



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

k-NN
+
Model Selection
+
Perceptron

Matt Gormley
Lecture 5
Feb. 15, 2021

Q&A

Q: Why don't my entropy calculations match those on the slides?

A: $H(Y)$ is conventionally reported in “bits” and computed using log base 2.
e.g., $H(Y) = -P(Y=0) \log_2 P(Y=0) - P(Y=1) \log_2 P(Y=1)$

Q: Why is entropy based on a sum of $p(\cdot) \log p(\cdot)$ terms?

A: We don't have time for a full treatment of why it *has* to be this, but we can develop the right intuition with a few examples...

Q&A

Q: How do we define a distance function when the features are categorical (e.g. weather takes values {sunny, rainy, overcast})?

A: Step 1: Convert from categorical attributes to numeric features (e.g. binary)
Step 2: Select an appropriate distance function (e.g. Hamming distance)

Reminders

- **Homework 2: Decision Trees**
 - Out: Wed, Feb. 10
 - Due: **Mon**, Feb. 22 at 11:59pm
- **Today's Poll:**
 - <http://poll.mlcourse.org>
 - fill out first two questions about HW1

Moss Cheat Checker

What is Moss?

- Moss (Measure Of Software Similarity): is an automatic system for determining the similarity of programs. To date, the main application of Moss has been in detecting plagiarism in programming classes.
- Moss reports:
 - The Andrew IDs associated with the file submissions
 - The number of lines matched
 - The percent lines matched
 - Color coded submissions where similarities are found

What is Moss?

At first glance, the submissions may look different

```
# Python program to find ordered words
import requests

# Scrapes the words from the URL below and stores
# them in a list
def getWords():

# contains about 2500 words
url = "http://www.puzzlers.org/pub/wordlists/unixdict.txt"
fetchData = requests.get(url)

# extracts the content of the webpage
wordList = fetchData.content

# decodes the UTF-8 encoded text and splits the
# string to turn it into a list of words
wordList = wordList.decode("utf-8").split()

return wordList

# function to determine whether a word is ordered or not
def isOrdered():

# fetching the wordList
collection = getWords()

# since the first few of the elements of the
# dictionary are numbers, getting rid of those
# numbers by slicing off the first 17 elements
collection = collection[16:]
word = ''

for word in collection:
    result = 'Word is ordered'
    i = 0
    l = len(word) - 1

    if (len(word) < 3): # skips the 1 and 2 lettered strings
        continue

    # traverses through all characters of the word in pairs
    while i < l:
        if (ord(word[i]) > ord(word[i+1])):
            result = 'Word is not ordered'
            break
        else:
            i += 1

# only printing the ordered words
if (result == 'Word is ordered'):
    print(word, ': ', result)

# execute isOrdered() function
if __name__ == '__main__':
    isOrdered()
```

```
import requests

def Ordered():
    coll = getWs()
    coll = coll[16:]
    word = ''
    for word in coll:
        r = 'Word is ordered'
        a = 0
        length = len(word) - 1
        if (len(word) < 3):
            continue
        while a < length:
            if (ord(word[a]) > ord(word[a+1])):
                r = 'Word is not ordered'
                break
            else:
                a += 1
        if (r == 'Word is ordered'):
            print(word, ': ', r)

def getWs():
    url = "http://www.puzzlers.org/pub/wordlists/unixdict.txt"
    fetch = requests.get(url)
    words = fetch.content
    words = words.decode("utf-8").split()
    return words

if __name__ == '__main__':
    Ordered()
```

What is Moss?

Moss can quickly find the similarities

```
''' File: leetcodefindword.com.edu_1_bundin.v
# Python program to find ordered words
import requests

# Scrapes the words from the URL, below and stores
# them in a list

def getWords():

# contains about 2100 words
url = "https://www.postleone.org/pub/wordlists/lexicollat.txt"
fetchData = requests.get(url)

# extracts the content of the webpage
wordList = fetchData.content

# decodes the UTF-8 encoded text and splits the
# string to turn it into a list of words
wordList = wordList.decode("utf-8").split()

return wordList

# Function to determine whether a word is ordered or not
def isOrdered():

# fetching the wordList
collection = getWords()

# since the first few of the elements of the
# dictionary are numbers, getting rid of those
# numbers by slicing off the first 17 elements
collection = collection[17:]
word = ""

for word in collection:
    result = "Word is ordered"
    i = 0
    l = len(word) - 1

    if (len(word) < 3): # skips the 1 and 2 lettered strings
        continue

    # traverses through all characters of the word in pairs
    while i < l:
        if (ord(word[i]) > ord(word[i+1])):
            result = "Word is not ordered"
            break
        else:
            i += 1

# only printing the ordered words
if (result == "Word is ordered"):
    print(word, ", ", end='')

# executes isOrdered() function
if __name__ == "__main__":
    isOrdered()
```

```
''' File: leetcodefindword.com.edu_1_bundin.v
import requests

def isOrdered():
    url = getURL()
    url = url[17:]
    word = ""
    for word in url:
        i = "Word is ordered"
        a = 0
        length = len(word) - 1
        if (len(word) < 3):
            continue
        while a < length:
            if (ord(word[a]) > ord(word[a+1])):
                i = "Word is not ordered"
                break
            else:
                a += 1
        if (i == "Word is ordered"):
            print(word, ", ", end='')

def getURL():
    url = "https://www.postleone.org/pub/wordlists/lexicollat.txt"
    fetch = requests.get(url)
    words = fetch.content
    words = words.decode("utf-8").split()
    return words

if __name__ == "__main__":
    isOrdered()
```


DECISION BOUNDARIES

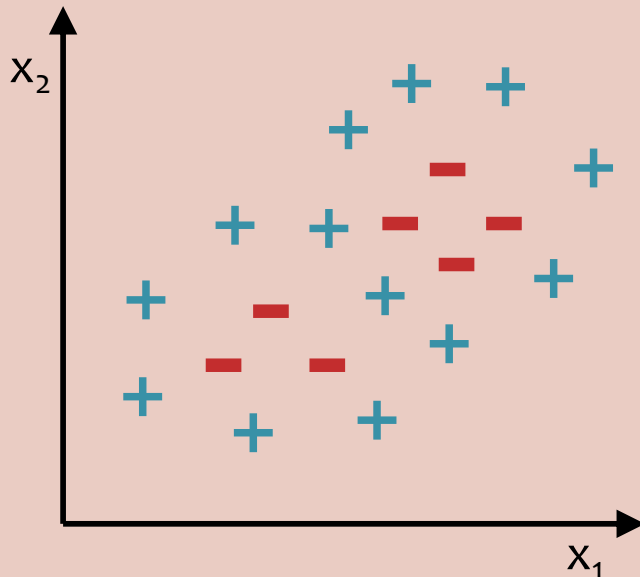
Decision Boundary Example

Dataset: Outputs $\{+, -\}$; Features x_1 and x_2

In-Class Exercise

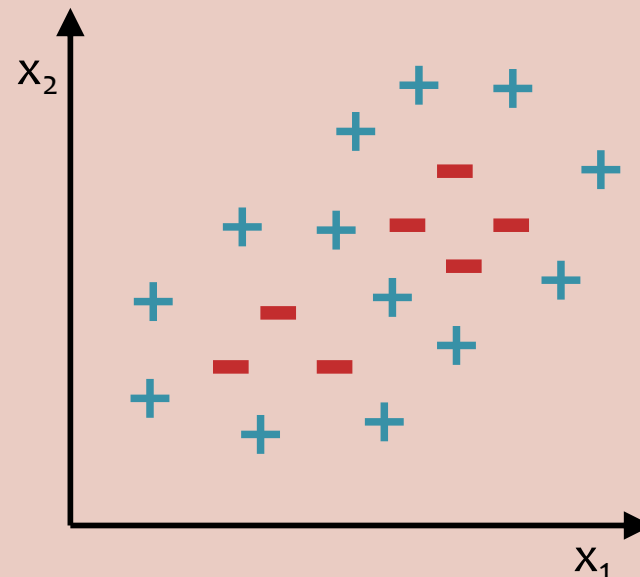
Question:

- A. Can a **k-Nearest Neighbor classifier with $k=1$** achieve **zero training error** on this dataset?
- B. If **'Yes'**, draw the learned decision boundary. If **'No'**, why not?



Question:

- A. Can a **Decision Tree classifier** achieve **zero training error** on this dataset?
- B. If **'Yes'**, draw the learned decision boundary. If **'No'**, why not?



k-Nearest Neighbors

Whiteboard:

- Decision Tree boundary with continuous features

KNN ON FISHER IRIS DATA



Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7

Fisher Iris Dataset

Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width	Petal Length	Petal Width
0	4.3	3.0	1.1	0.1
0	4.9	3.6	1.4	0.1
0	5.3	3.7	1.5	0.2
1	4.9	2.4	3.3	1.0
1	5.7	2.8	4.1	1.3
1	6.3	3.3	4.7	1.6
1	6.7	3.0	5.0	1.7

Fisher Iris Dataset

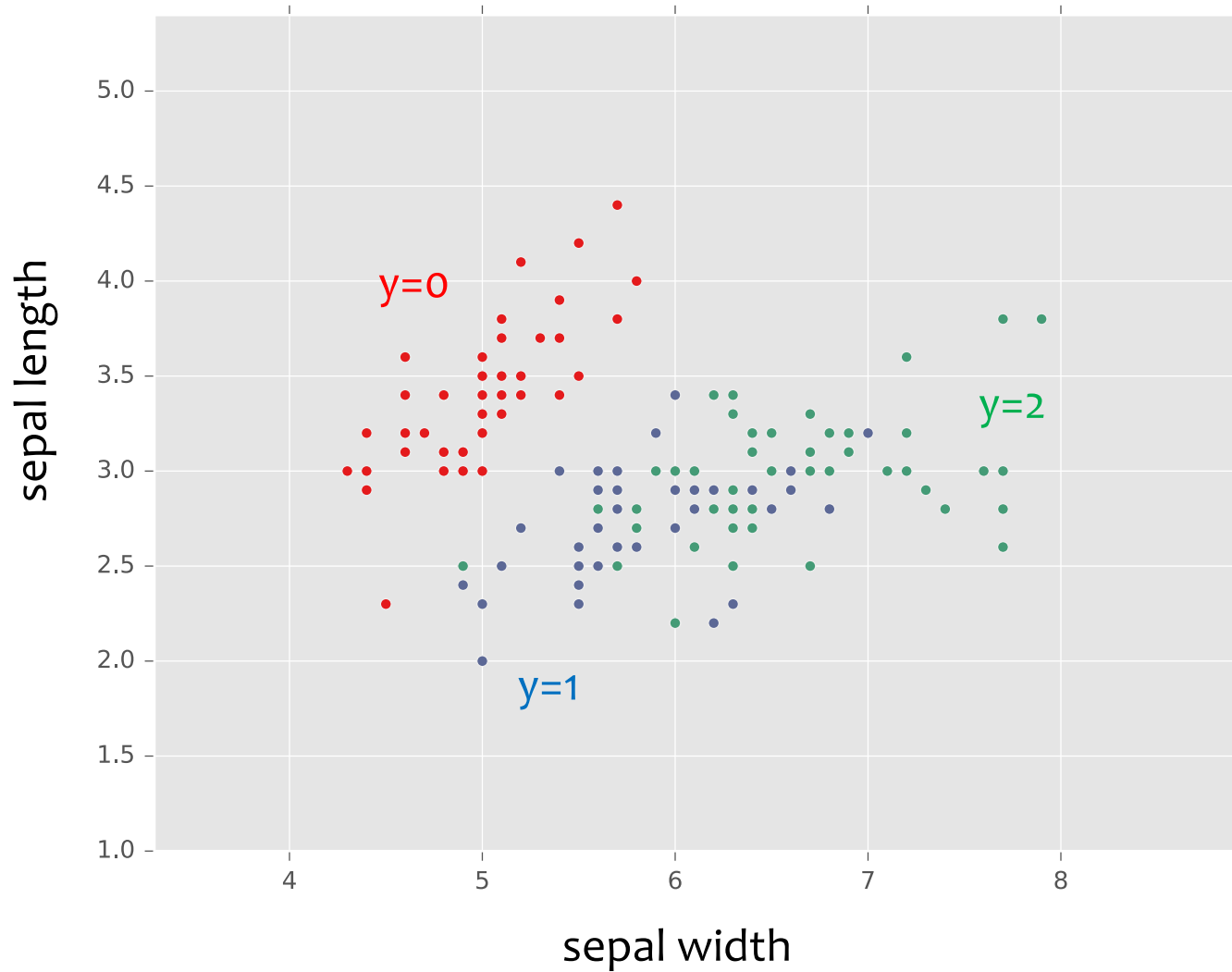
Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

Species	Sepal Length	Sepal Width
0	4.3	3.0
0	4.9	3.6
0	5.3	3.7
1	4.9	2.4
1	5.7	2.8
1	6.3	3.3
1	6.7	3.0

