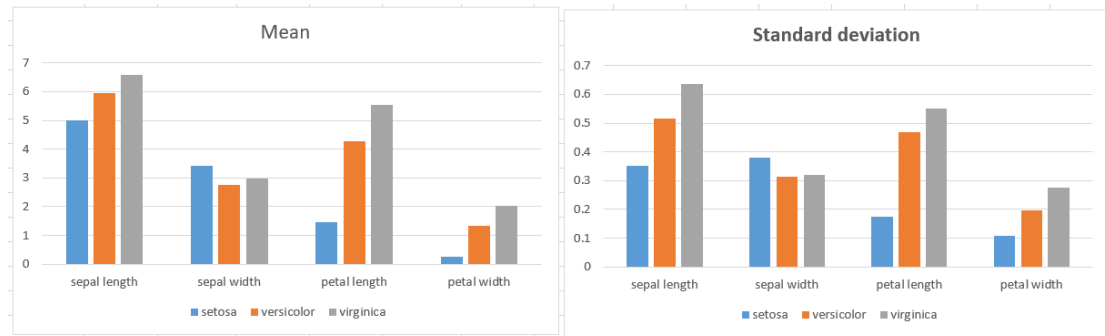
Iris

Data input:

用 pandas 的 `read_csv()` 讀取資料 `iris.data`

Data Visualization:

使用 Excel:



Data preprocessing:

打亂順序

Model construction:

$$(\mathbf{q}) = \operatorname{argmax}_{Y \in \mathbb{T}} [\log P(Y) + \sum_{i=1}^m \log P(X_i | Y)]$$

$$P(X_i | Y) = \frac{1}{(\sigma \sqrt{2\pi})} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Train-and-spilt:

分為 train : test = 7 : 3

```
train = df.sample(frac=0.7, random_state=200)
```

```
test = df.drop(train.index)
```

Results:

	Predicted Setosa	Predicted Virginica	Predicted Versicolour
Actual Setosa	14	0	0
Actual Virginica	0	18	1
Actual Versicolour	0	12	0

Confusion Matrix:

	Predicted Setosa	Predicted Virginica	Predicted Versicolour
Actual Setosa	14	0	0
Actual Virginica	0	18	1
Actual Versicolour	0	12	0

Accuracy: 0.711111

Precision (Setosa = True): 1

Sensitivity (Setosa = True): 1

Precision (Virginica = True): 0.947368

Sensitivity (Virginica = True): 0.6

Precision (Versicolour = True): 0

Sensitivity (Versicolour = Ture): 0

Comparison & Conclusion:

在這個 model 中 Versicolour 會被誤認為是 Virginica，可能是因為 Versicolour 這個品種的花和 Virginica 沒有太突出的特徵。