Deleted two of the four features, so that input space is 2D

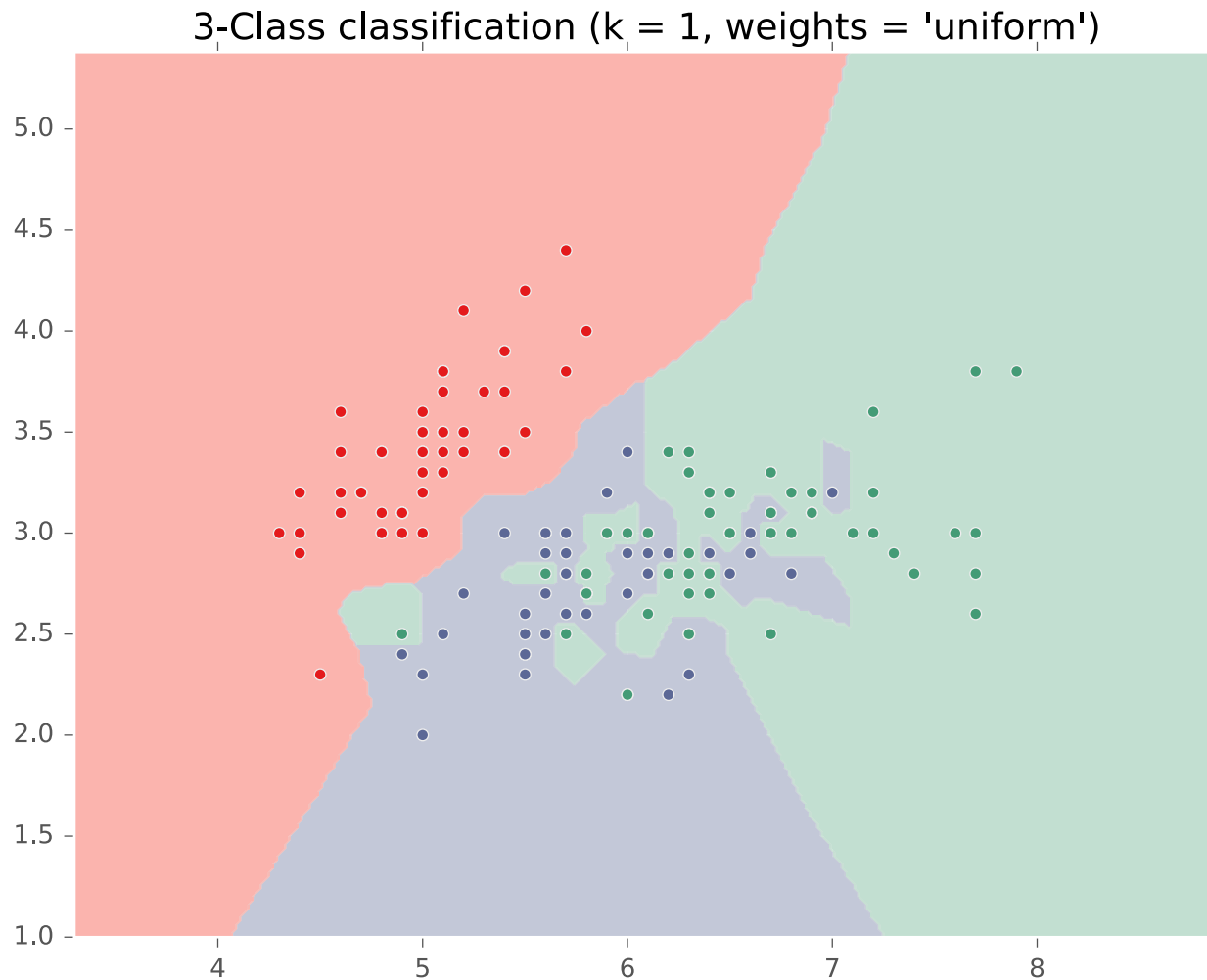


KNN on Fisher Iris Data

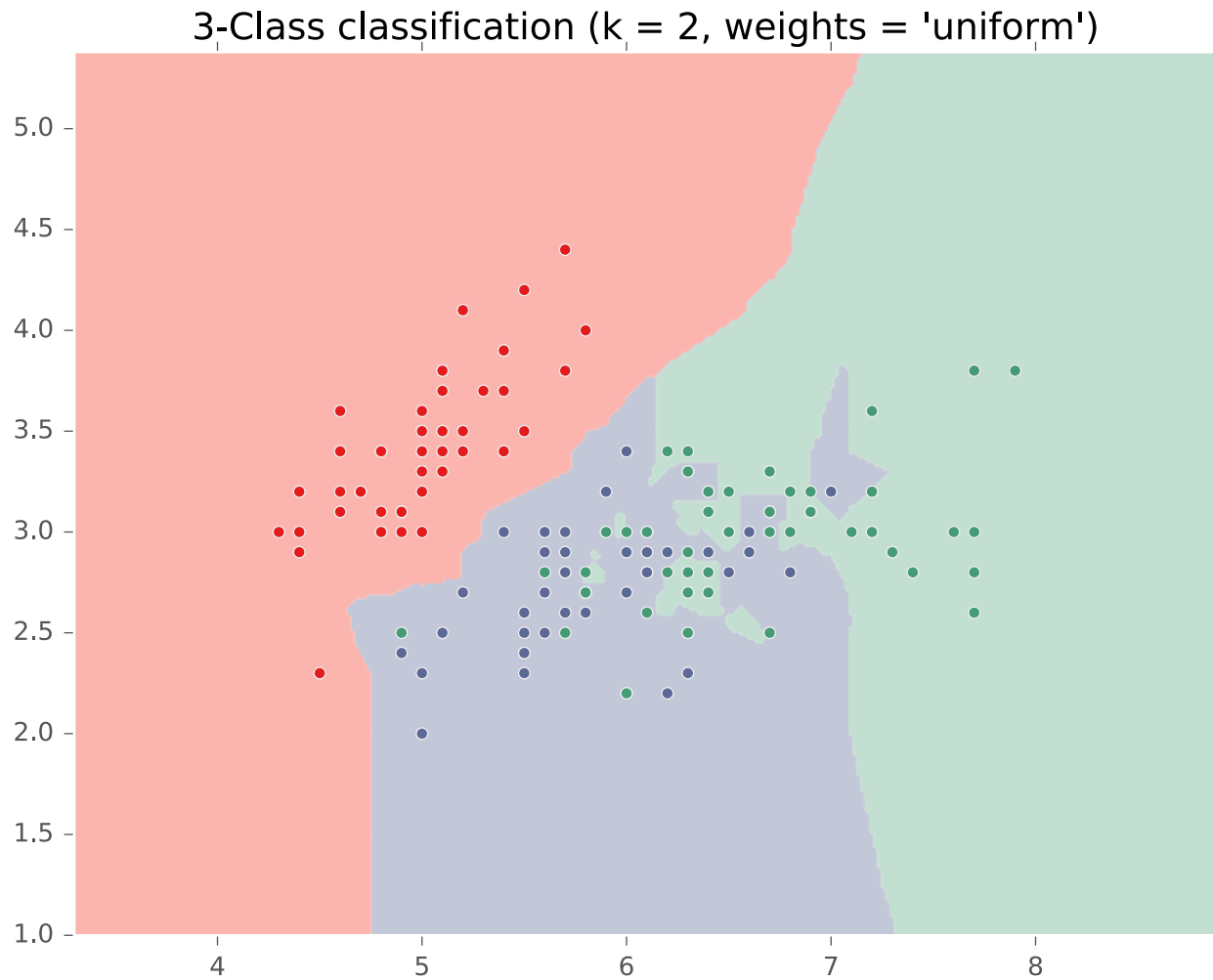


KNN on Fisher Iris Data

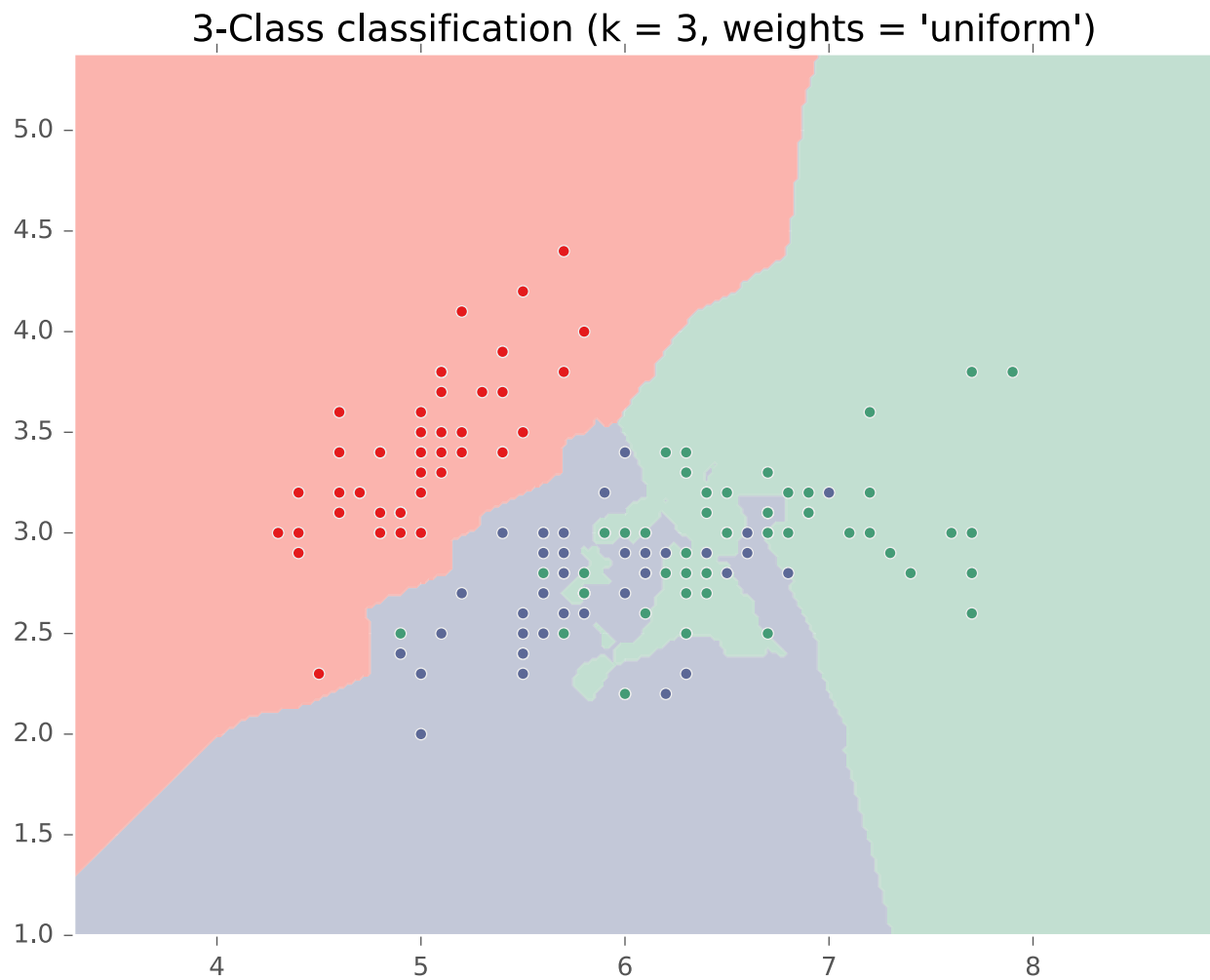
Special Case: Nearest Neighbor



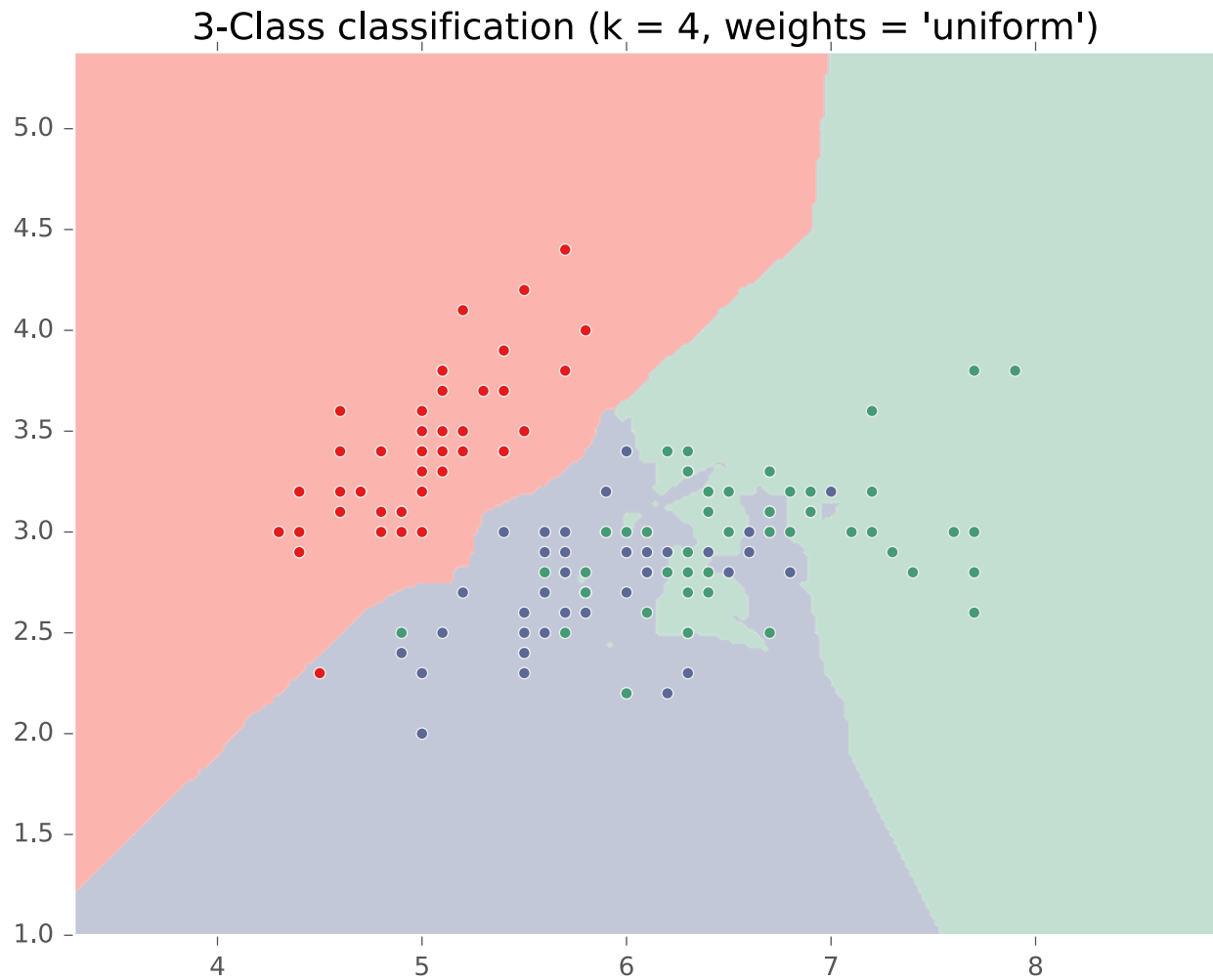
KNN on Fisher Iris Data



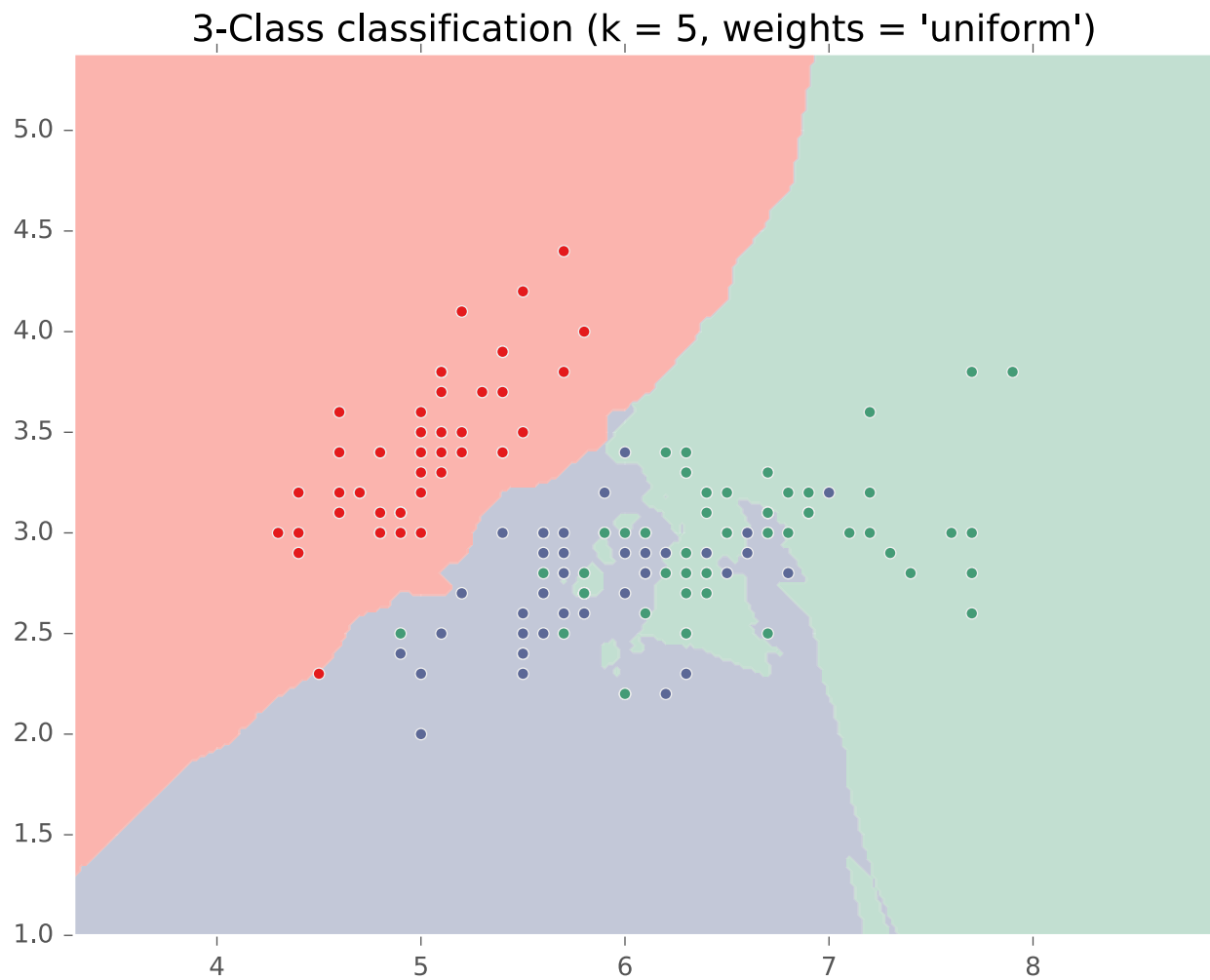
KNN on Fisher Iris Data



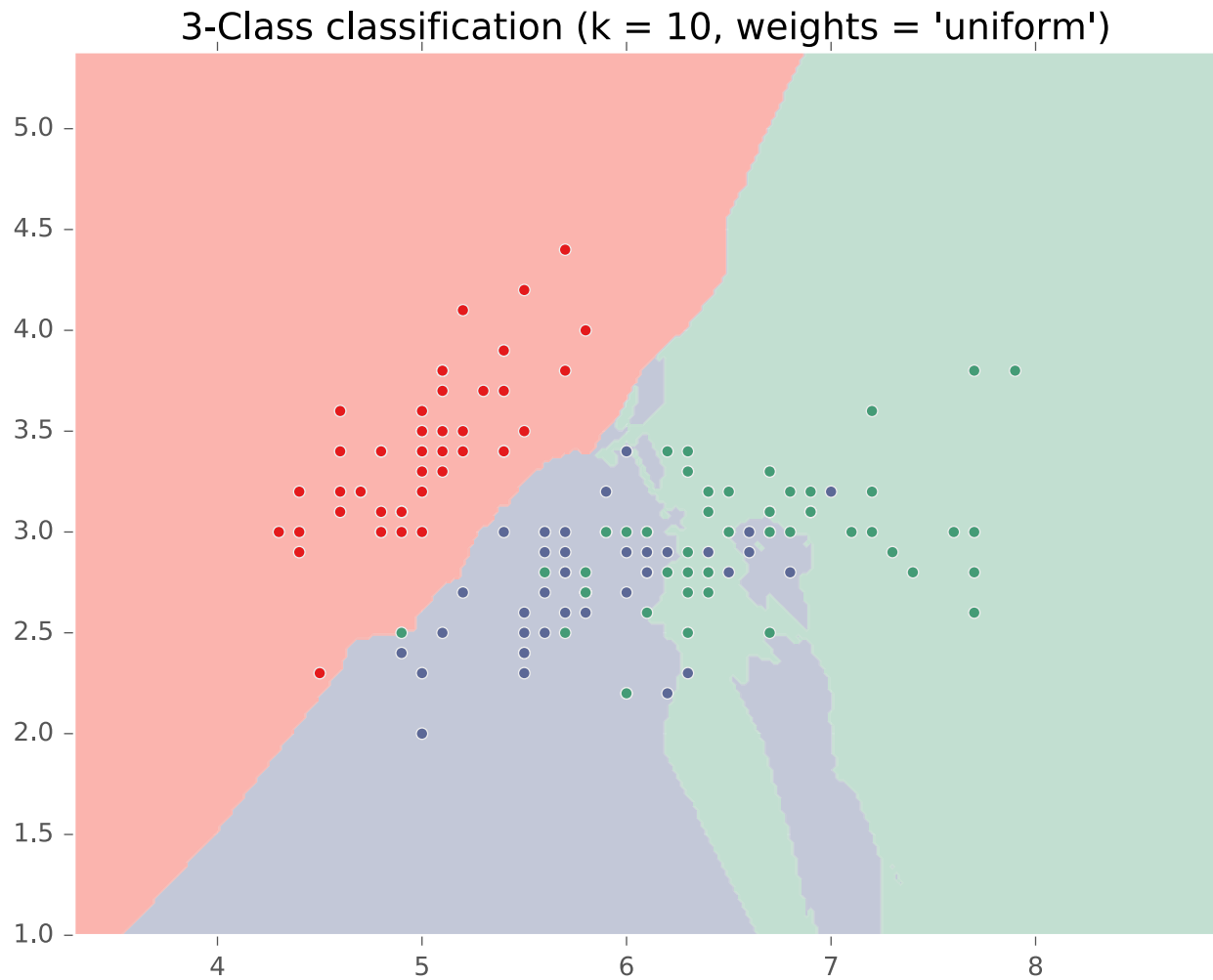
KNN on Fisher Iris Data



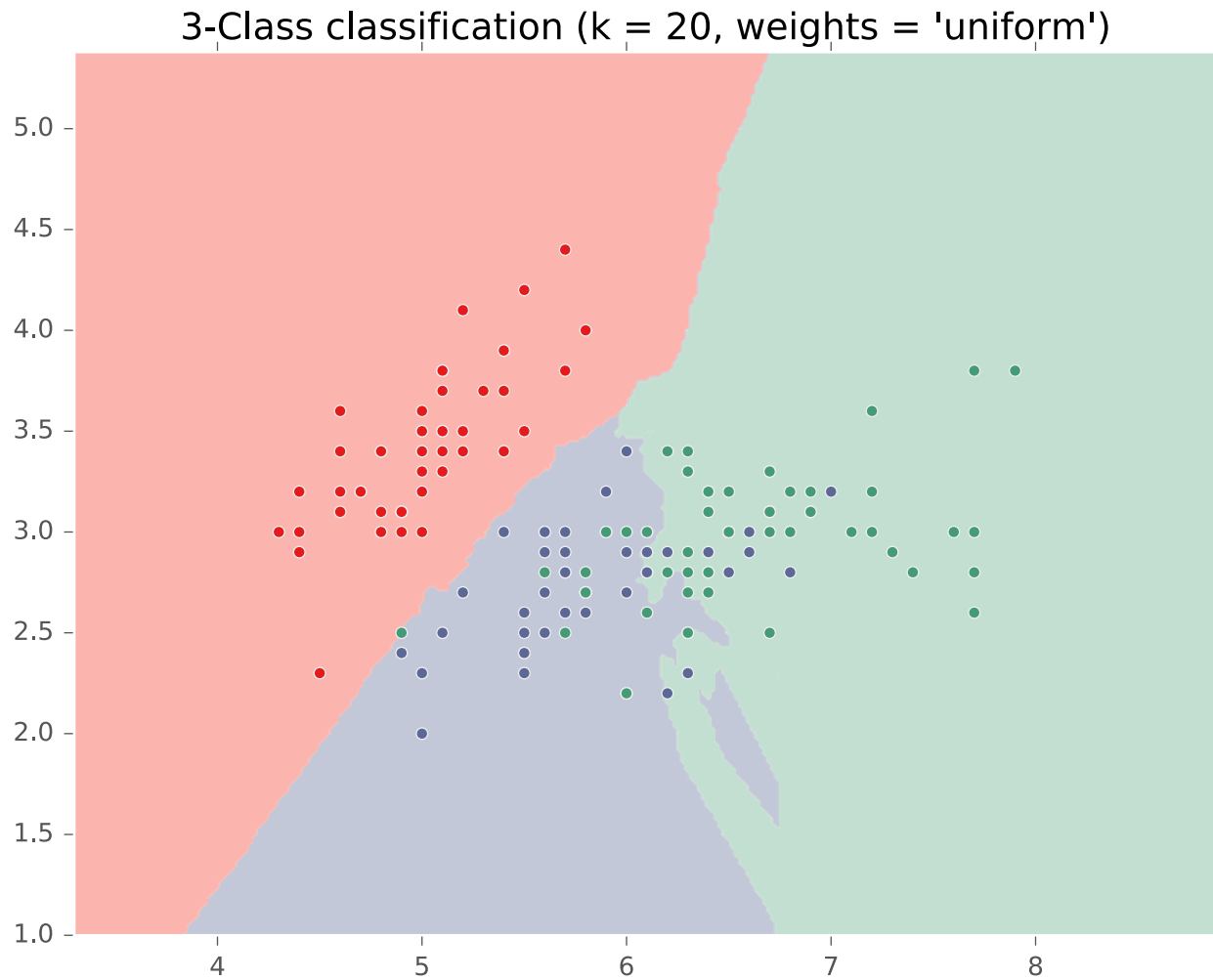
KNN on Fisher Iris Data



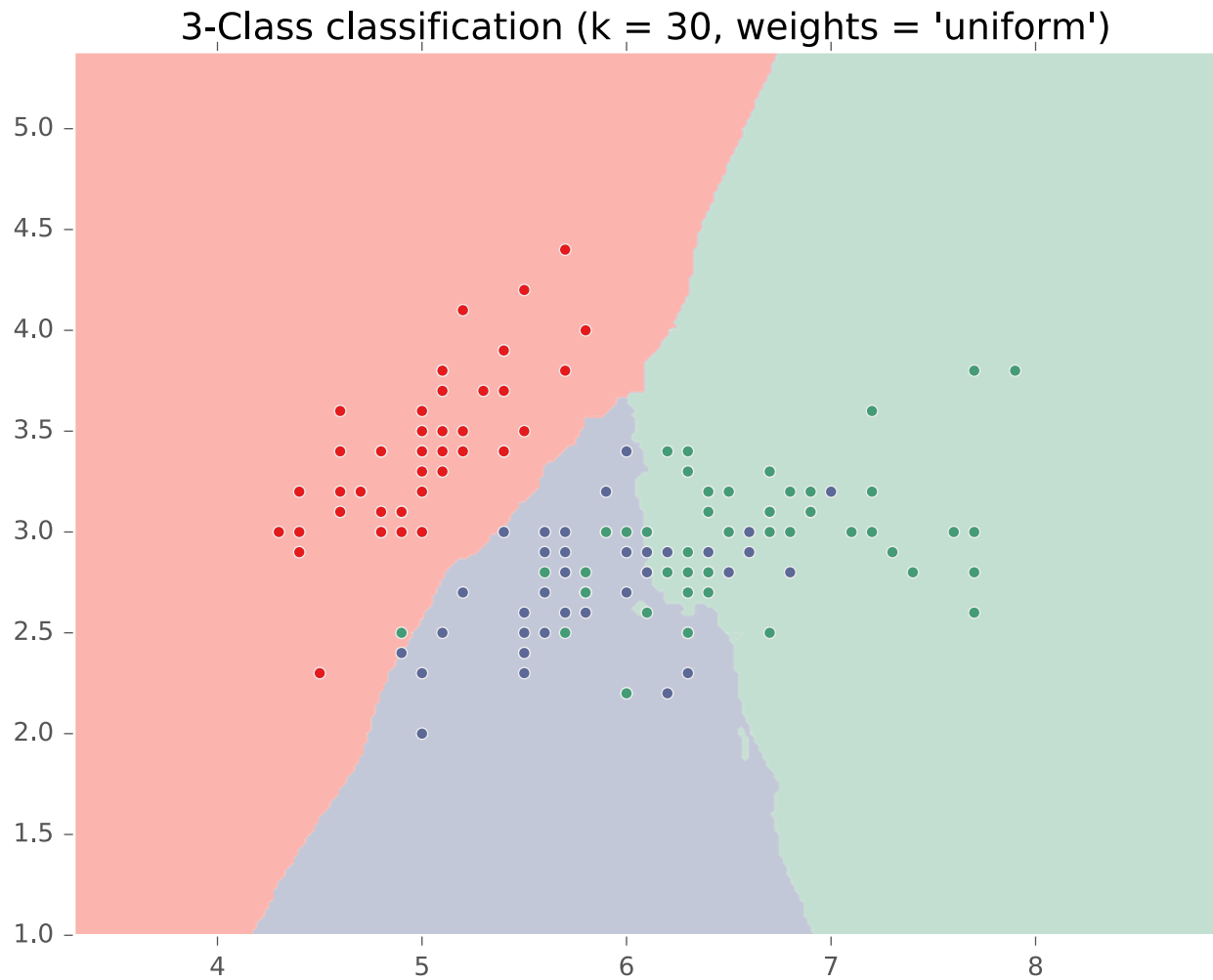
KNN on Fisher Iris Data



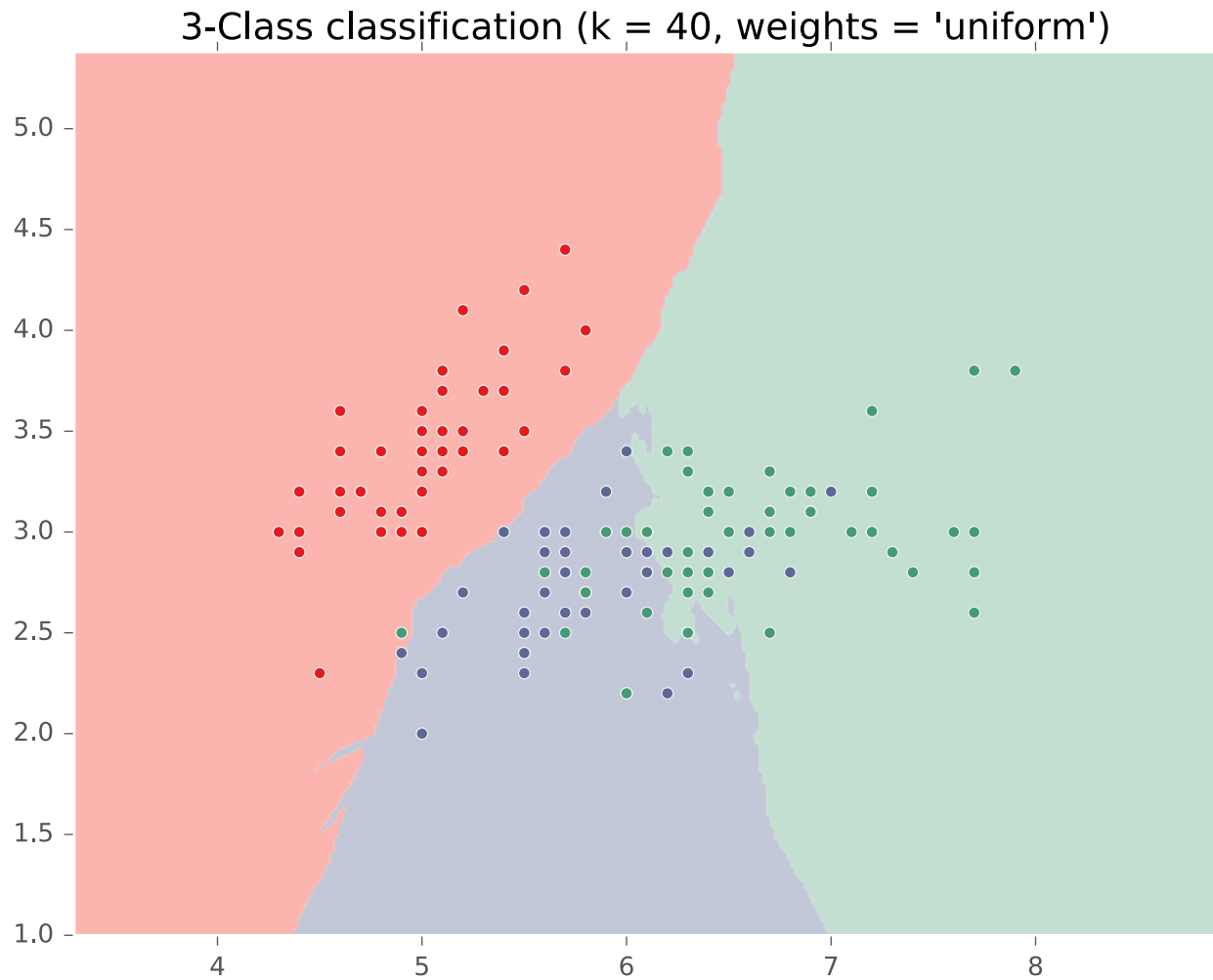
KNN on Fisher Iris Data



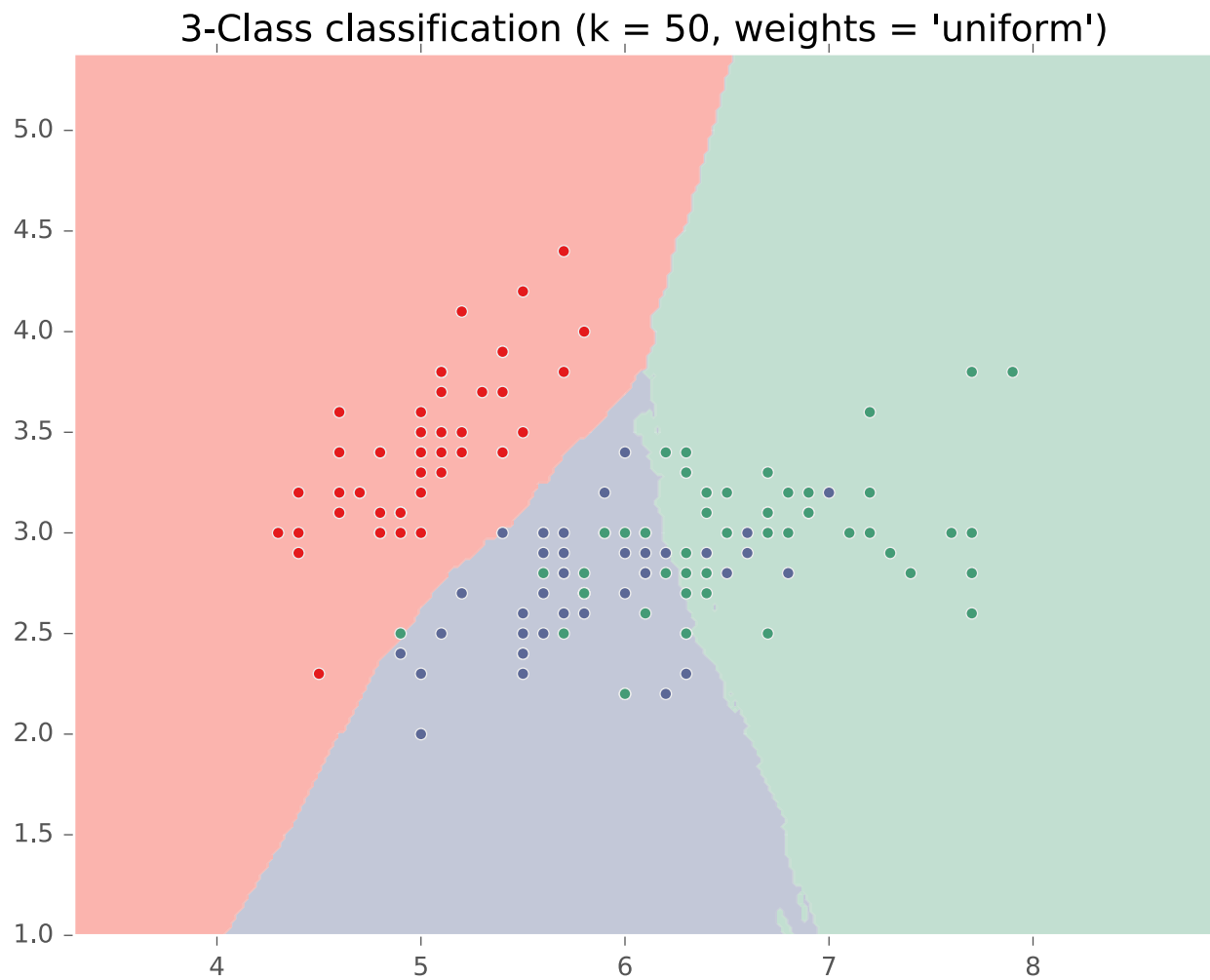
KNN on Fisher Iris Data



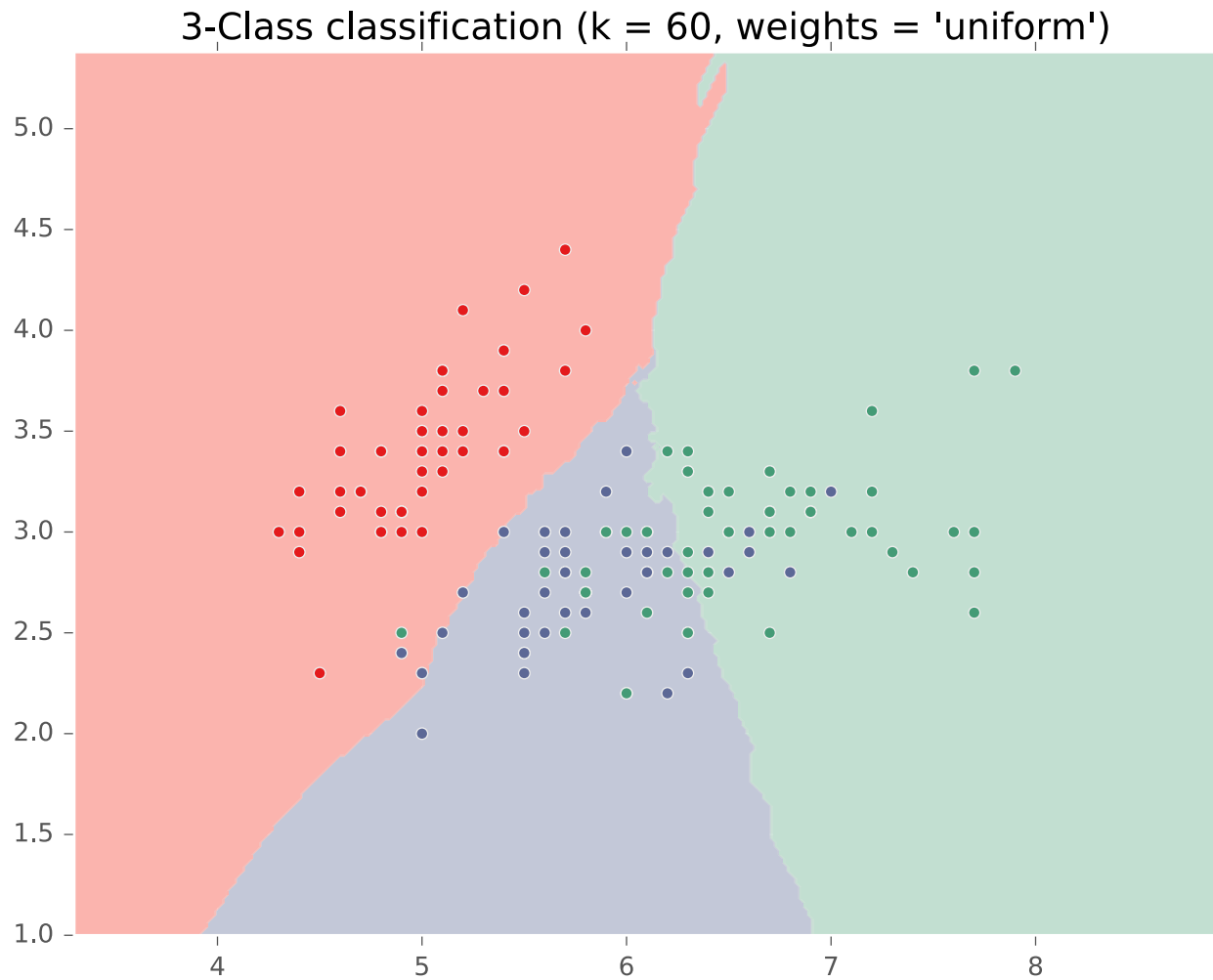
KNN on Fisher Iris Data



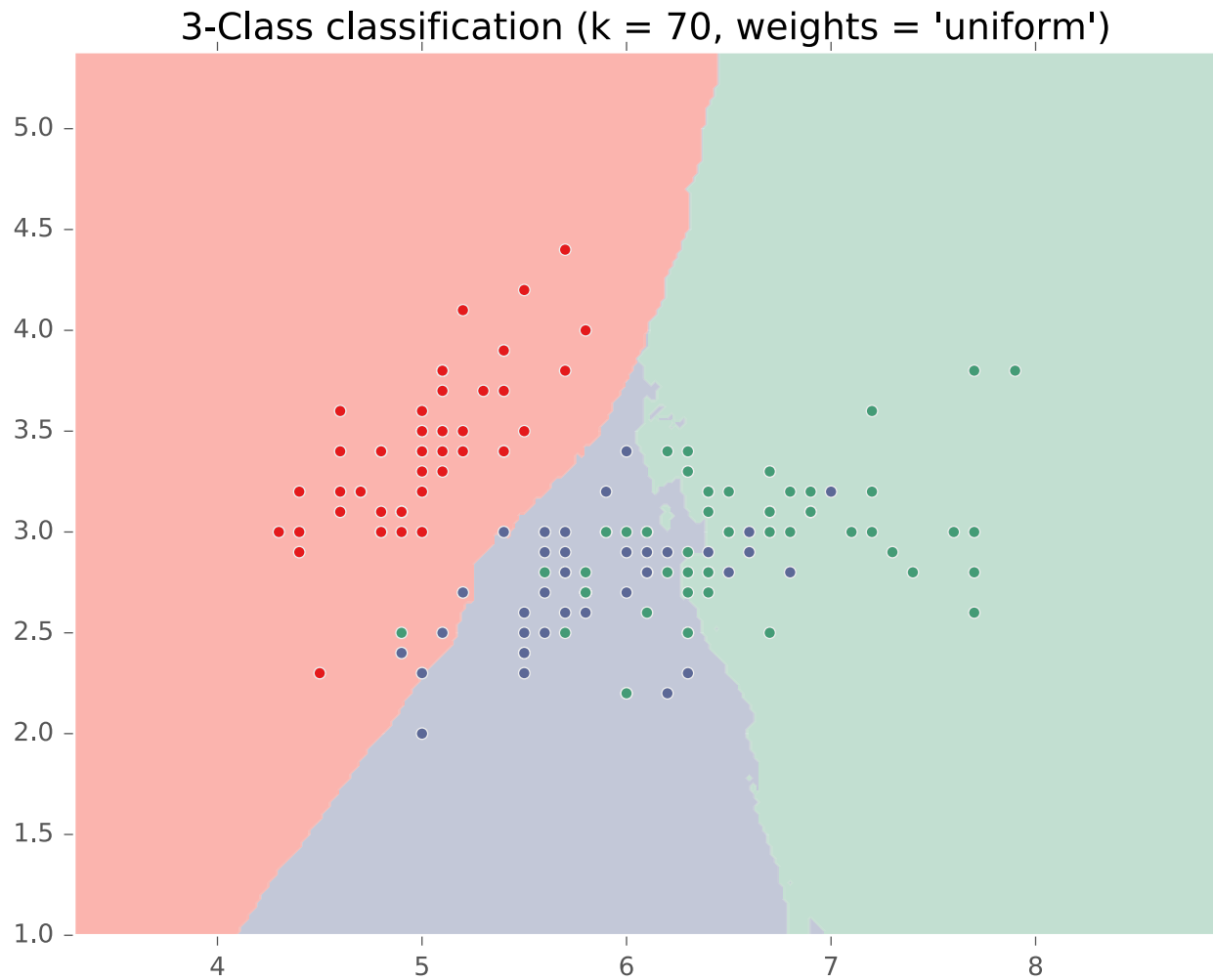
KNN on Fisher Iris Data



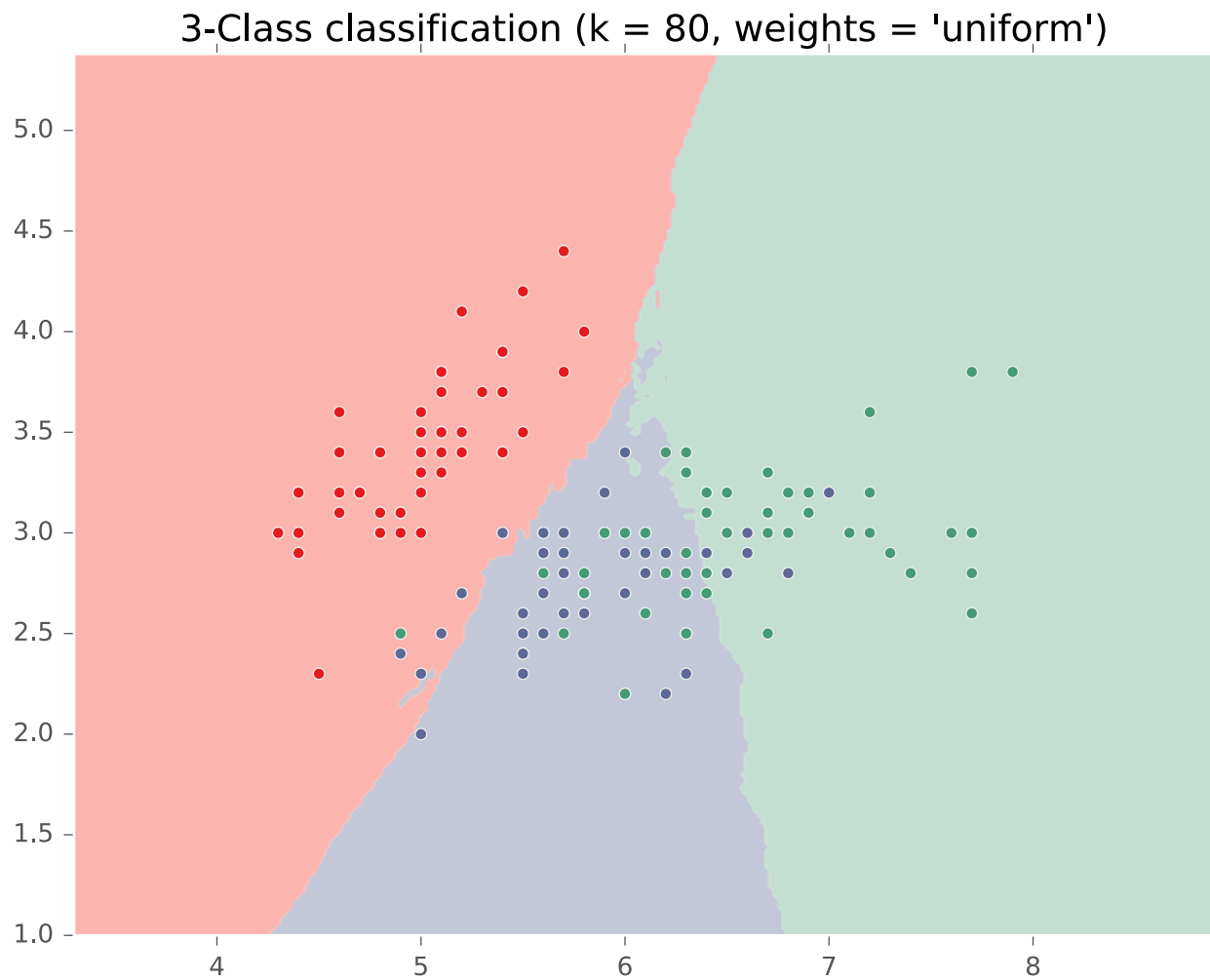
KNN on Fisher Iris Data



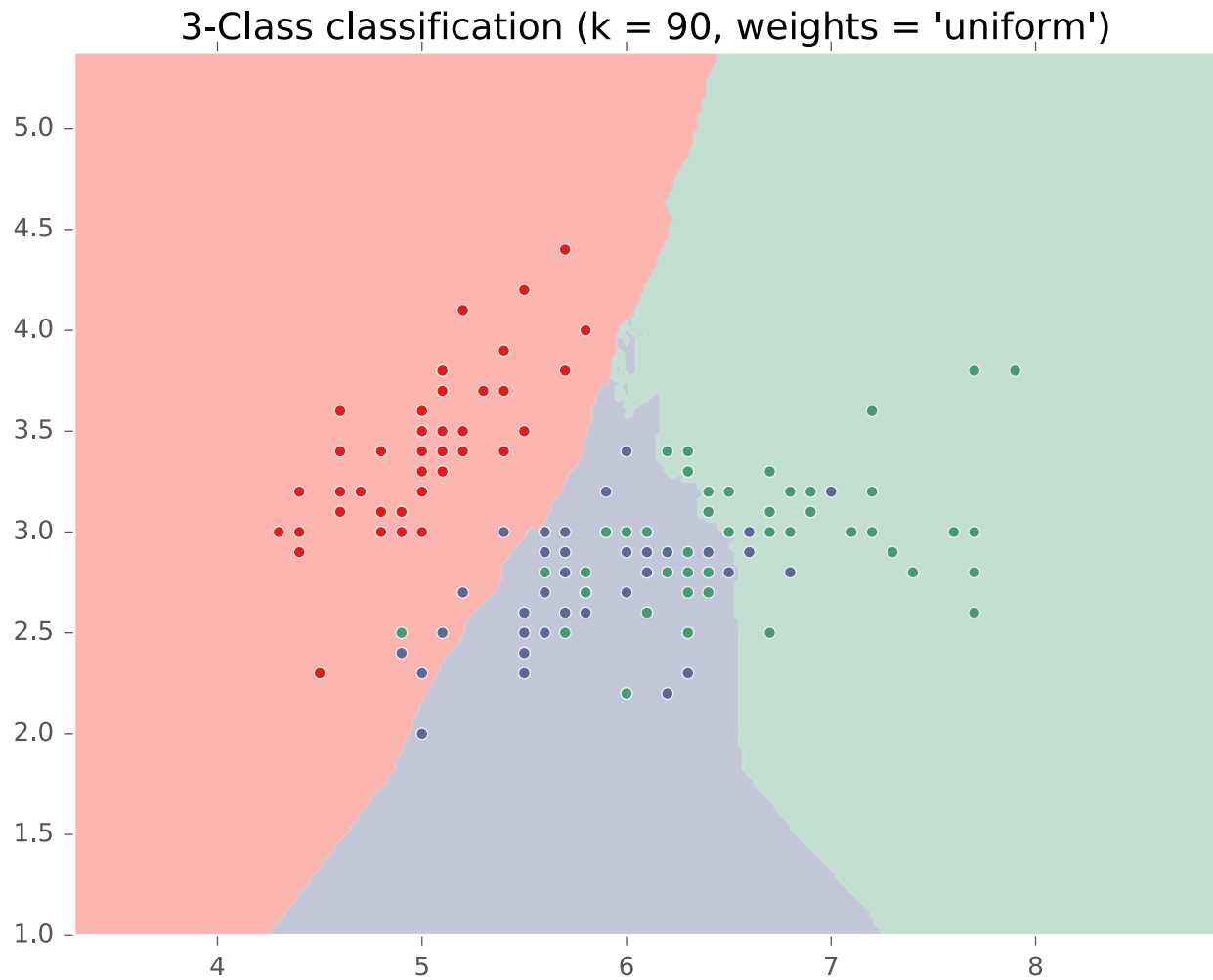
KNN on Fisher Iris Data



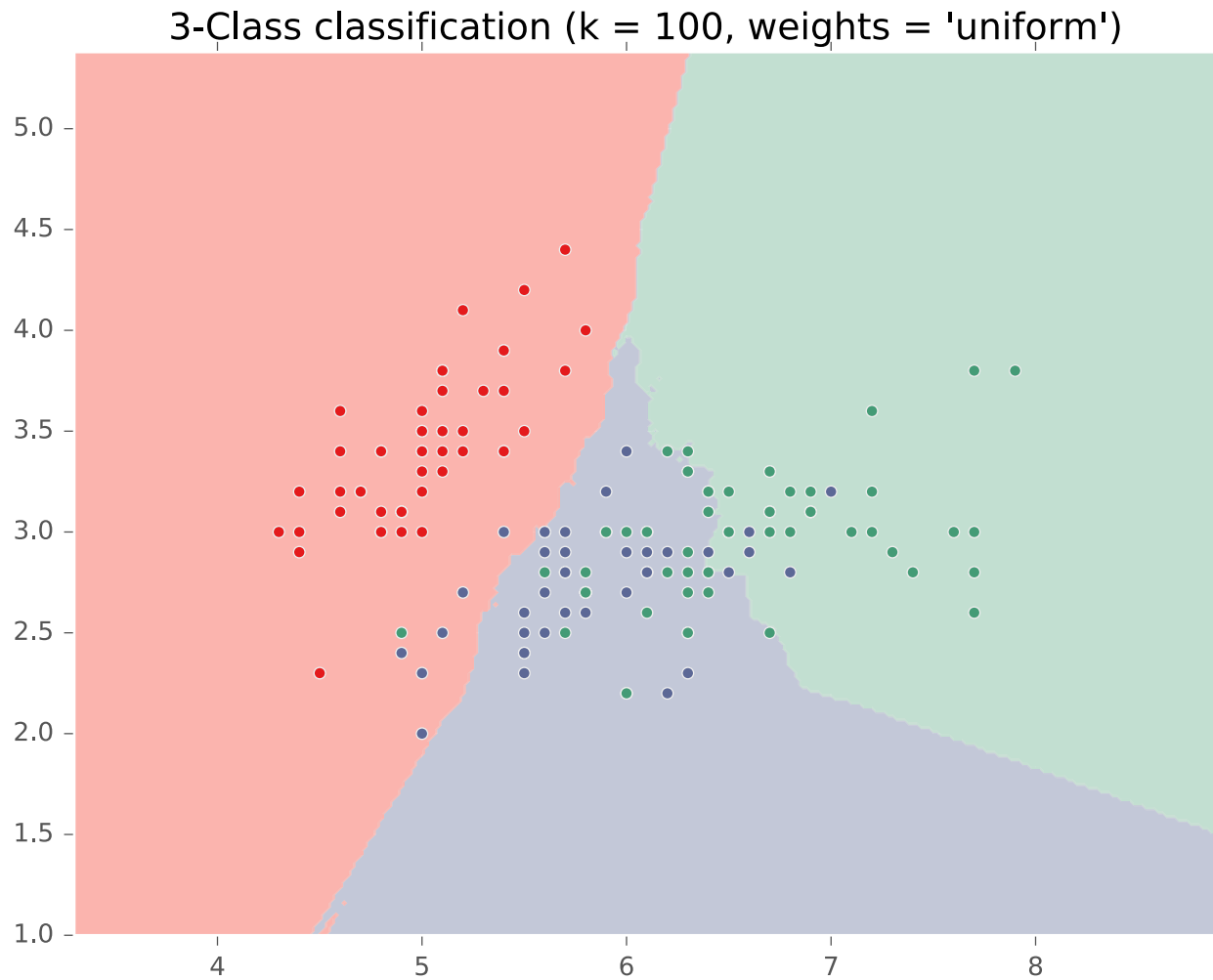
KNN on Fisher Iris Data



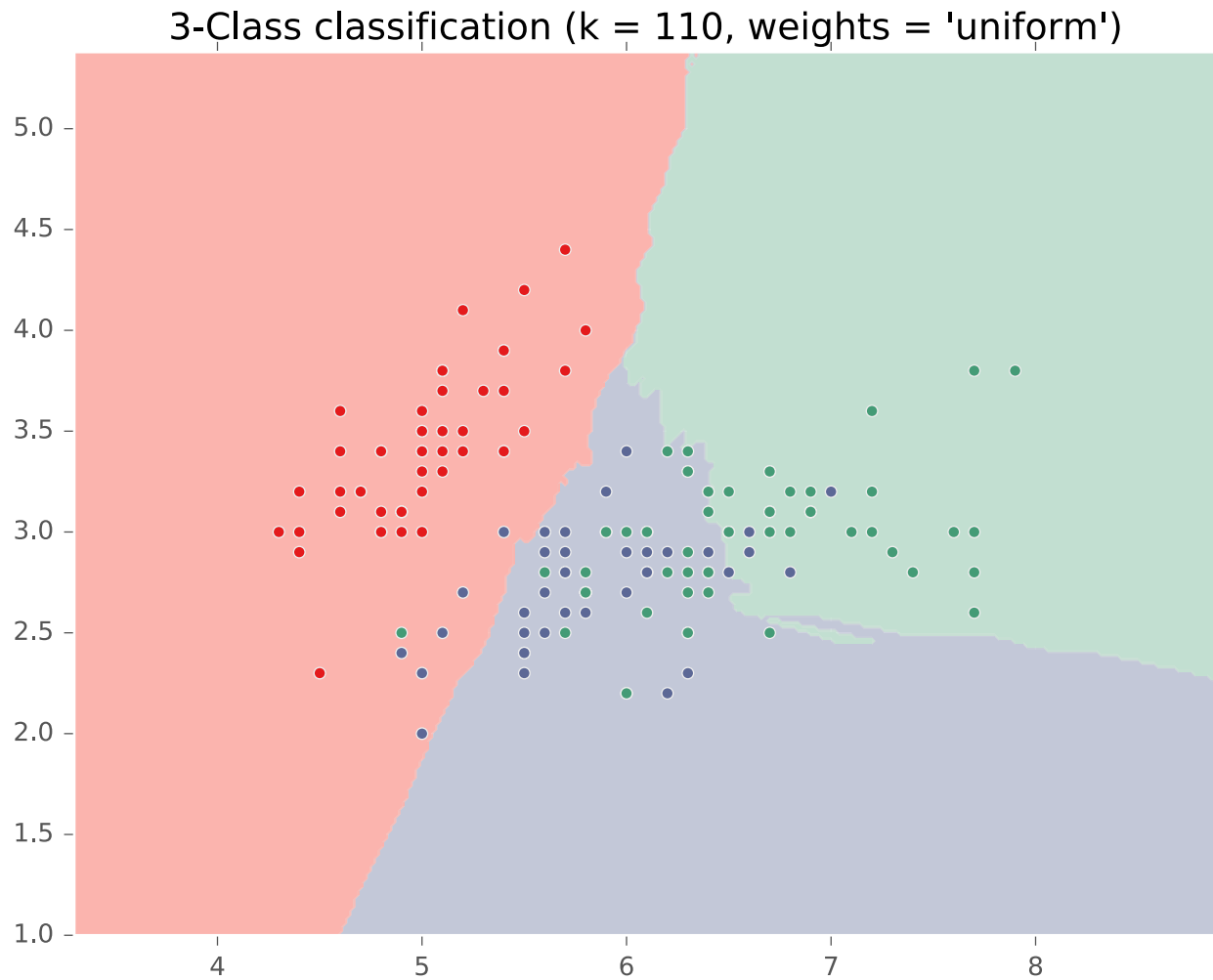
KNN on Fisher Iris Data



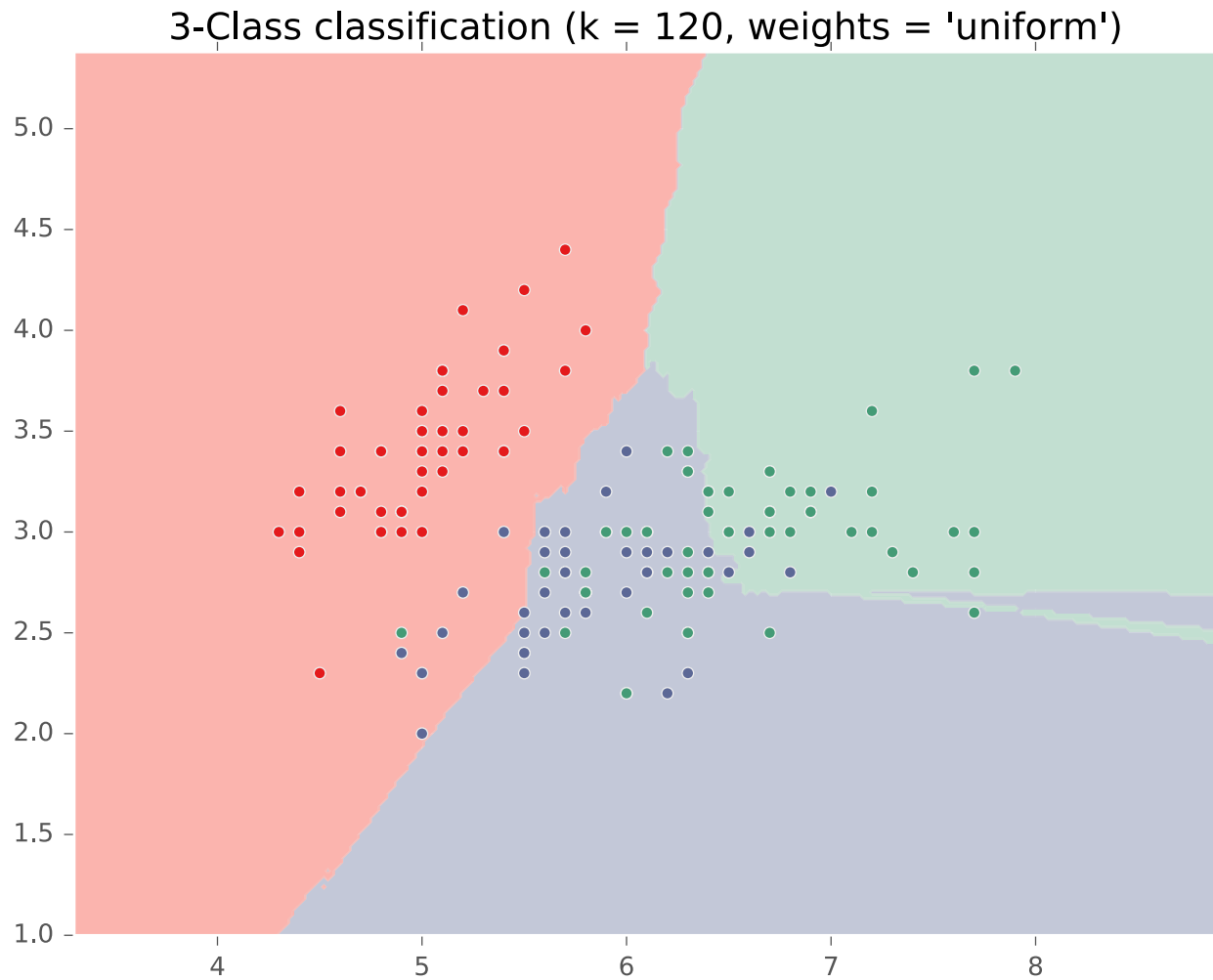
KNN on Fisher Iris Data



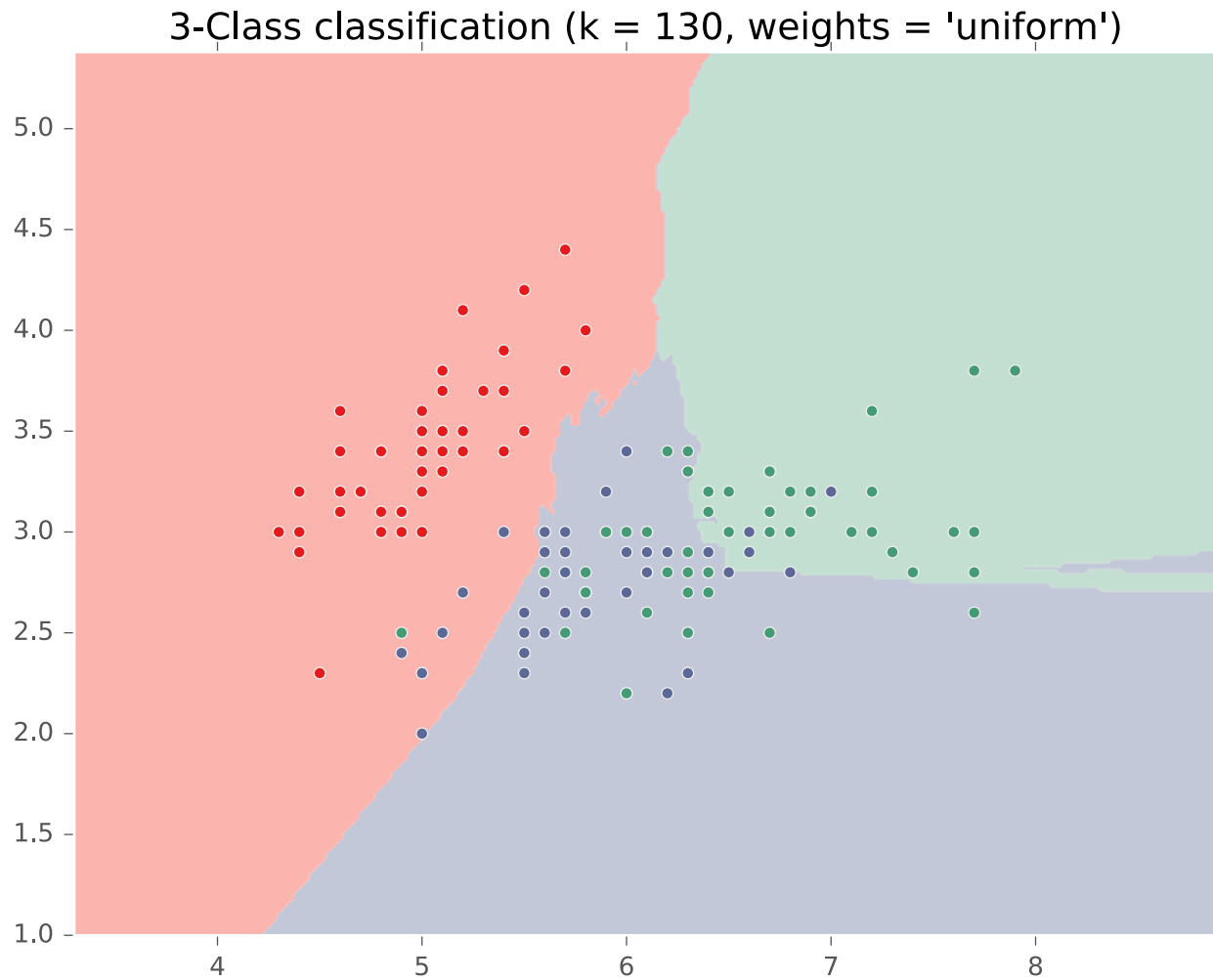
KNN on Fisher Iris Data



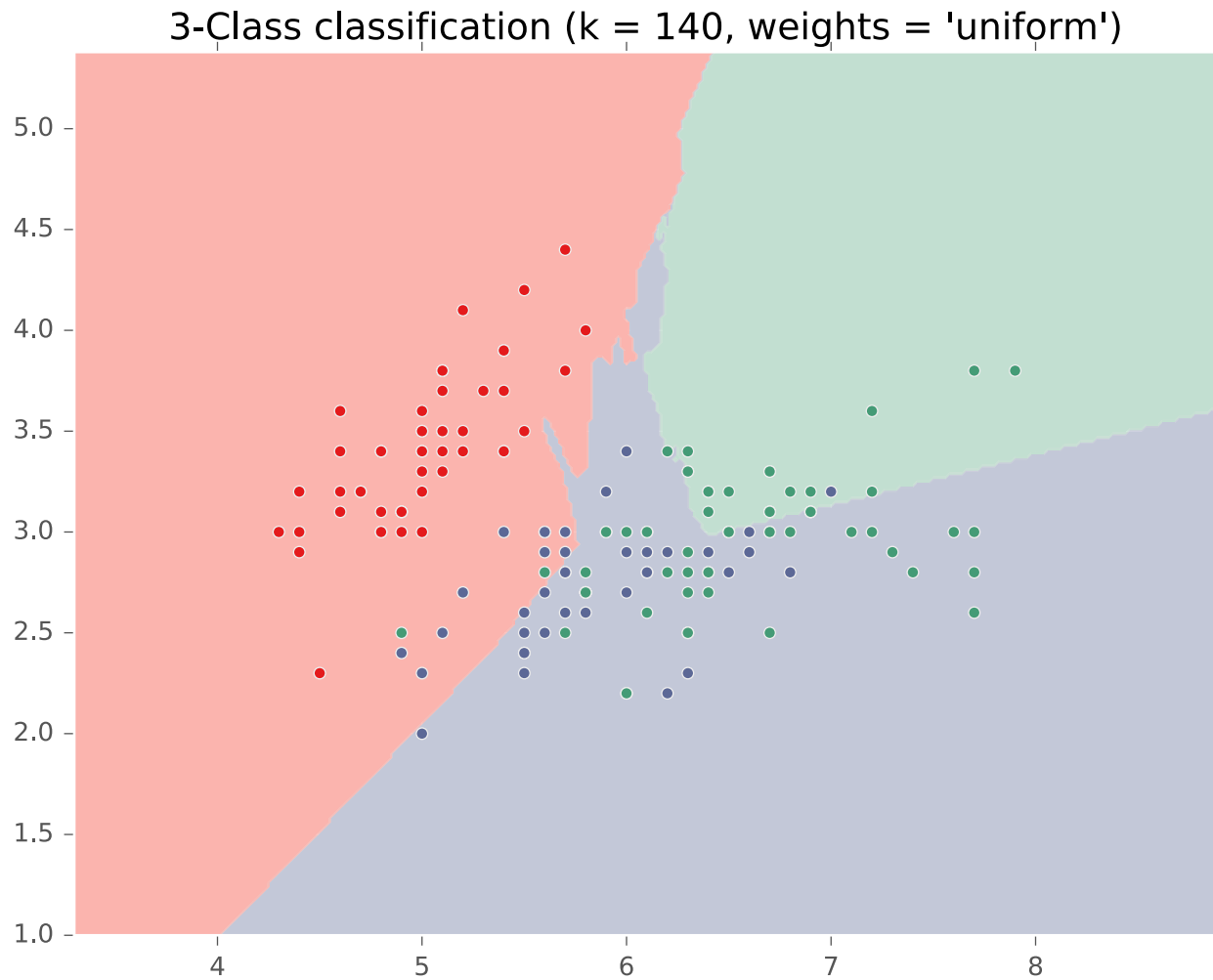
KNN on Fisher Iris Data



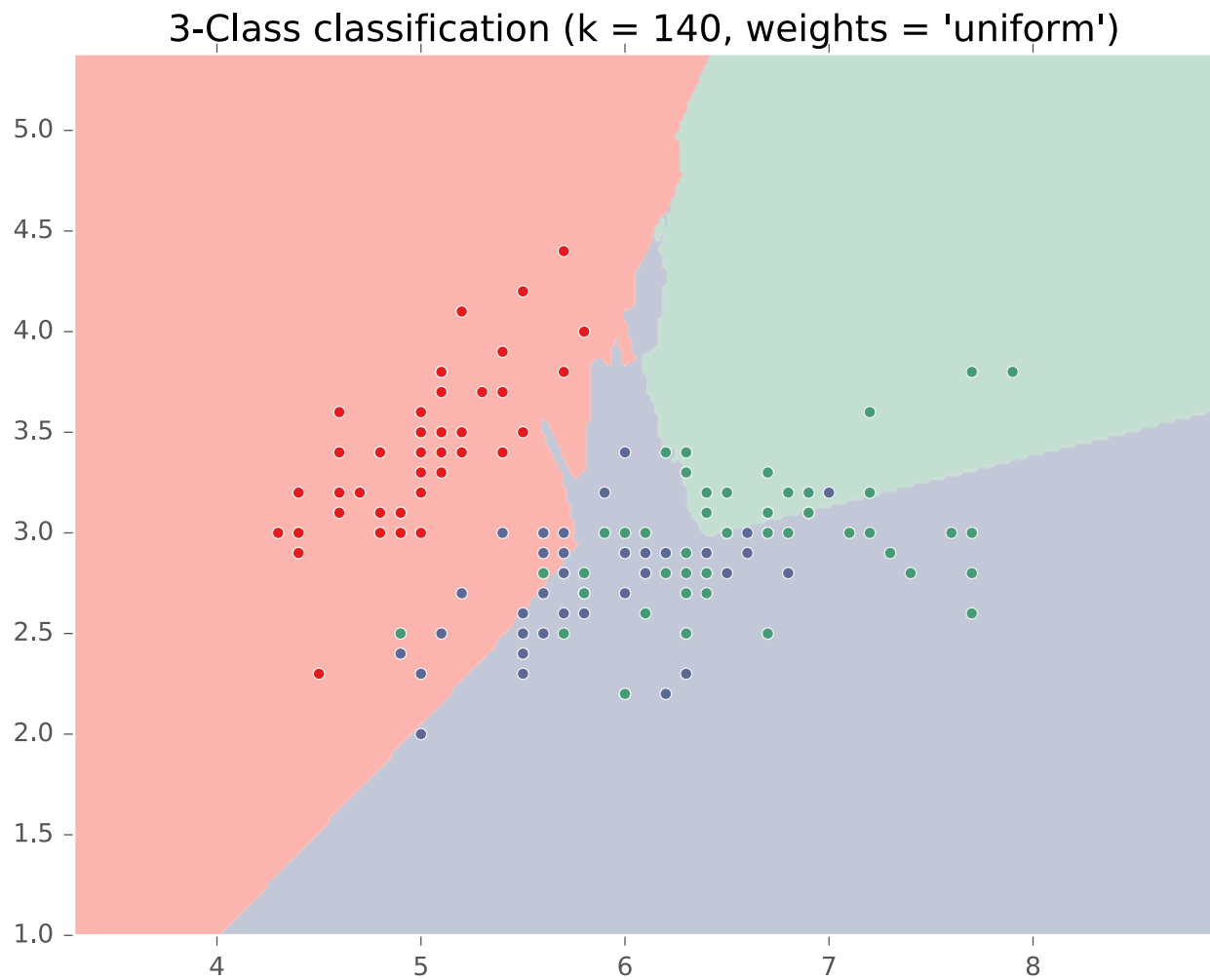
KNN on Fisher Iris Data



KNN on Fisher Iris Data

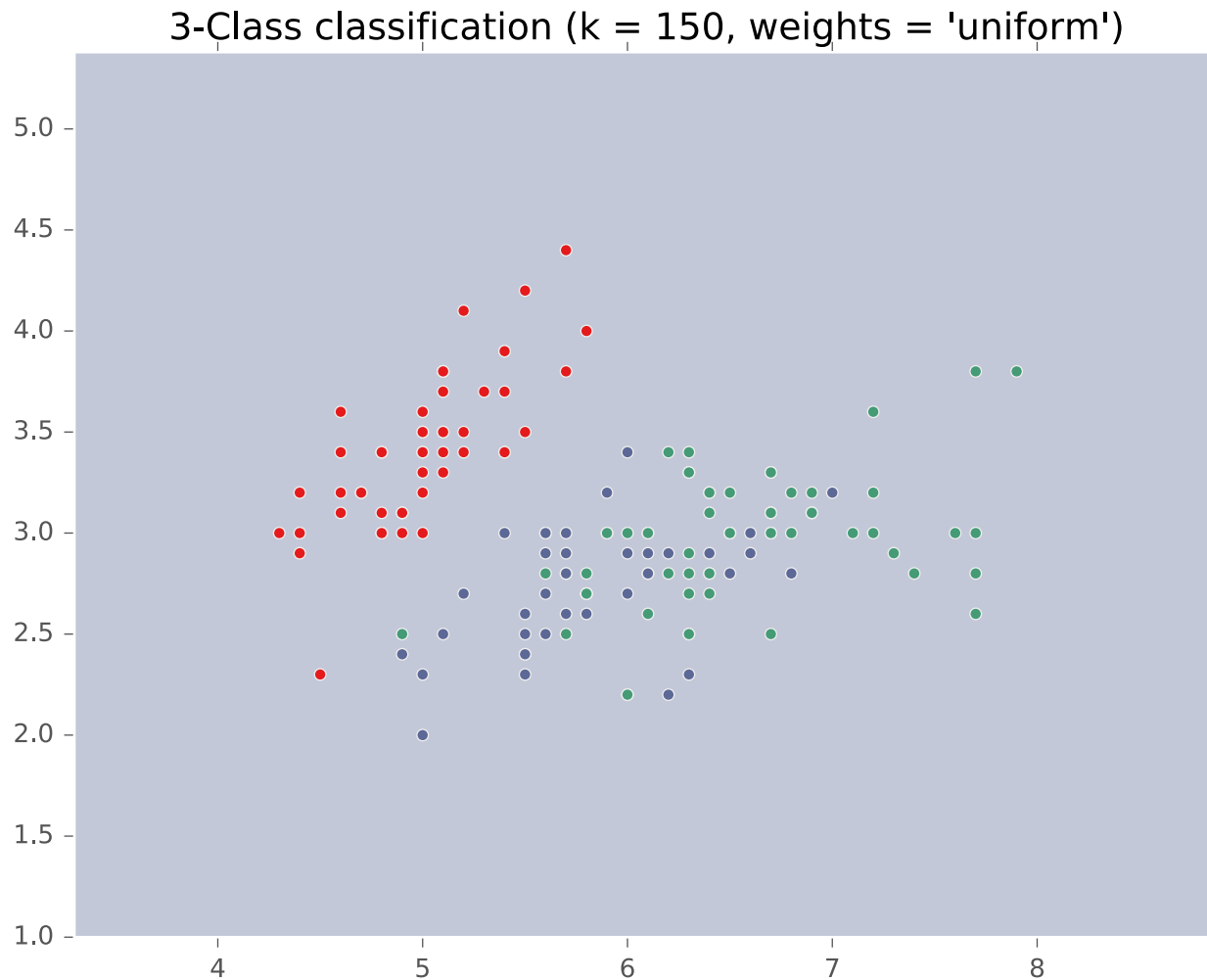


KNN on Fisher Iris Data



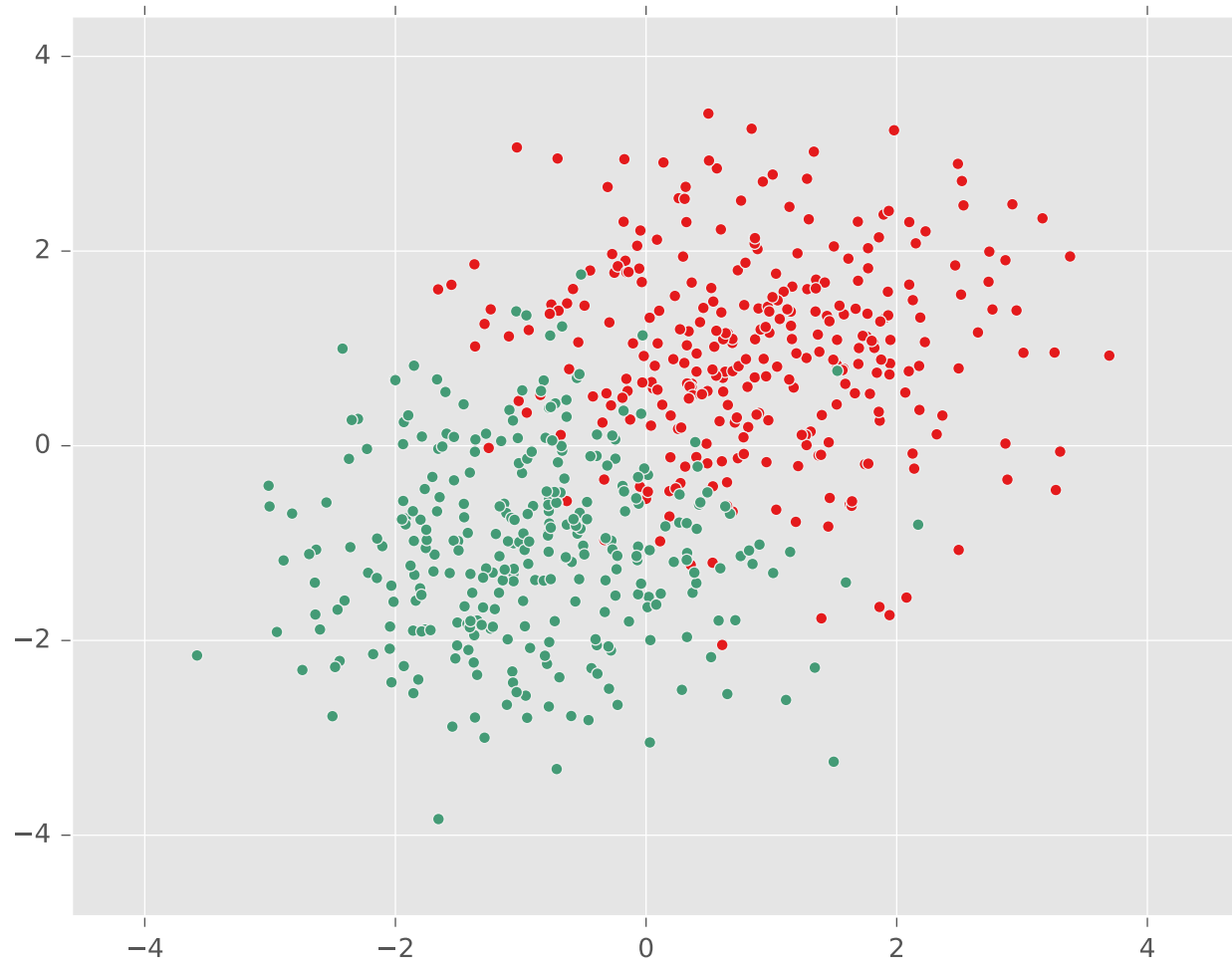
KNN on Fisher Iris Data

Special Case: Majority Vote

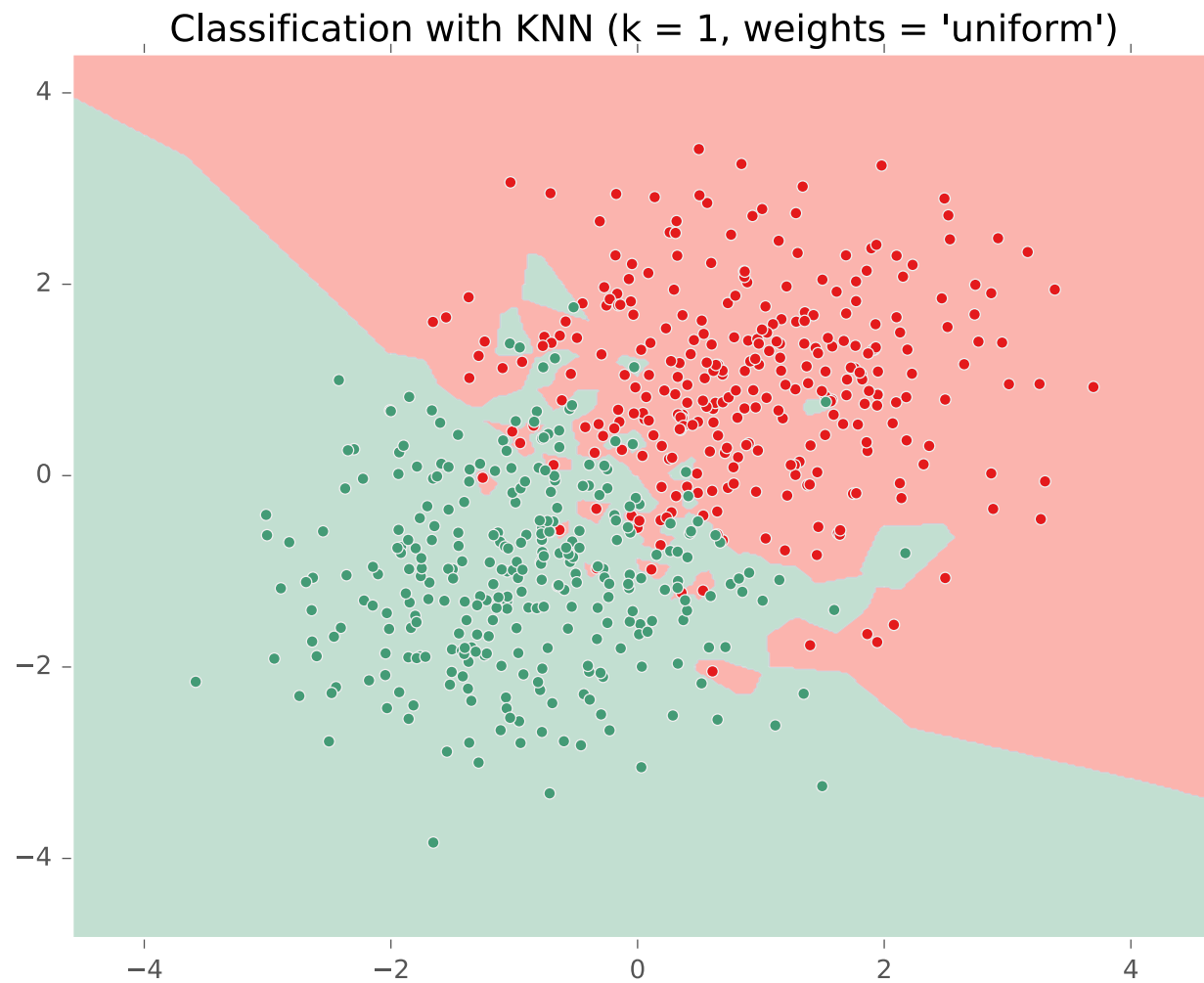


KNN ON GAUSSIAN DATA

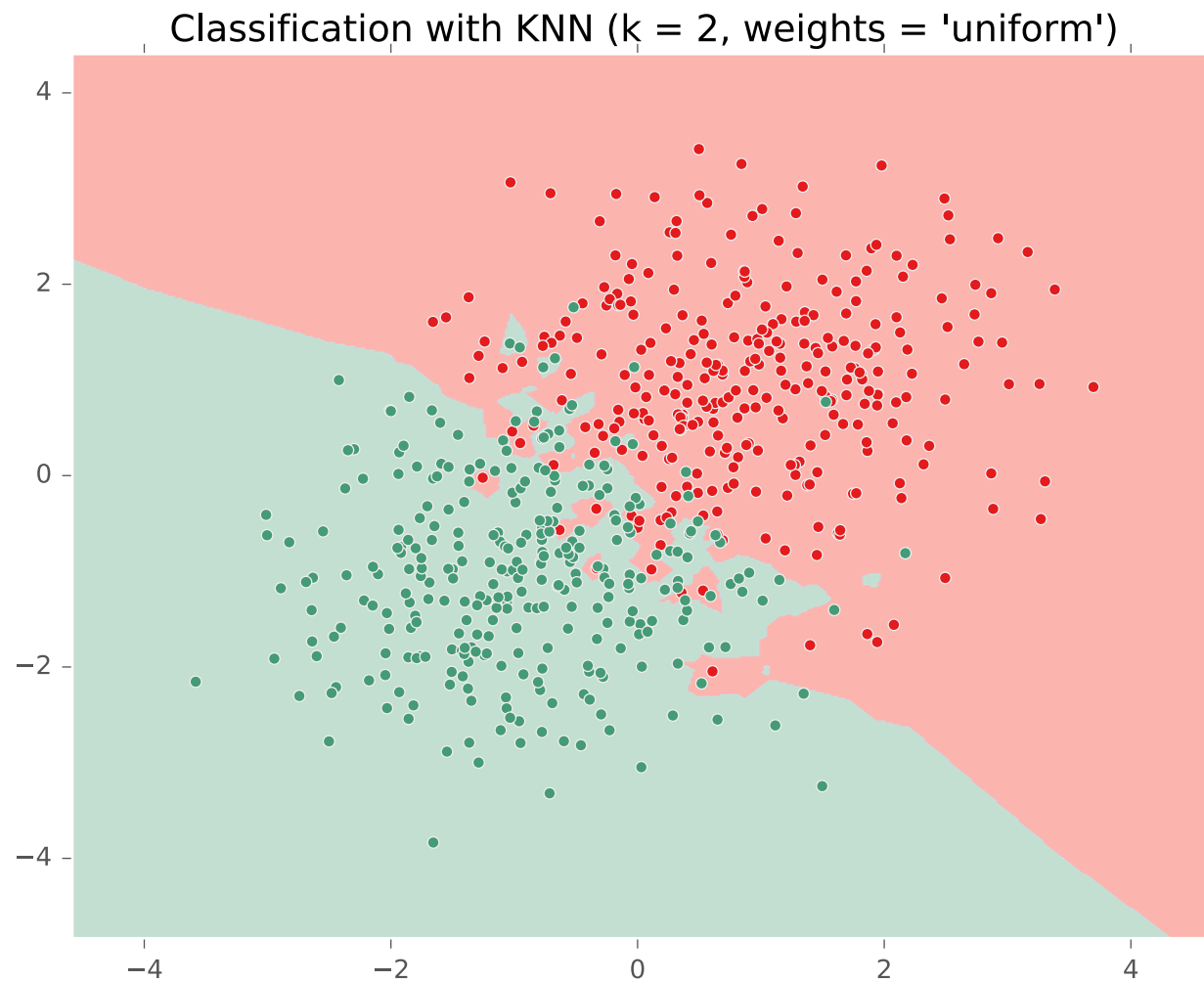
KNN on Gaussian Data



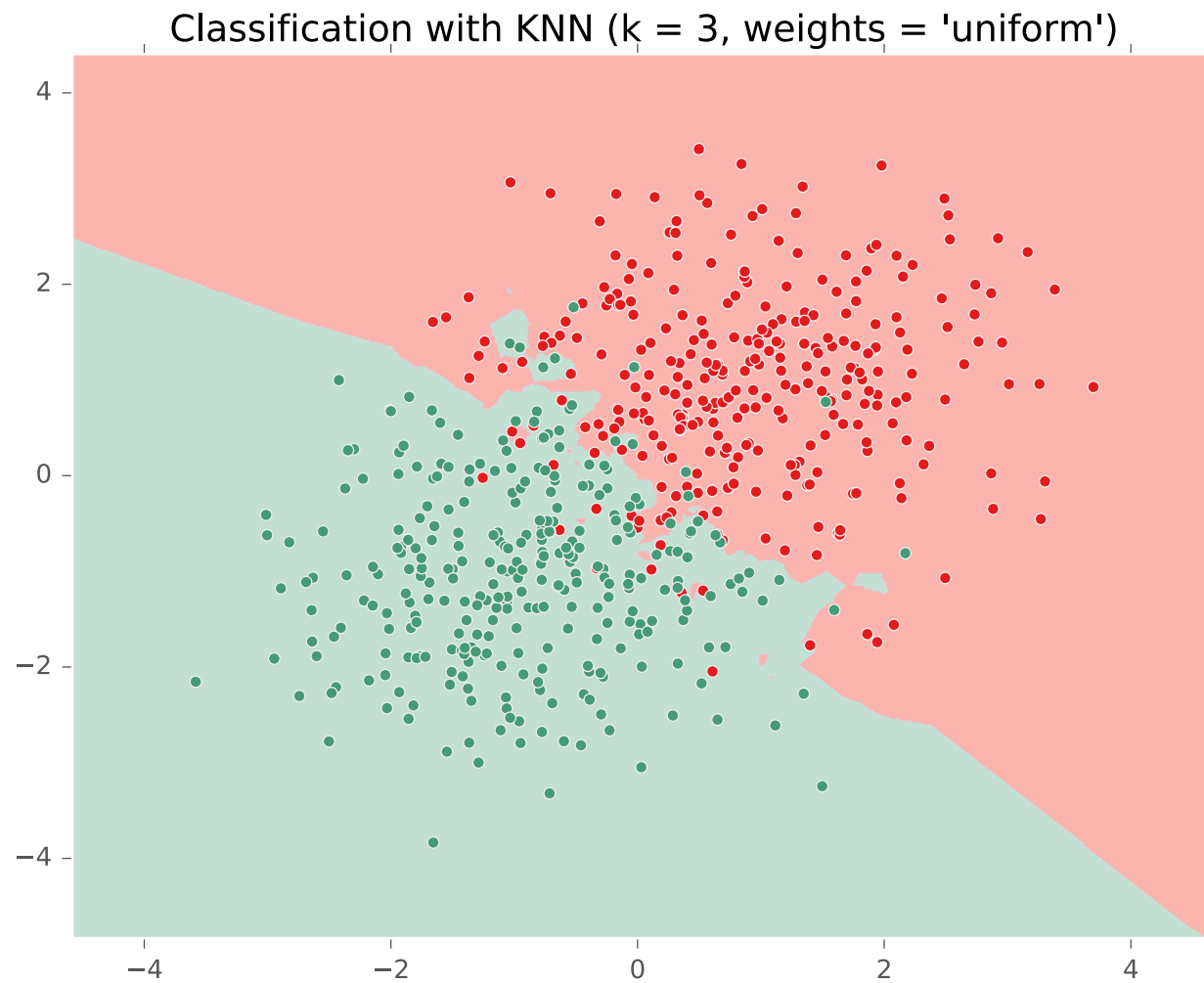
KNN on Gaussian Data



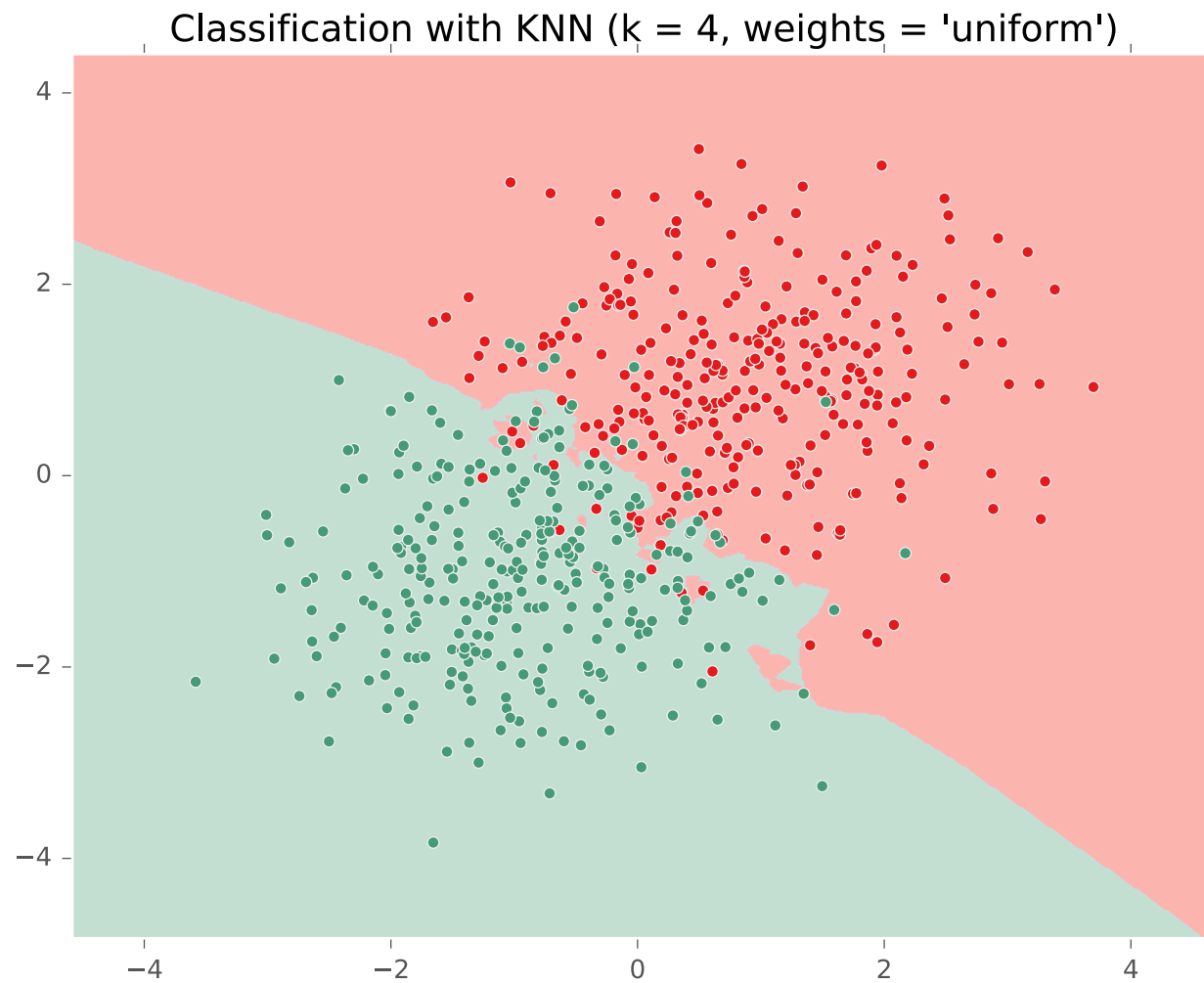
KNN on Gaussian Data



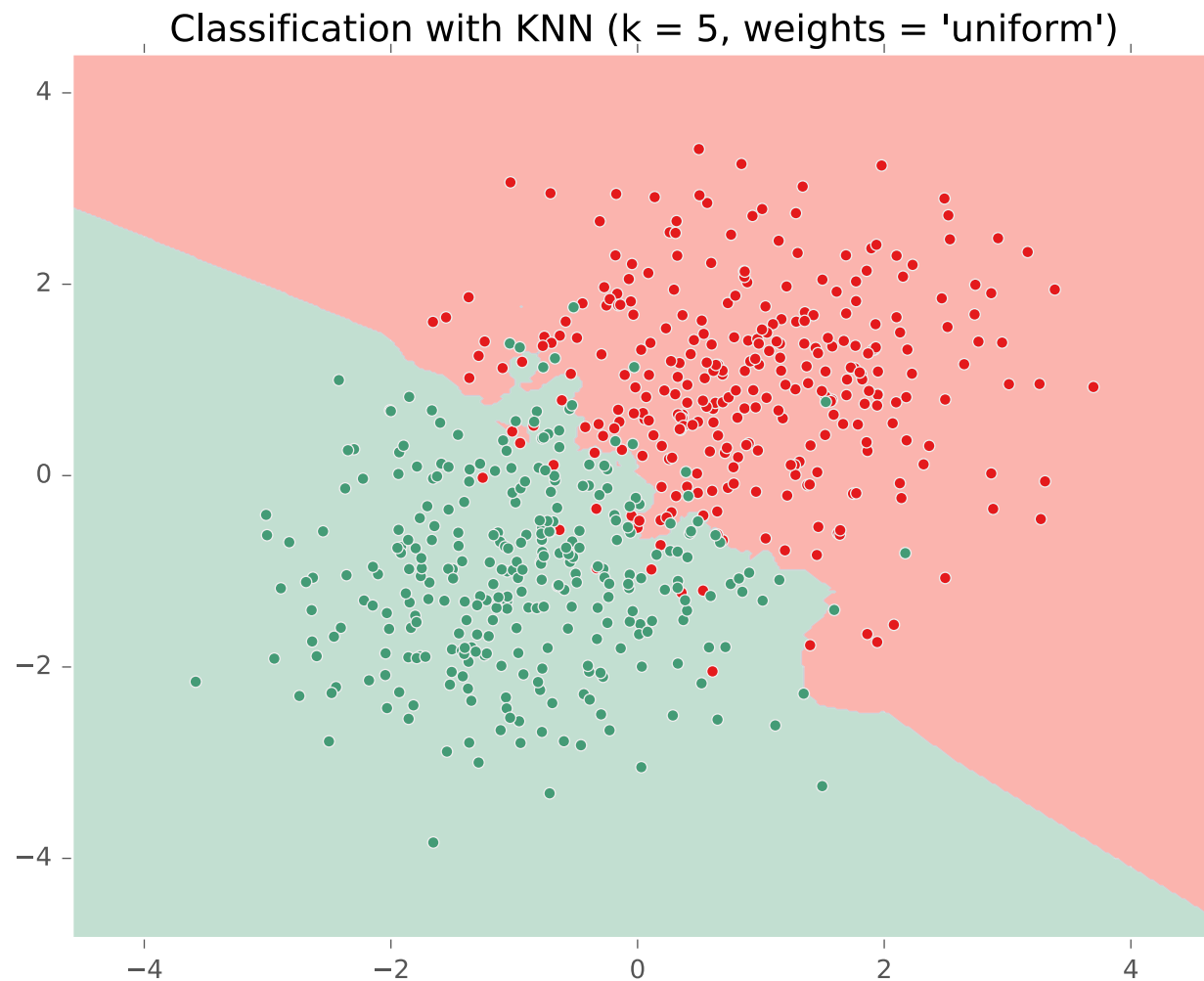
KNN on Gaussian Data



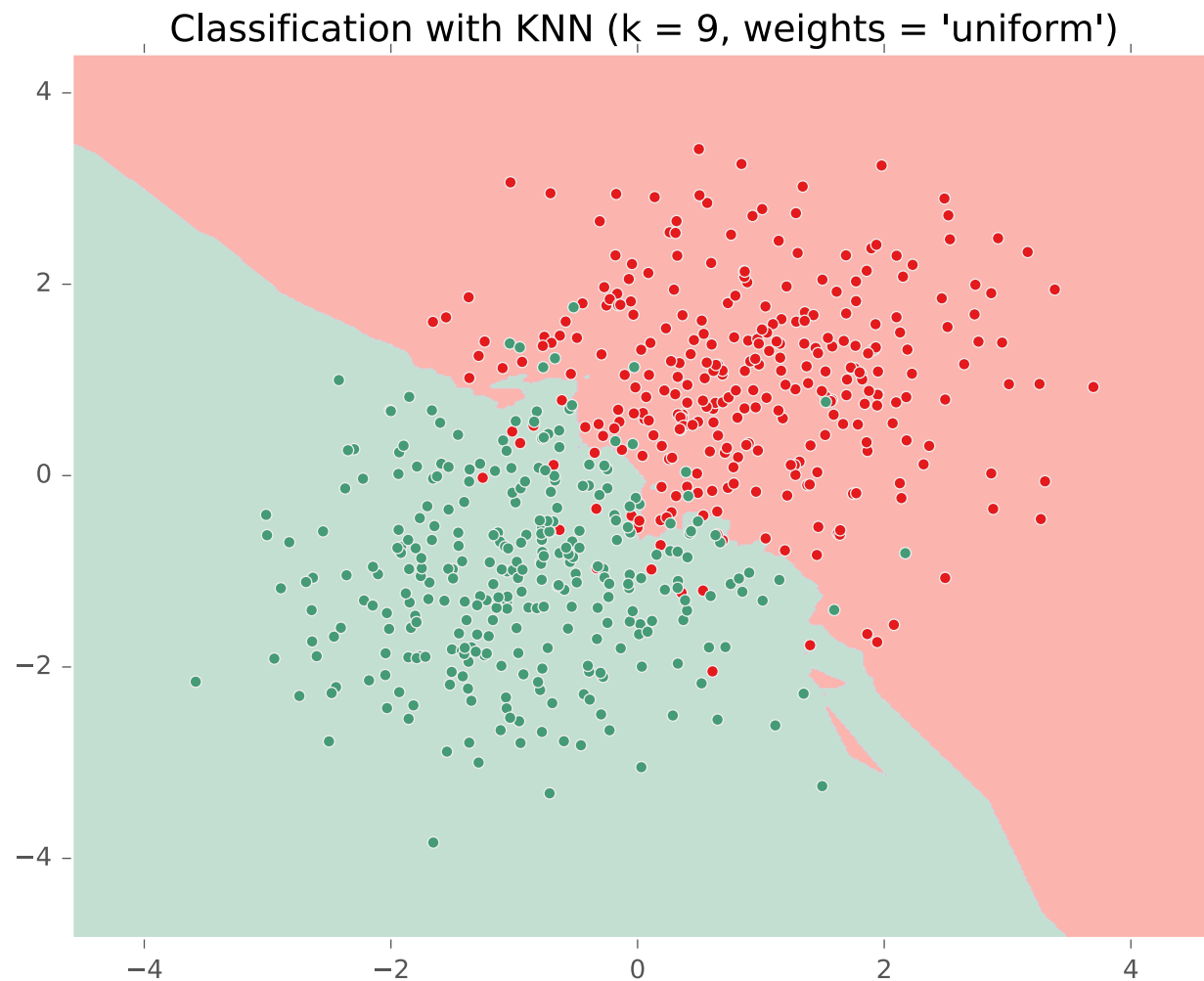
KNN on Gaussian Data



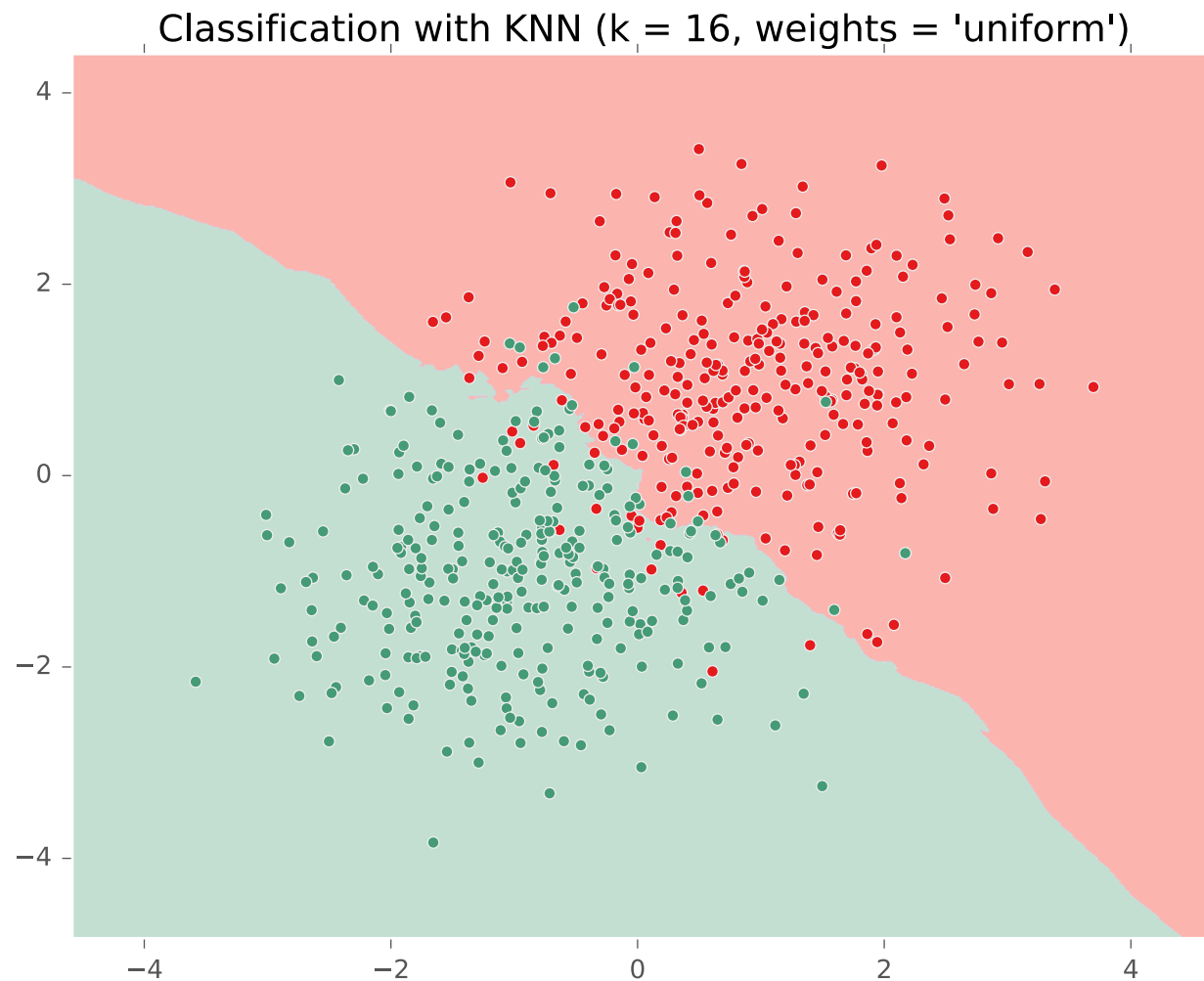
KNN on Gaussian Data



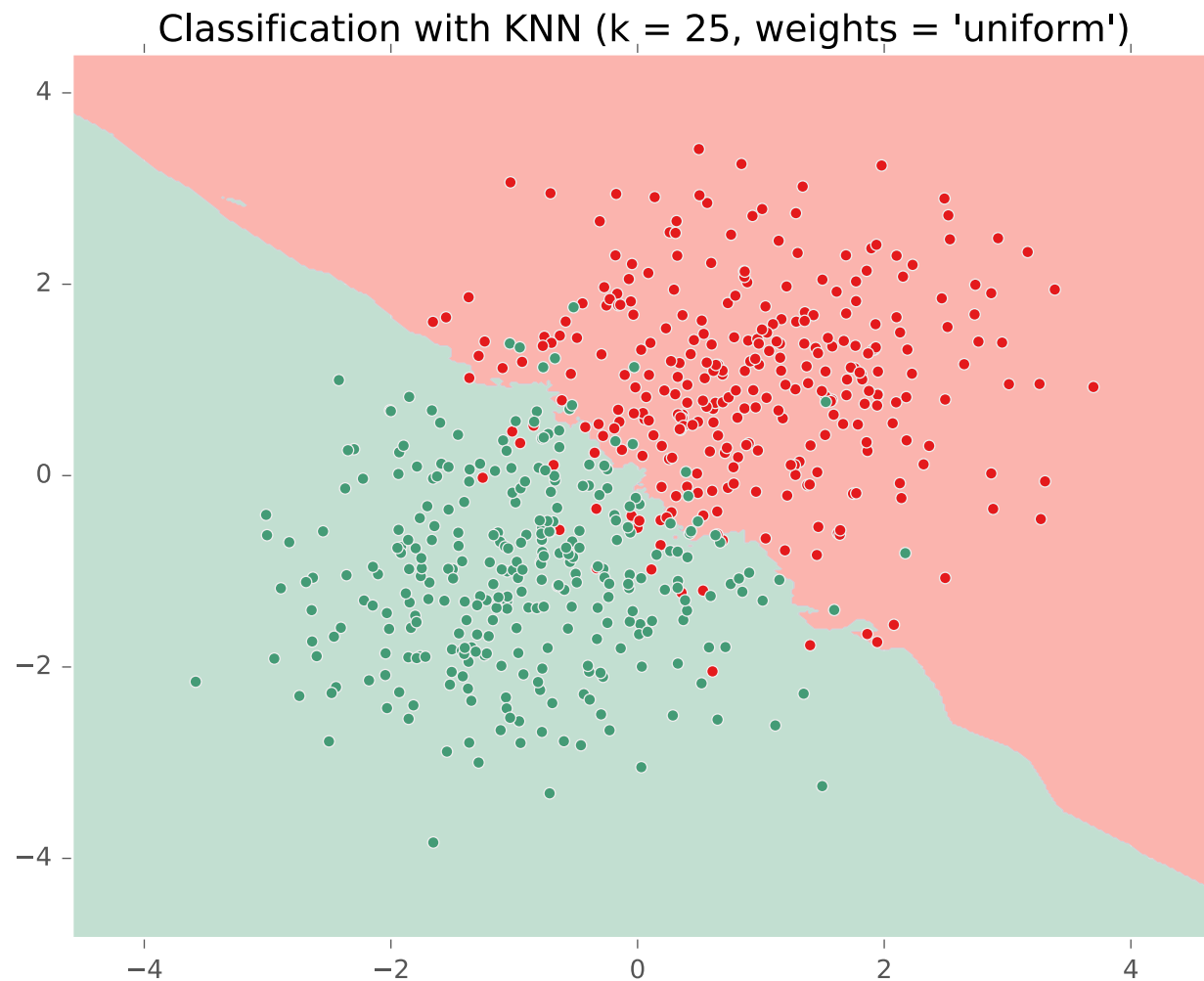
KNN on Gaussian Data



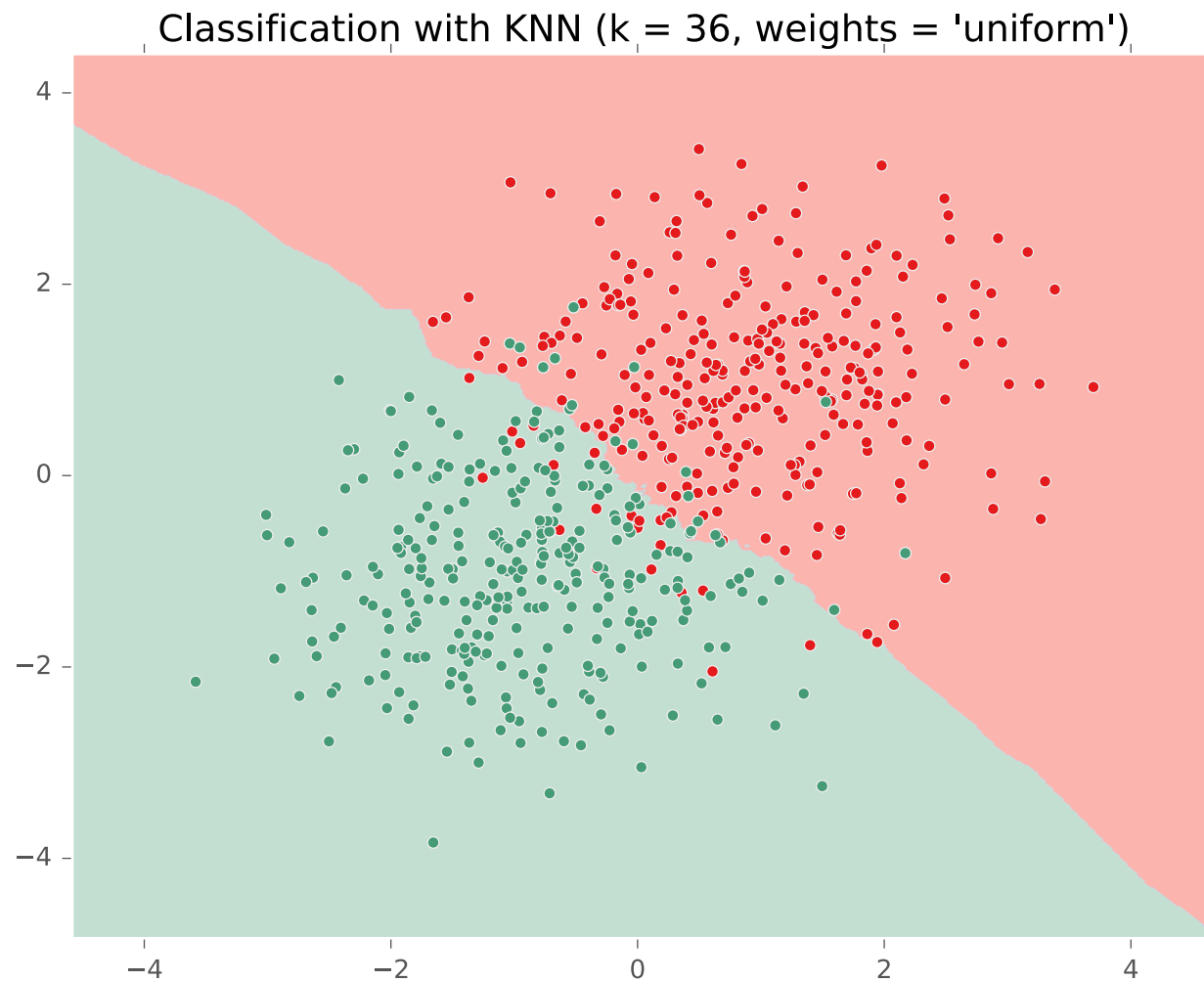
KNN on Gaussian Data



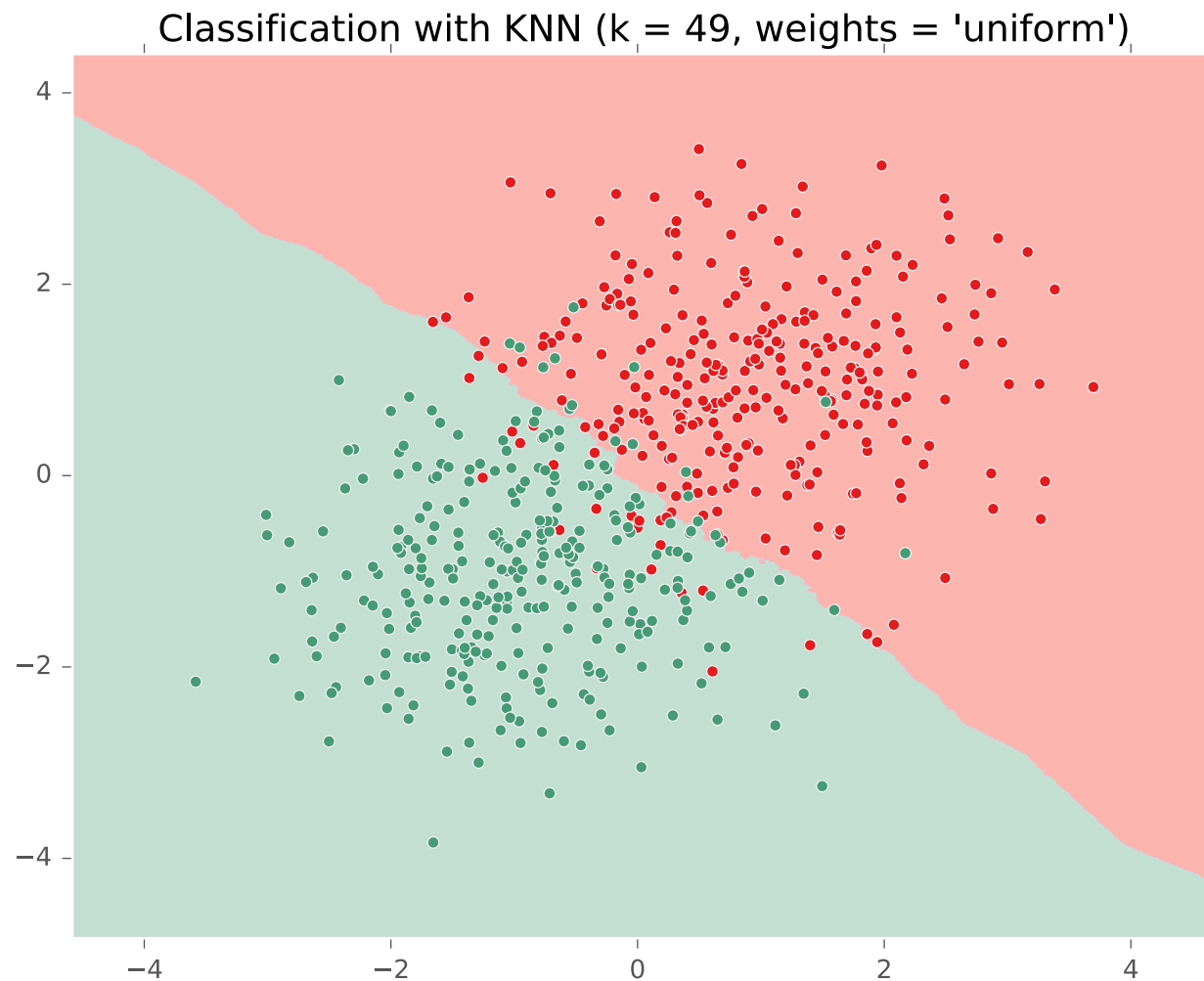
KNN on Gaussian Data



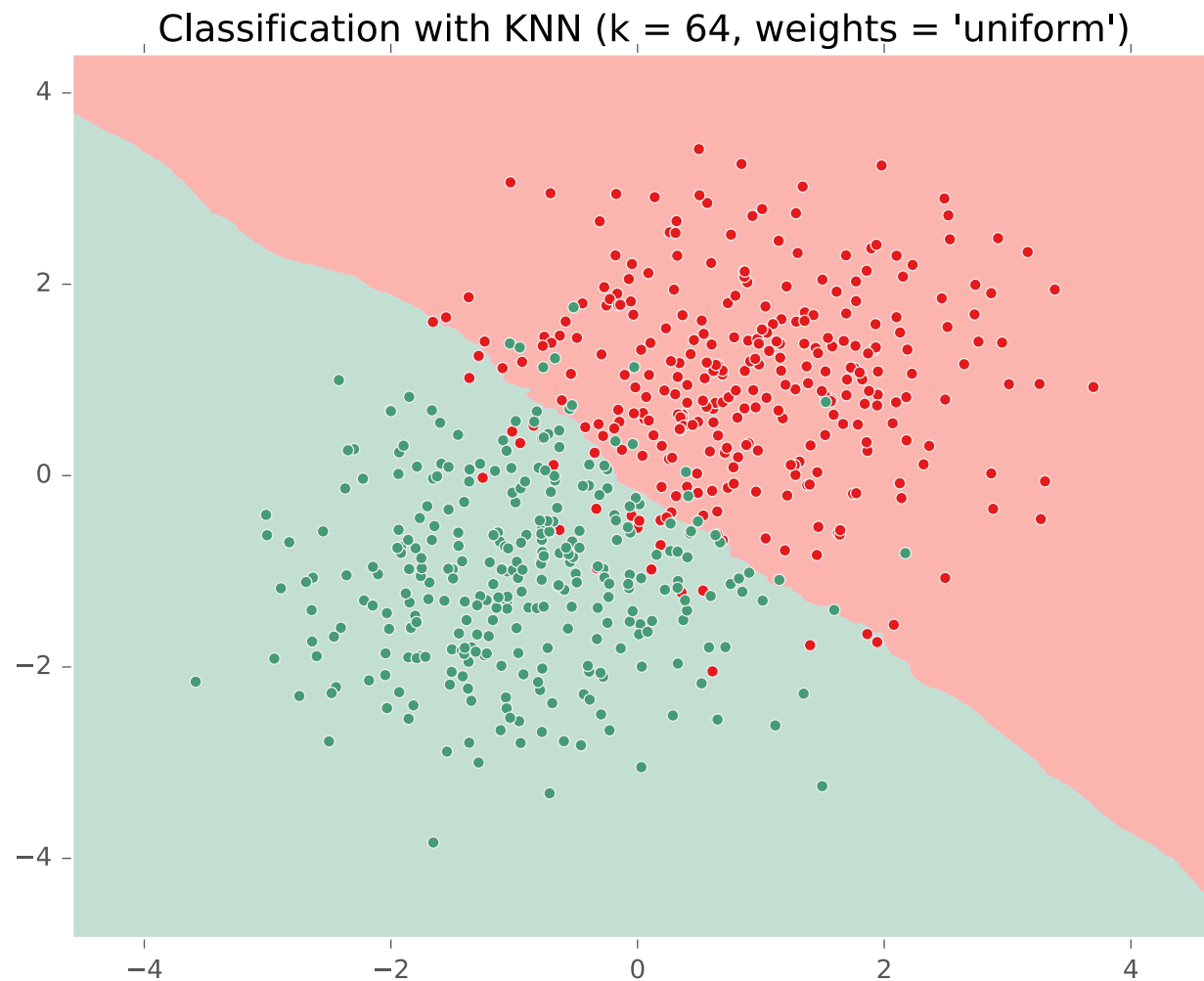
KNN on Gaussian Data



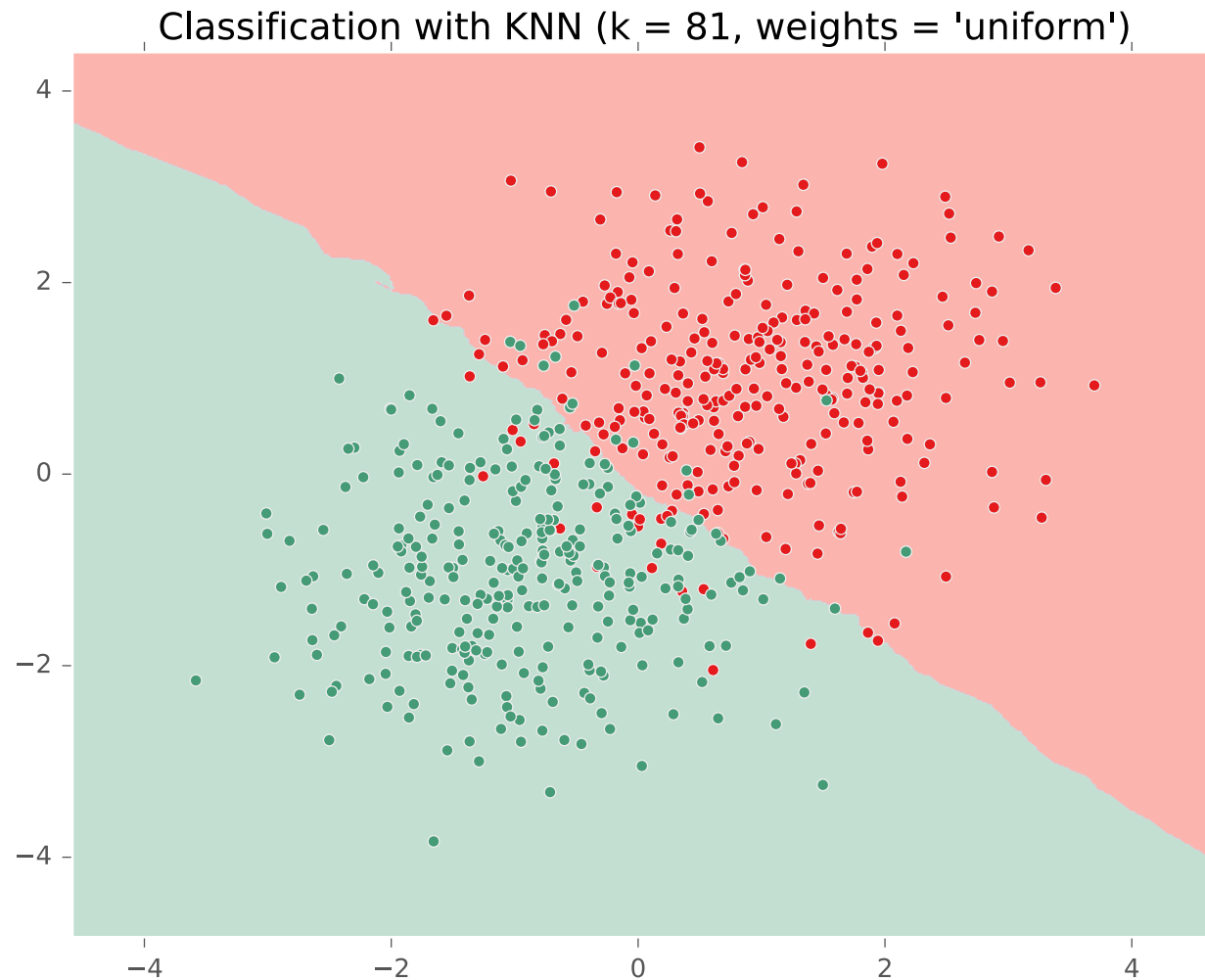
KNN on Gaussian Data



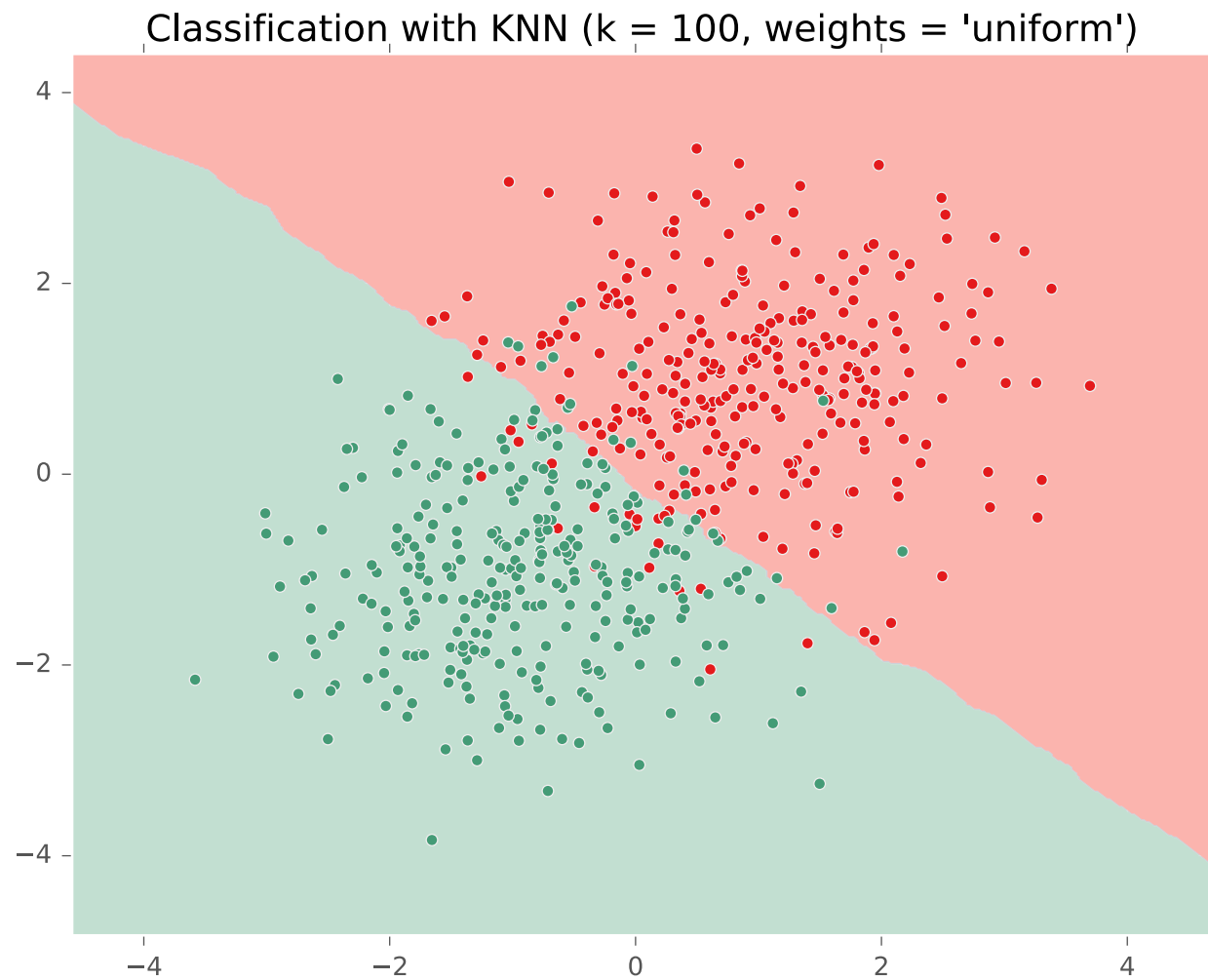
KNN on Gaussian Data



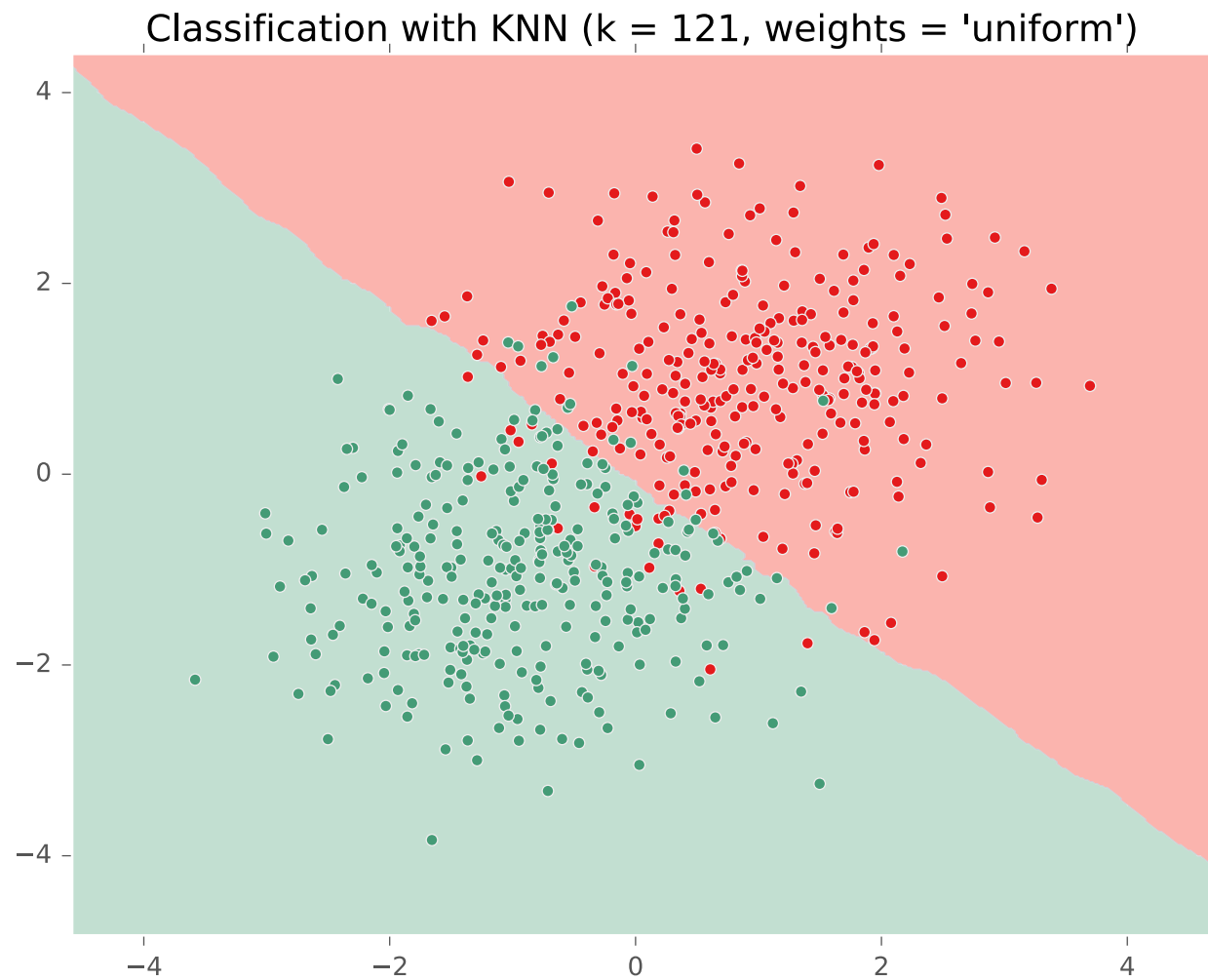
KNN on Gaussian Data



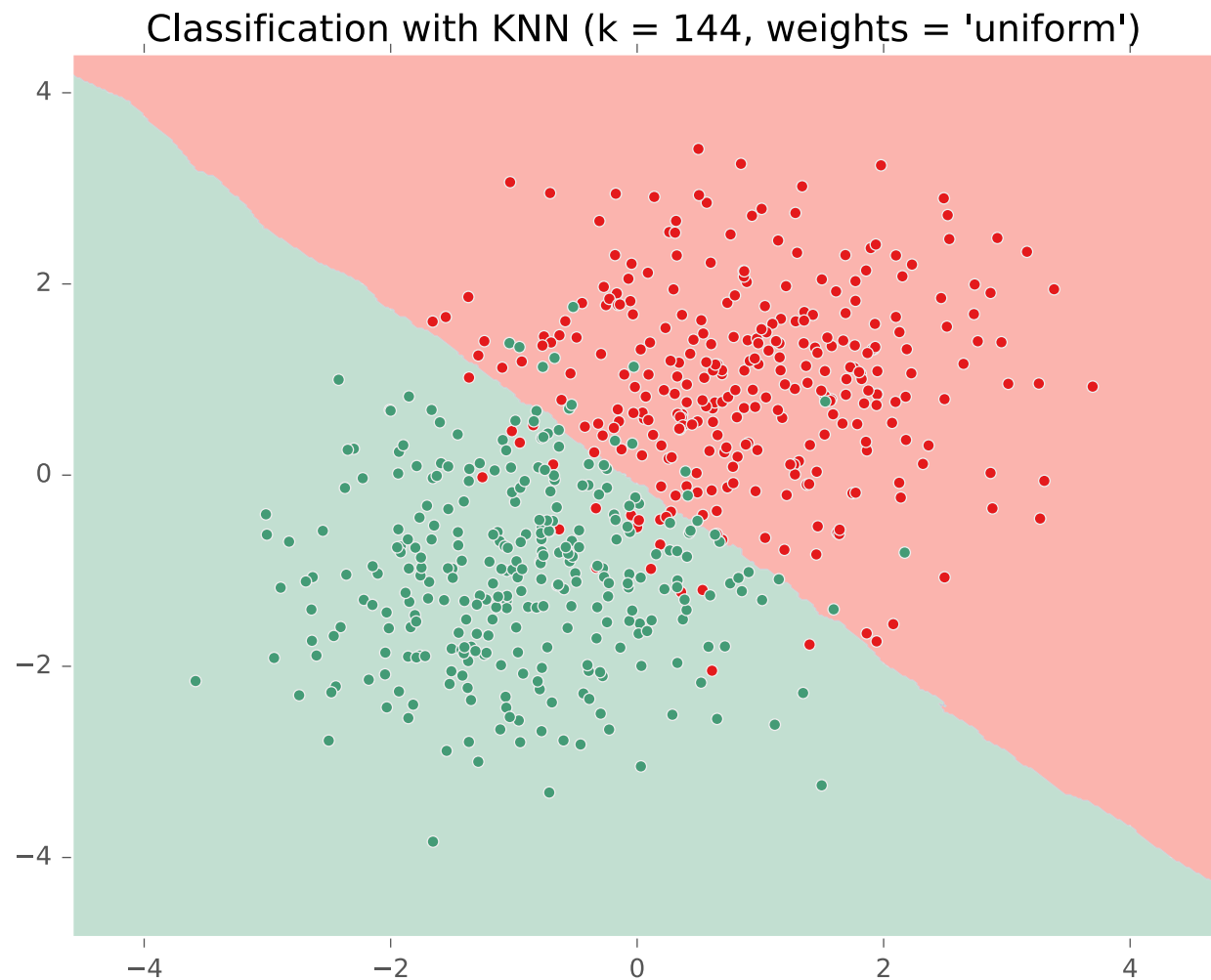
KNN on Gaussian Data



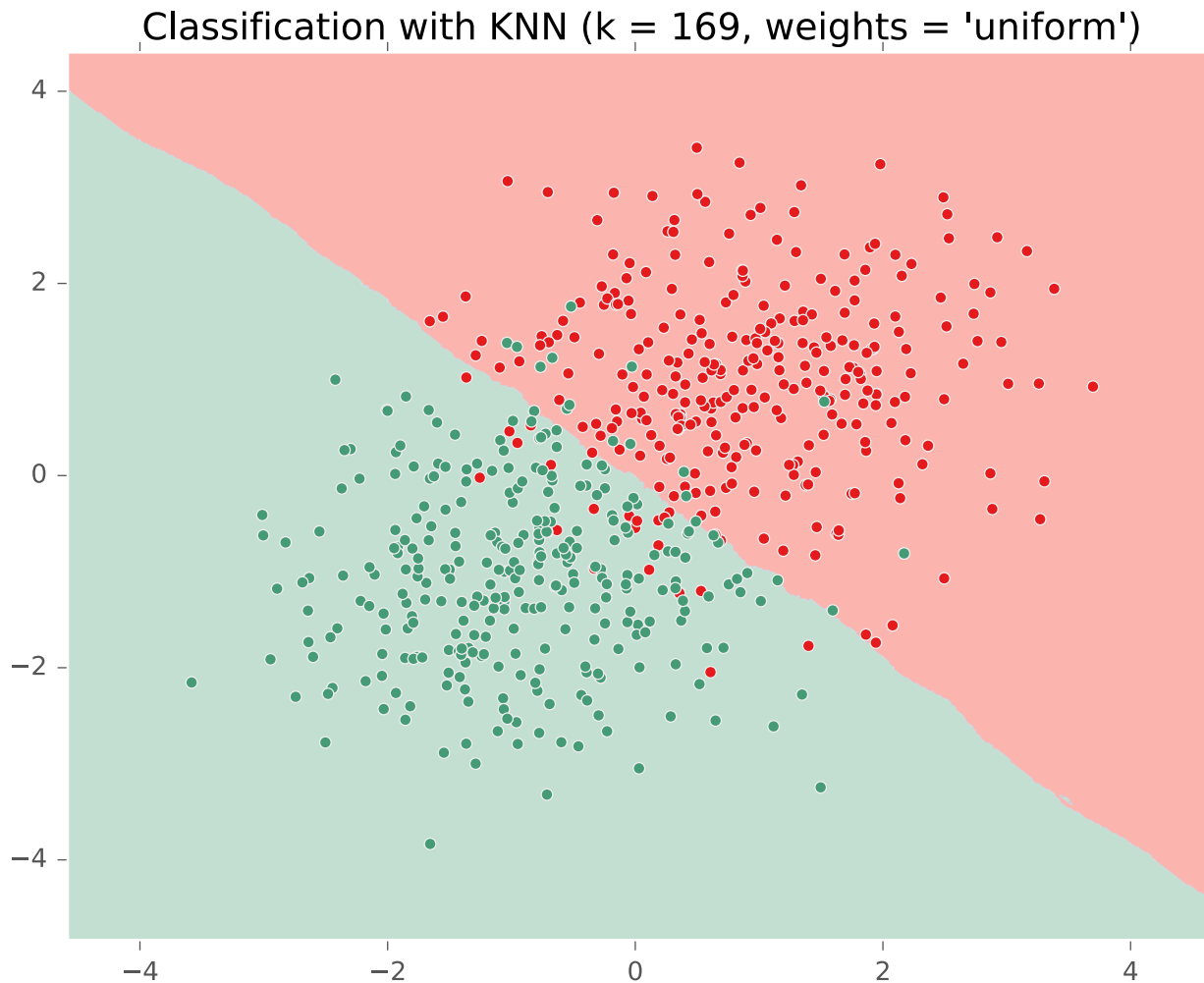
KNN on Gaussian Data



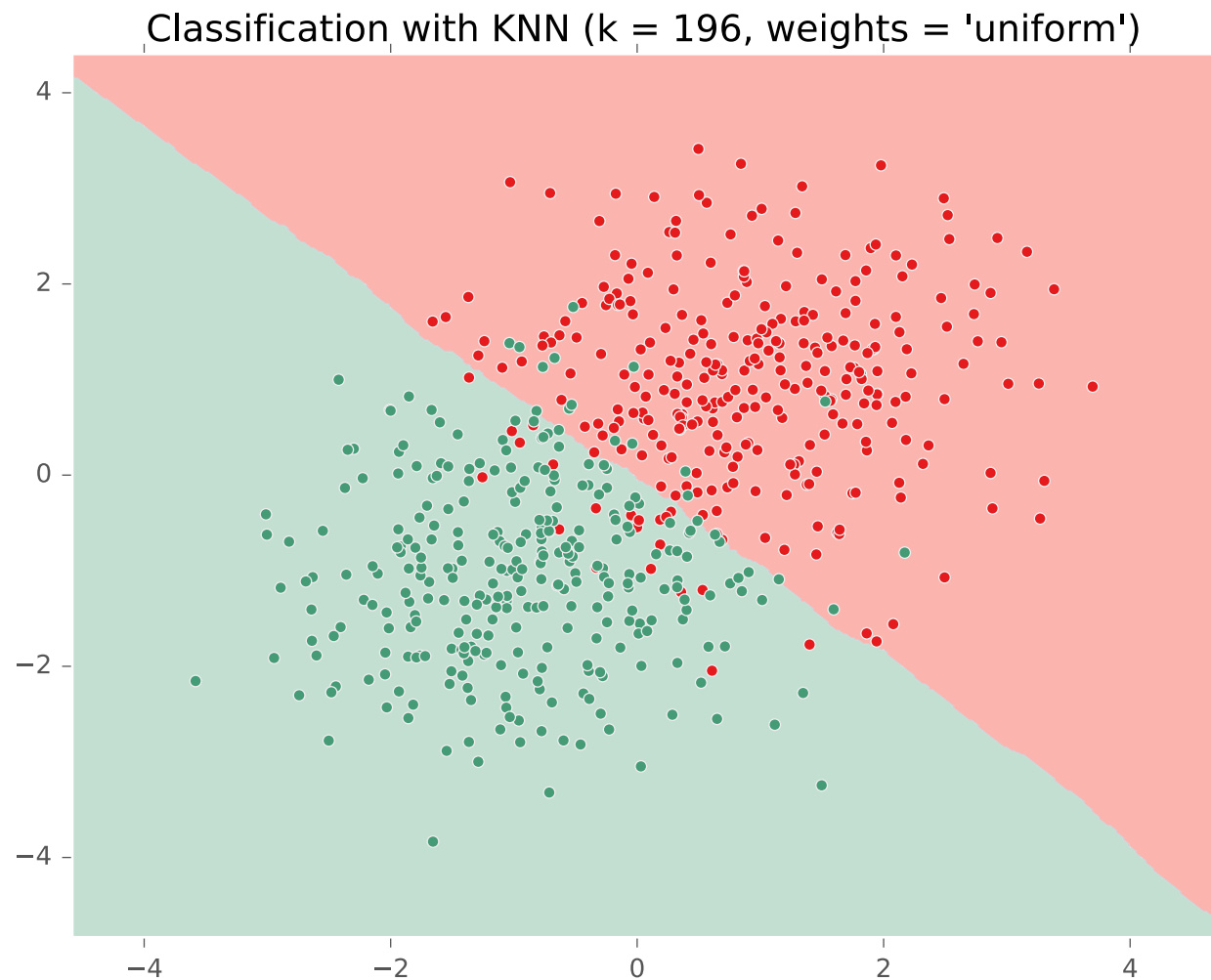
KNN on Gaussian Data



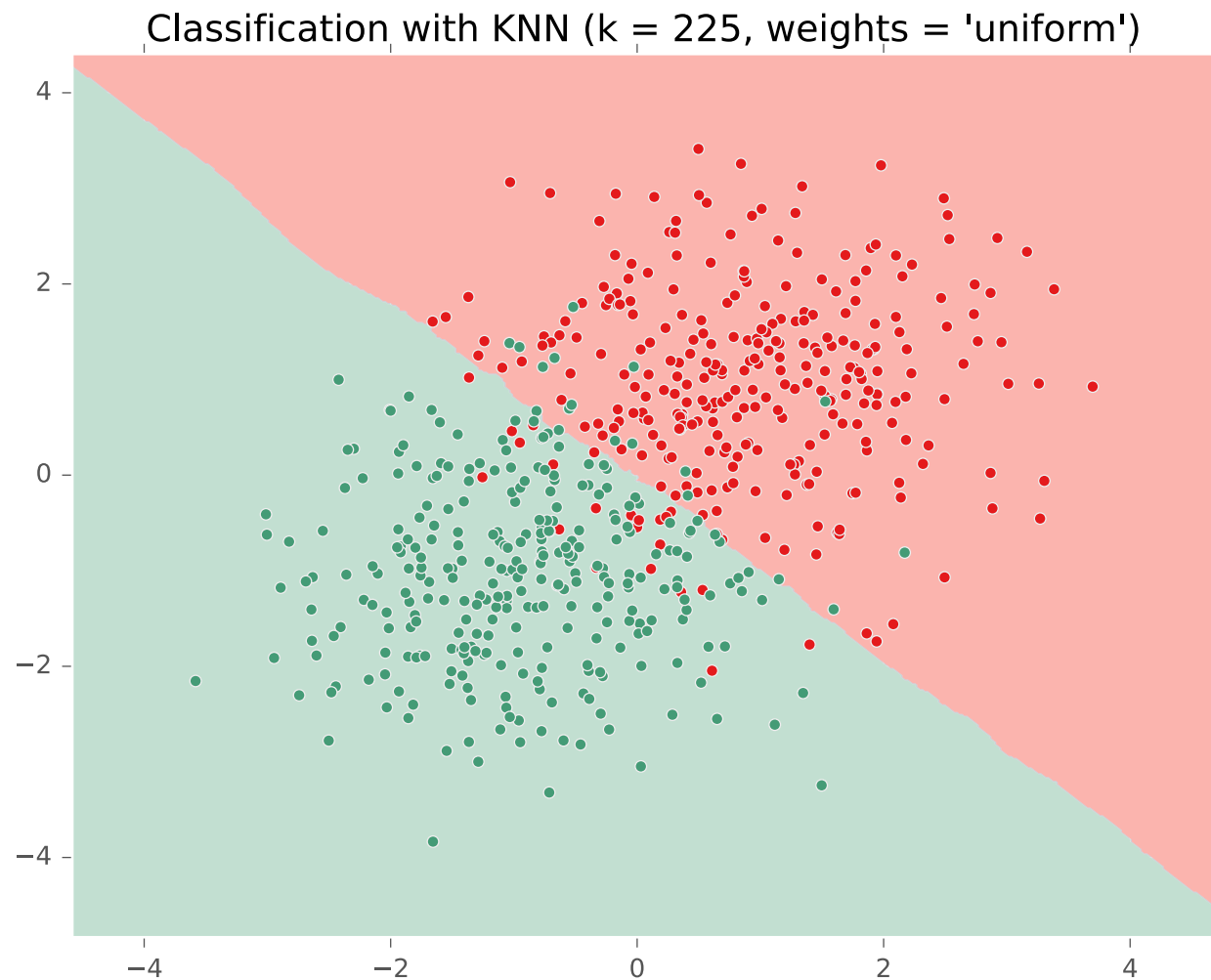
KNN on Gaussian Data



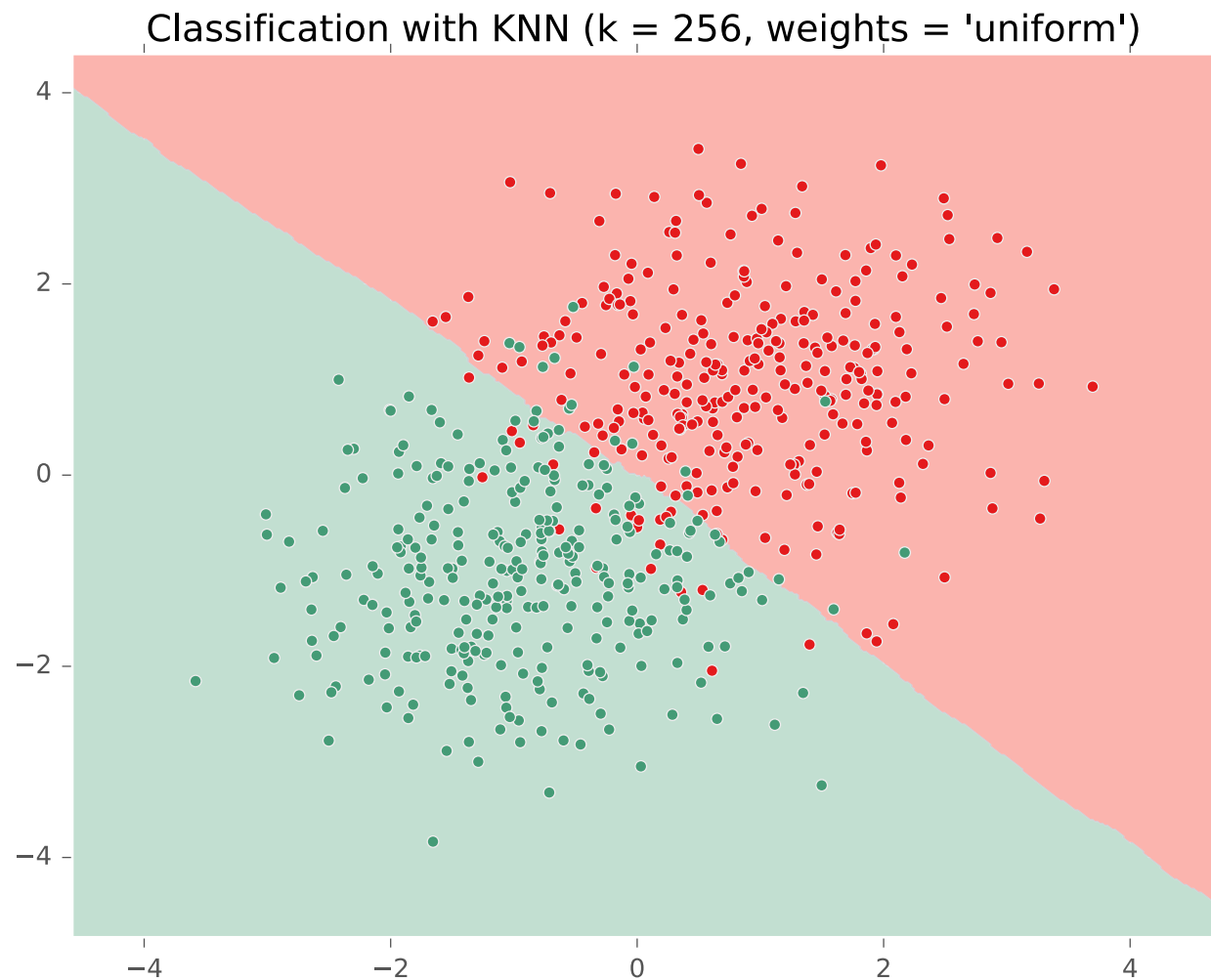
KNN on Gaussian Data



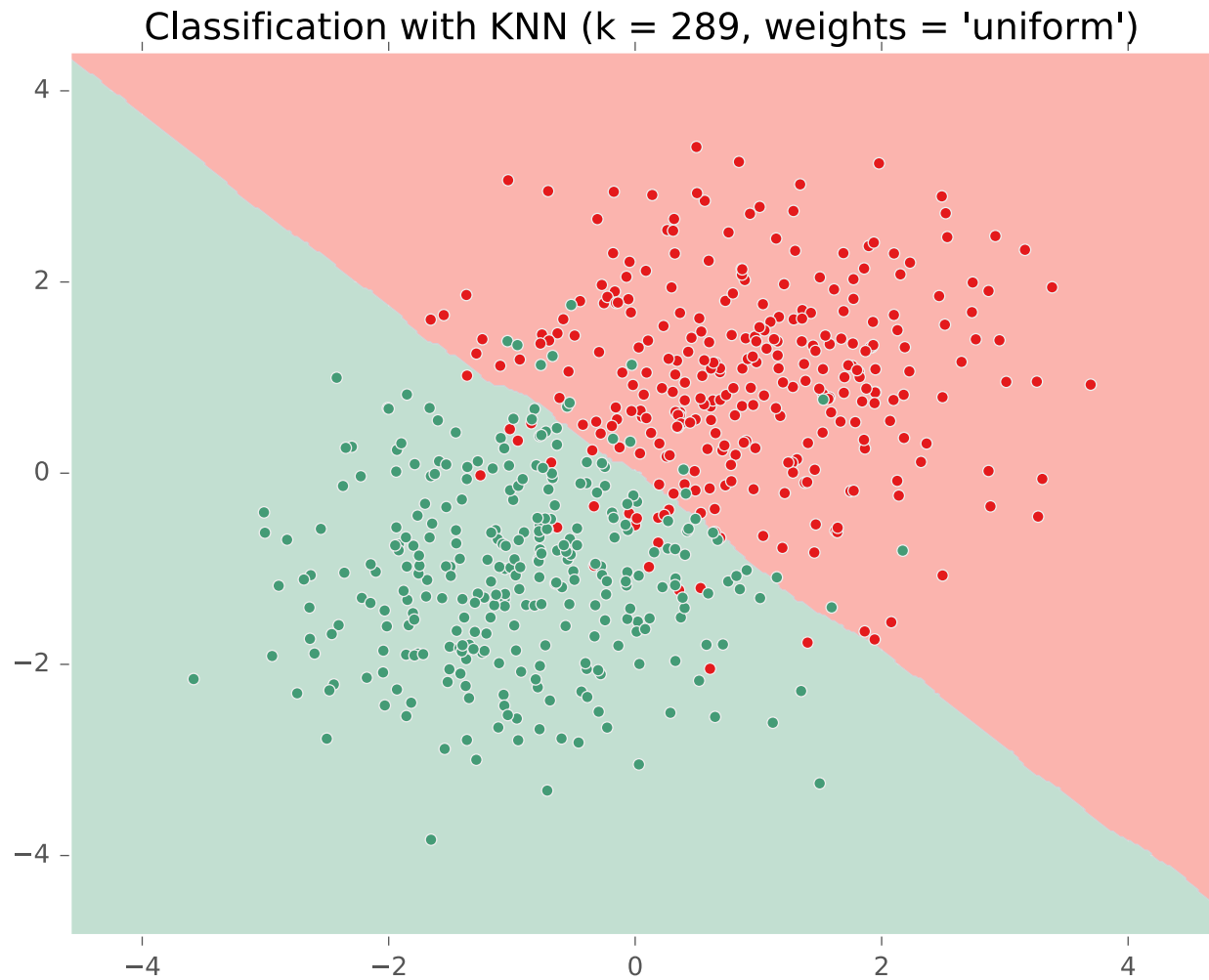
KNN on Gaussian Data



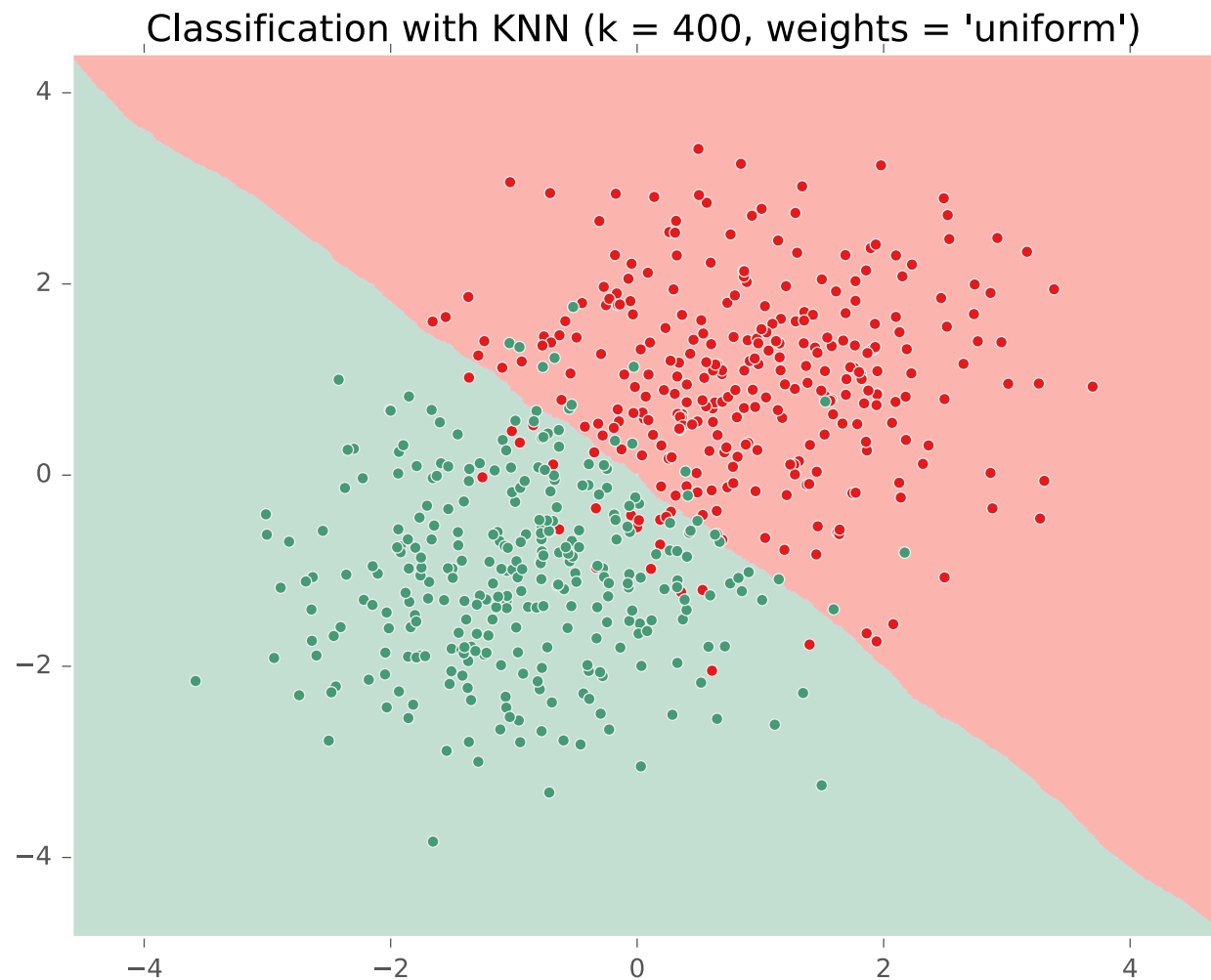
KNN on Gaussian Data



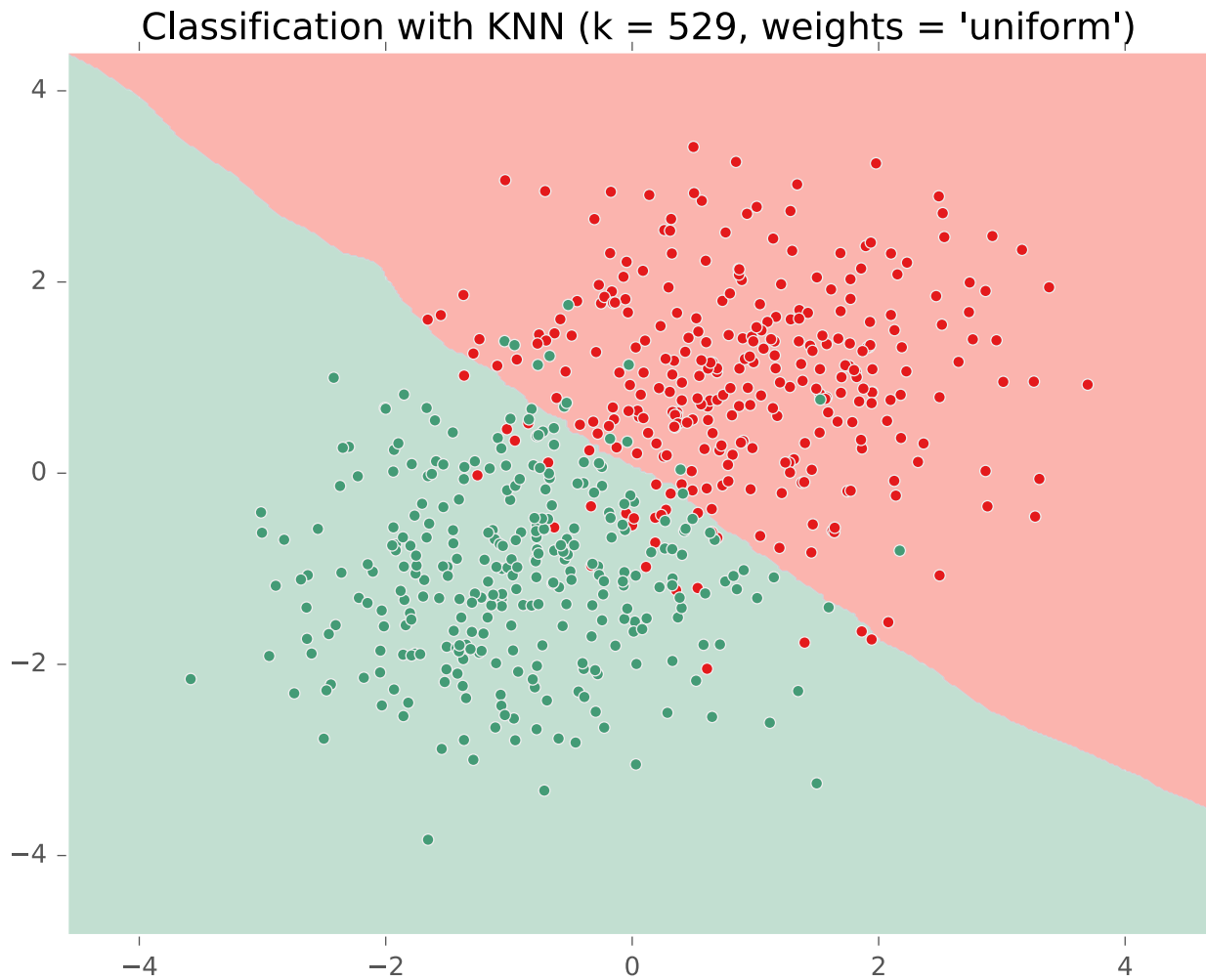
KNN on Gaussian Data



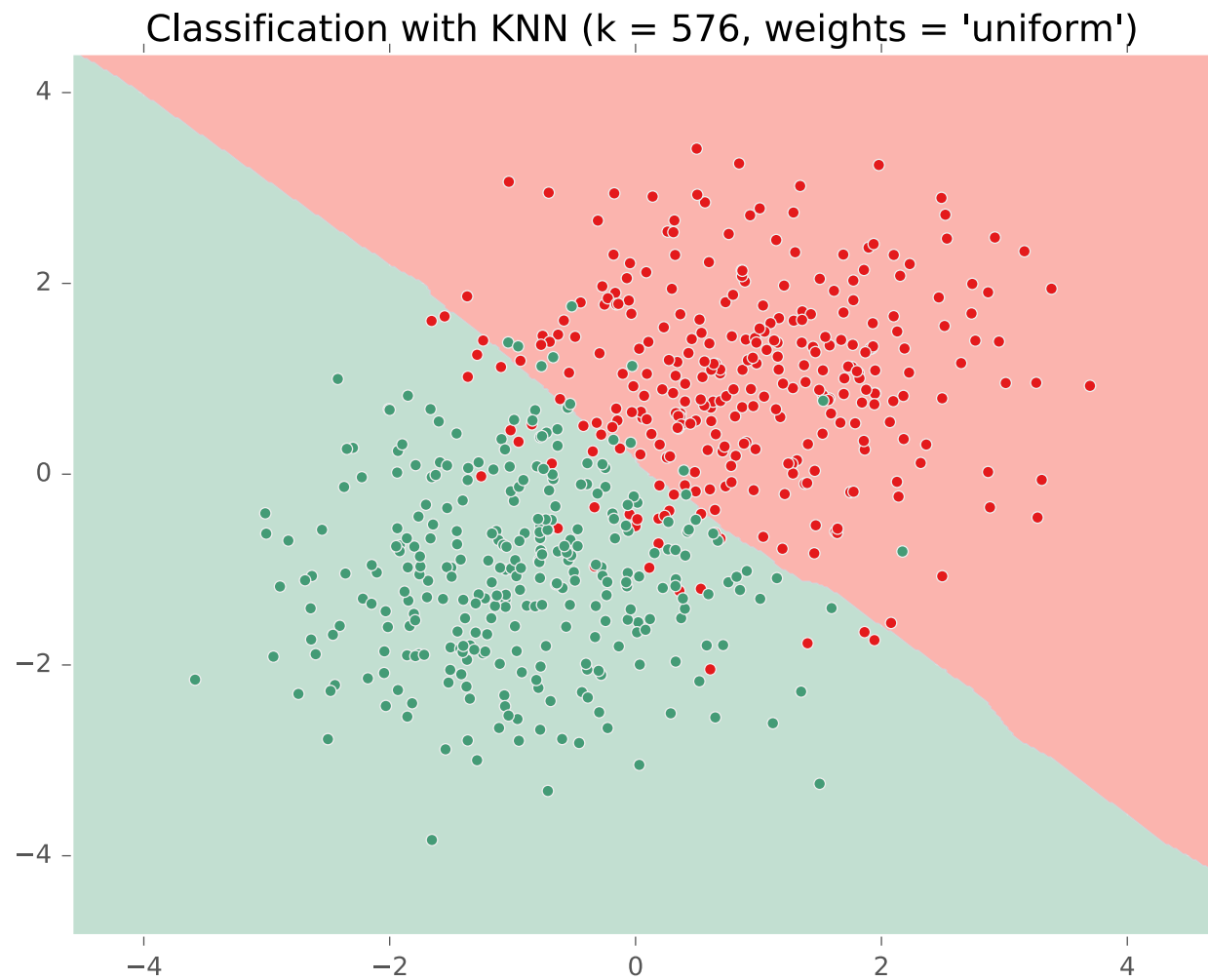
KNN on Gaussian Data



KNN on Gaussian Data



KNN on Gaussian Data



K-NEAREST NEIGHBORS

Questions

- How could k-Nearest Neighbors (KNN) be applied to **regression**?
- Can we do better than majority vote? (e.g. **distance-weighted KNN**)
- Where does the Cover & Hart (1967) **Bayes error rate bound** come from?

KNN Learning Objectives

You should be able to...

- Describe a dataset as points in a high dimensional space [CIML]
- Implement k-Nearest Neighbors with $O(N)$ prediction
- Describe the inductive bias of a k-NN classifier and relate it to feature scale [a la. CIML]
- Sketch the decision boundary for a learning algorithm (compare k-NN and DT)
- State Cover & Hart (1967)'s large sample analysis of a nearest neighbor classifier
- Invent "new" k-NN learning algorithms capable of dealing with even k
- Explain computational and geometric examples of the curse of dimensionality

MODEL SELECTION

Model Selection

WARNING:

- In some sense, our discussion of model selection is premature.
- The models we have considered thus far are fairly simple.
- The models and the many decisions available to the data scientist wielding them will grow to be much more complex than what we've seen so far.

Model Selection

Statistics

- *Def:* a **model** defines the data generation process (i.e. a set or family of parametric probability distributions)
- *Def:* **model parameters** are the values that give rise to a particular probability distribution in the model family
- *Def:* **learning** (aka. estimation) is the process of finding the parameters that best fit the data
- *Def:* **hyperparameters** are the parameters of a prior distribution over parameters

Machine Learning

- *Def:* (loosely) a **model** defines the hypothesis space over which learning performs its search
- *Def:* **model parameters** are the numeric values or structure selected by the learning algorithm that give rise to a hypothesis
- *Def:* the **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- *Def:* **hyperparameters** are the tunable aspects of the model, that the learning algorithm does *not* select

MODEL SELECTION

Model Selection

WARNING:

- In some sense, our discussion of model selection is premature.
- The models we have considered thus far are fairly simple.
- The models and the many decisions available to the data scientist wielding them will grow to be much more complex than what we've seen so far.

Model Selection

Statistics

- *Def:* a **model** defines the data generation process (i.e. a set or family of parametric probability distributions)
- *Def:* **model parameters** are the values that give rise to a particular probability distribution in the model family
- *Def:* **learning** (aka. estimation) is the process of finding the parameters that best fit the data
- *Def:* **hyperparameters** are the parameters of a prior distribution over parameters

Machine Learning

- *Def:* (loosely) a **model** defines the hypothesis space over which learning performs its search
- *Def:* **model parameters** are the numeric values or structure selected by the learning algorithm that give rise to a hypothesis
- *Def:* the **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- *Def:* **hyperparameters** are the tunable aspects of the model, that the learning algorithm does *not* select

Model Selection

Example: Decision Tree

- model = set of all possible trees, possibly restricted by some hyperparameters (e.g. max depth)
- parameters = structure of a specific decision tree
- learning algorithm = ID3, CART, etc.
- hyperparameters = max-depth, threshold for splitting criterion, etc.

Machine Learning

- *Def:* (loosely) a **model** defines the hypothesis space over which learning performs its search
- *Def:* **model parameters** are the numeric values or structure selected by the learning algorithm that give rise to a hypothesis
- *Def:* the **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- *Def:* **hyperparameters** are the tunable aspects of the model, that the learning algorithm does *not* select

Model Selection

Example: k-Nearest Neighbors

- model = set of all possible nearest neighbors classifiers
- parameters = none (KNN is an instance-based or non-parametric method)
- learning algorithm = for naïve setting, just storing the data
- hyperparameters = k , the number of neighbors to consider

Machine Learning

- *Def:* (loosely) a **model** defines the hypothesis space over which learning performs its search
- *Def:* **model parameters** are the numeric values or structure selected by the learning algorithm that give rise to a hypothesis
- *Def:* the **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- *Def:* **hyperparameters** are the tunable aspects of the model, that the learning algorithm does *not* select

Model Selection

Example: Perceptron

- model = set of all linear separators
- parameters = vector of weights (one for each feature)
- learning algorithm = mistake based updates to the parameters
- hyperparameters = none (unless using some variant such as averaged perceptron)

Machine Learning

- *Def:* (loosely) a **model** defines the hypothesis space over which learning performs its search
- *Def:* **model parameters** are the numeric values or structure selected by the learning algorithm that give rise to a hypothesis
- *Def:* the **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- *Def:* **hyperparameters** are the tunable aspects of the model, that the learning algorithm does *not* select

Model Selection

Statistics

- Def: a **model** defines the data generation process (i.e. a family of parameter distributions)
- Def: **model parameters** are the values that give rise to a particular probability distribution in the model
- Def: **learning** (aka. estimation) is the process of finding the parameters that best fit the data
- Def: **hyperparameters** are the parameters of a prior distribution over parameters

Machine Learning

- Def: (loosely) a **model** defines the hypothesis space over which the learning algorithm performs its search
- **hyperparameters** are the tunable aspects of the model, that the learning algorithm does not select
- **model parameters** are the values or structure that the learning algorithm selects to form a hypothesis
- **learning algorithm** defines the data-driven search over the hypothesis space (i.e. search for good parameters)
- Def: **hyperparameters** are the tunable aspects of the model, that the learning algorithm does not select

If “learning” is all about picking the best parameters how do we pick the best hyperparameters?



Model Selection

- Two *very* similar definitions:
 - *Def: model selection* is the process by which we choose the “best” model from among a set of candidates
 - *Def: hyperparameter optimization* is the process by which we choose the “best” hyperparameters from among a set of candidates (**could be called a special case of model selection**)
- **Both** assume access to a function capable of measuring the quality of a model
- **Both** are typically done “outside” the main training algorithm --- typically training is treated as a black box

Experimental Design

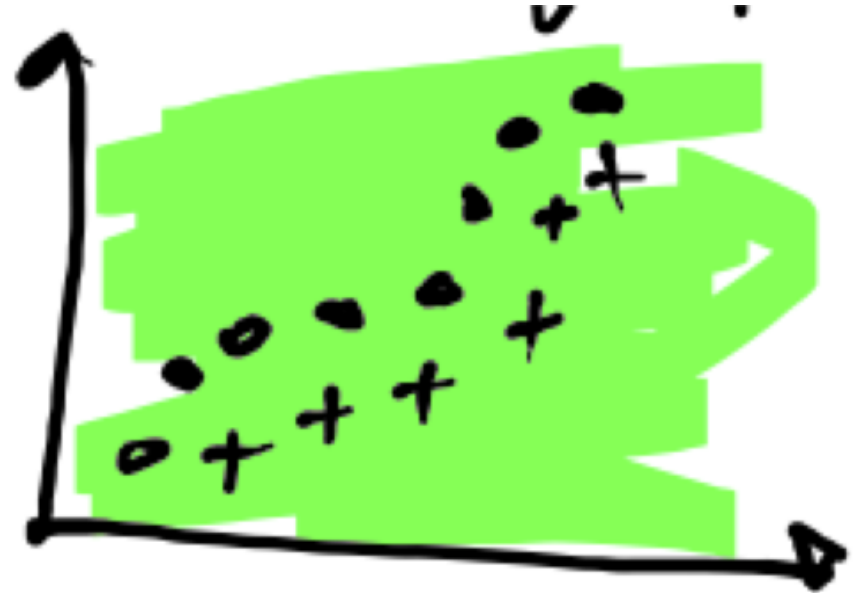
	Input	Output	Notes
Training	<ul style="list-style-type: none">• training dataset• hyperparameters	<ul style="list-style-type: none">• best model parameters	We pick the best model parameters by learning on the training dataset for a fixed set of hyperparameters
Hyperparameter Optimization	<ul style="list-style-type: none">• training dataset• validation dataset	<ul style="list-style-type: none">• best hyperparameters	We pick the best hyperparameters by learning on the training data and evaluating error on the validation error
Testing	<ul style="list-style-type: none">• test dataset• hypothesis (i.e. fixed model parameters)	<ul style="list-style-type: none">• test error	We evaluate a hypothesis corresponding to a decision rule with fixed model parameters on a test dataset to obtain test error

Special Cases of k-NN

k=1: Nearest Neighbor



k=N: Majority Vote



Example of Hyperparameter Opt.

Whiteboard:

- Special cases of k-Nearest Neighbors
- Choosing k with validation data
- Choosing k with cross-validation

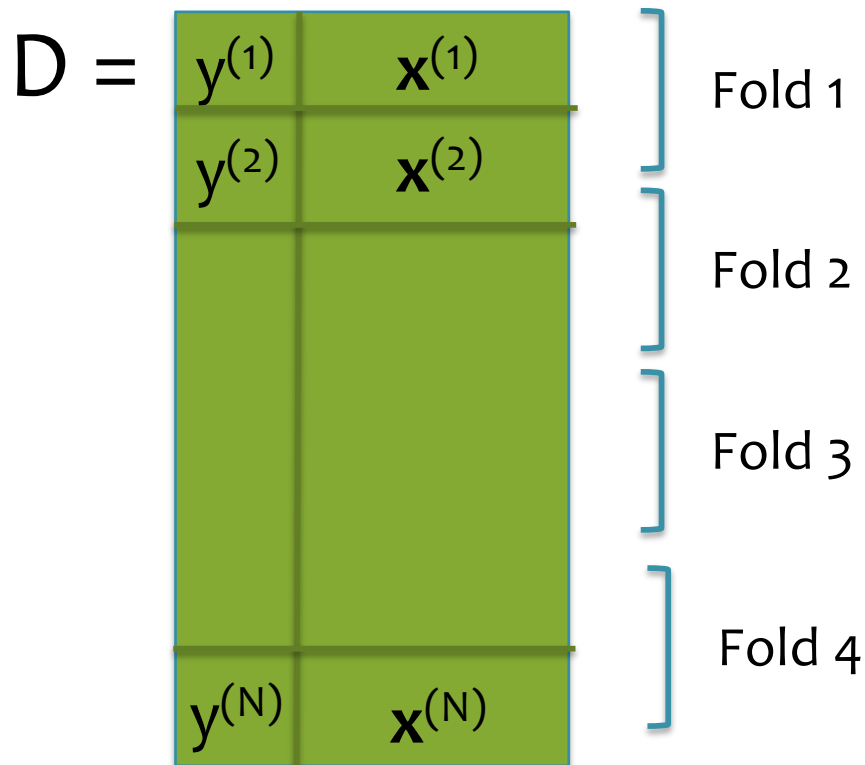
Cross-Validation

Cross validation is a method of estimating loss on held out data

Input: training data, learning algorithm, loss function (e.g. 0/1 error)

Output: an estimate of loss function on held-out data

Key idea: rather than just a single “validation” set, use many!
(Error is more stable. Slower computation.)



Algorithm:

Divide data into folds (e.g. 4)

1. Train on folds $\{1,2,3\}$ and predict on $\{4\}$
2. Train on folds $\{1,2,4\}$ and predict on $\{3\}$
3. Train on folds $\{1,3,4\}$ and predict on $\{2\}$
4. Train on folds $\{2,3,4\}$ and predict on $\{1\}$

Concatenate all the predictions and evaluate loss (*almost* equivalent to averaging loss over the folds)

Definition:
N-fold cross validation = cross validation with N folds

Experimental Design

	Input	Output	Notes
Training	<ul style="list-style-type: none">• training dataset• hyperparameters	<ul style="list-style-type: none">• best model parameters	We pick the best model parameters by learning on the training dataset for a fixed set of hyperparameters
Hyperparameter Optimization	<ul style="list-style-type: none">• training dataset• validation dataset	<ul style="list-style-type: none">• best hyperparameters	We pick the best hyperparameters by learning on the training data and evaluating error on the validation error
Cross-Validation	<ul style="list-style-type: none">• training dataset• validation dataset	<ul style="list-style-type: none">• cross-validation error	We estimate the error on held out data by repeatedly training on N-1 folds and predicting on the held-out fold
Testing	<ul style="list-style-type: none">• test dataset• hypothesis (i.e. fixed model parameters)	<ul style="list-style-type: none">• test error	We evaluate a hypothesis corresponding to a decision rule with fixed model parameters on a test dataset to obtain test error

Experimental Design

Q: We pick the best hyperparameters by learning on the training data and evaluating error on the validation error. For our final model, should we then learn from training + validation?

A: Yes.

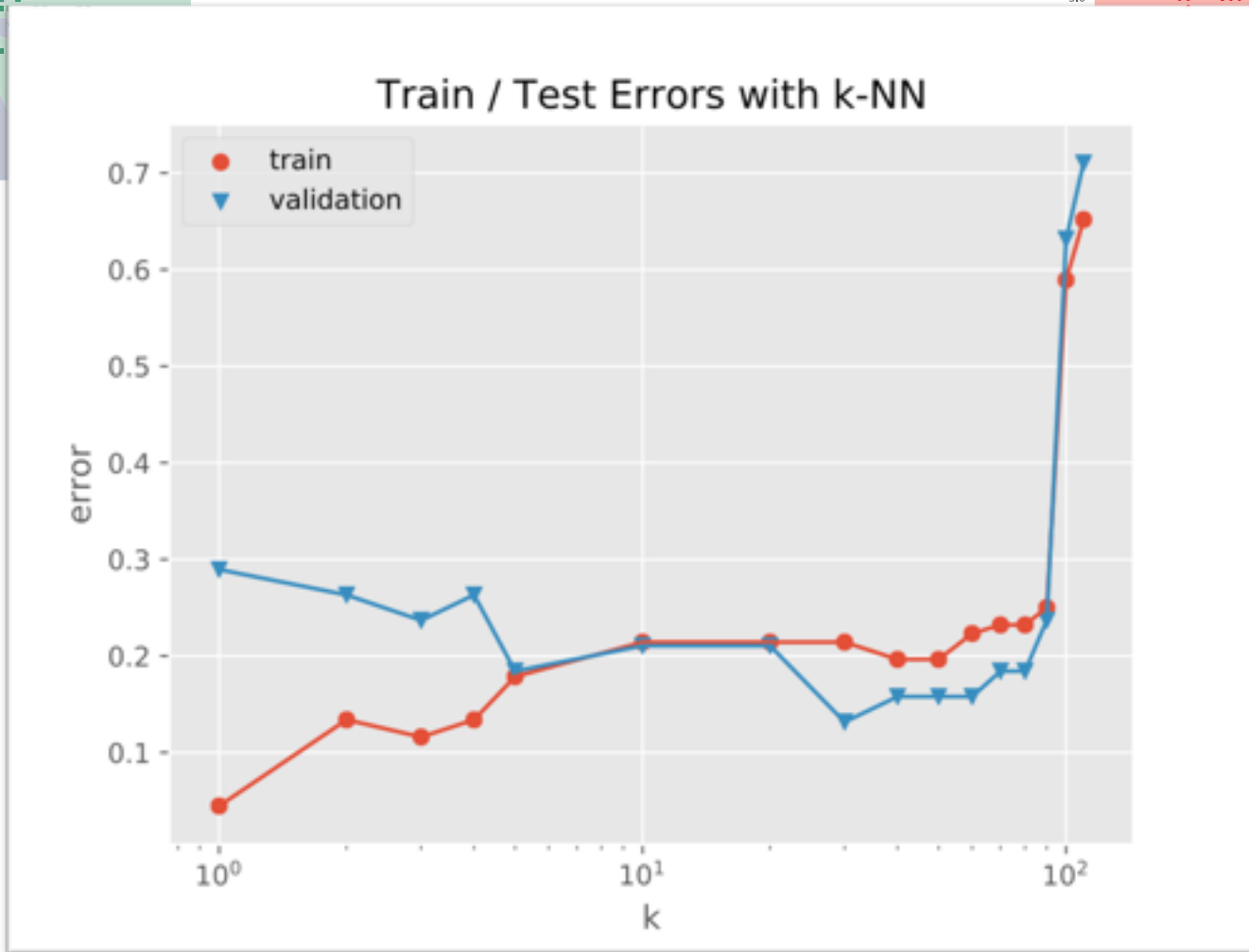
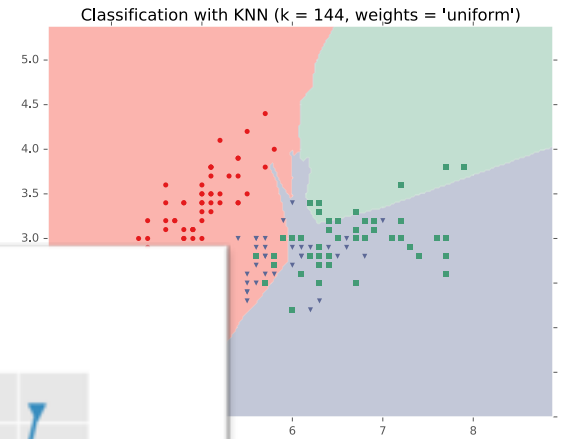
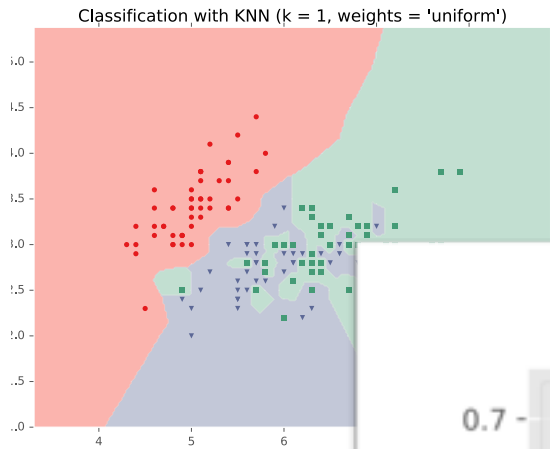
Let's assume that {train-original} is the original training data, and {test} is the provided test dataset.

1. Split {train-original} into {train-subset} and {validation}.
2. Pick the hyperparameters that when training on {train-subset} give the lowest error on {validation}. Call these hyperparameters {best-hyper}.
3. Retrain a new model using {best-hyper} on {train-original} = {train-subset} \cup {validation}.
4. Report test error by evaluating on {test}.

Alternatively, you could replace Steps 1-2 with the following:

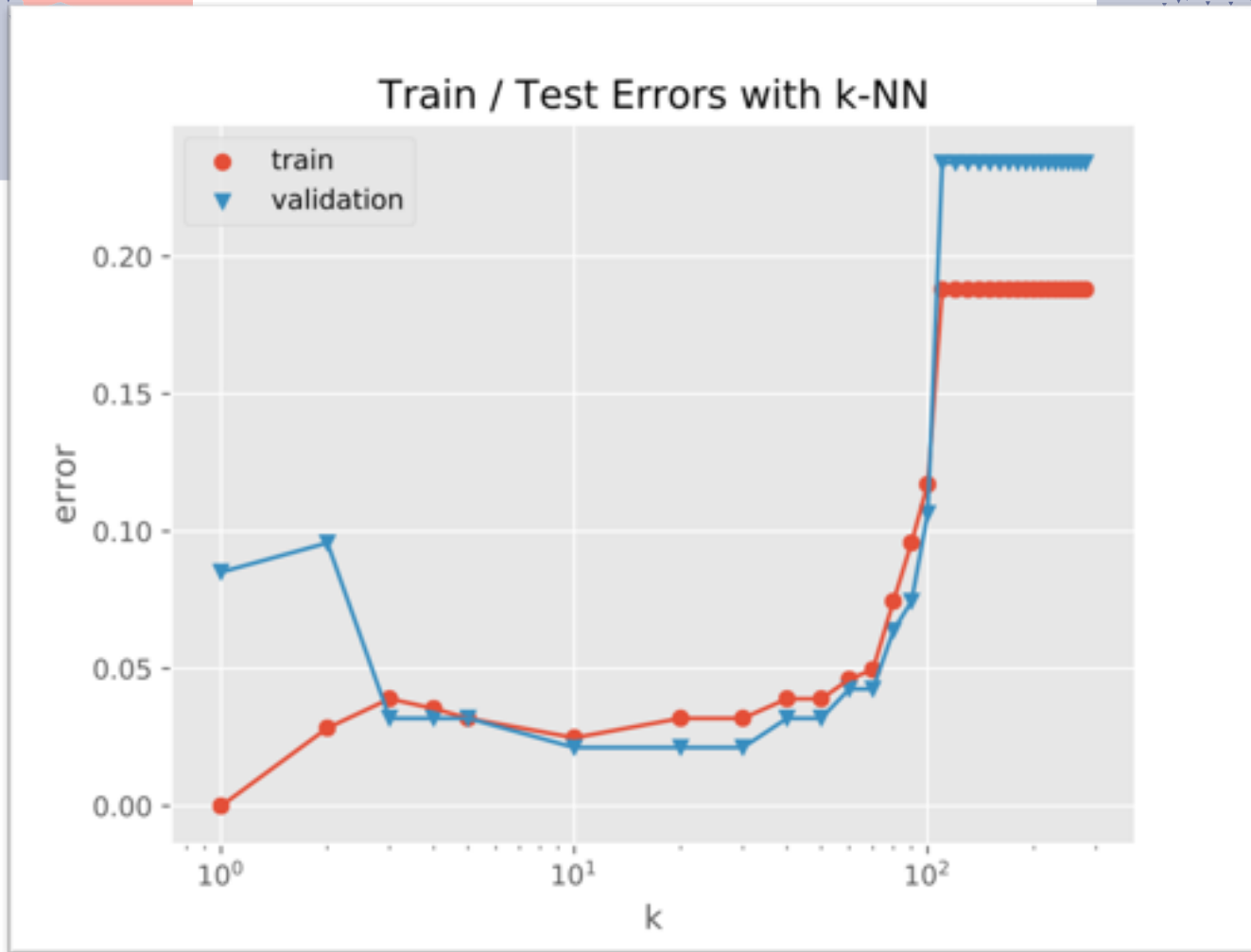
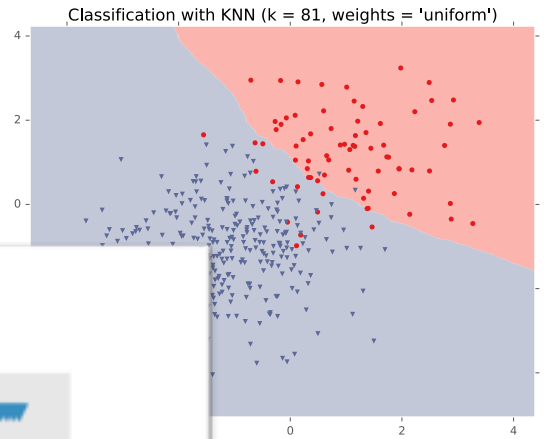
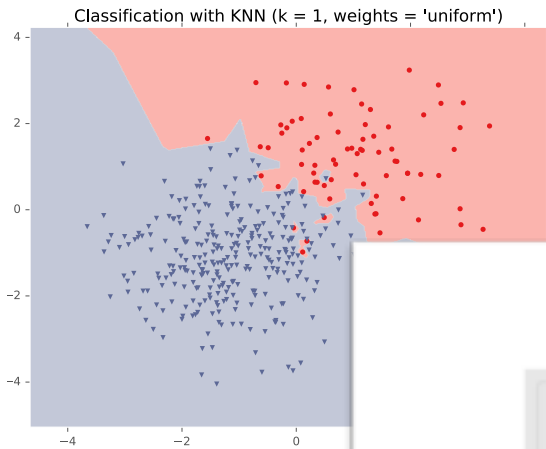
1. Pick the hyperparameters that give the lowest cross-validation error on {train-original}. Call these hyperparameters {best-hyper}.

k-NN: Choosing k



Fisher Iris Data: varying the value of k

k-NN: Choosing k



Gaussian Data: varying the value of k

Model Selection

WARNING (again):

- This section is only scratching the surface!
- Lots of methods for hyperparameter optimization: (to talk about later)
 - Grid search
 - Random search
 - Bayesian optimization
 - Graduate-student descent
 - ...

Main Takeaway:

- Model selection / hyperparameter optimization is just another form of learning

Hyperparameter Optimization

Setting: suppose we have hyperparameters α , β , and χ and we wish to pick the “best” values for each one

Algorithm 1: Grid Search

- Pick a set of values for each hyperparameter
 $\alpha \in \{a_1, a_2, \dots, a_n\}$, $\beta \in \{b_1, b_2, \dots, b_n\}$, and $\chi \in \{c_1, c_2, \dots, c_n\}$
- Run a grid search

```
for  $\alpha \in \{a_1, a_2, \dots, a_n\}$ :  
  for  $\beta \in \{b_1, b_2, \dots, b_n\}$ :  
    for  $\chi \in \{c_1, c_2, \dots, c_n\}$ :  
       $\theta = \text{train}(D_{\text{train}}; \alpha, \beta, \chi)$   
       $\text{error} = \text{predict}(D_{\text{validation}}; \theta)$ 
```

- return α , β , and χ with lowest validation error

Hyperparameter Optimization

Setting: suppose we have hyperparameters α , β , and χ and we wish to pick the “best” values for each one

Algorithm 2: Random Search

- Pick a range of values for each parameter
 $\alpha \in \{a_1, a_2, \dots, a_n\}$, $\beta \in \{b_1, b_2, \dots, b_n\}$, and $\chi \in \{c_1, c_2, \dots, c_n\}$
- Run a random search

for $t = 1, 2, \dots, T$:

sample α uniformly from $\{a_1, a_2, \dots, a_n\}$

sample β uniformly from $\{b_1, b_2, \dots, b_n\}$

sample χ uniformly from $\{c_1, c_2, \dots, c_n\}$

$\theta = \text{train}(D_{\text{train}}; \alpha, \beta, \chi)$

error = $\text{predict}(D_{\text{validation}}; \theta)$

- return α , β , and χ with lowest validation error

Hyperparameter Optimization

Question:

True or False: given a finite amount of computation time, **grid search** is more likely to find good values for hyperparameters than **random search**.

Answer:

Model Selection Learning Objectives

You should be able to...

- Plan an experiment that uses training, validation, and test datasets to predict the performance of a classifier on unseen data (without cheating)
- Explain the difference between (1) training error, (2) validation error, (3) cross-validation error, (4) test error, and (5) true error
- For a given learning technique, identify the model, learning algorithm, parameters, and hyperparameters
- Define "instance-based learning" or "nonparametric methods"
- Select an appropriate algorithm for optimizing (aka. learning) hyperparameters

THE PERCEPTRON ALGORITHM

Perceptron: History

Imagine you are trying to build a new machine learning technique... your name is Frank Rosenblatt... and the year is 1957

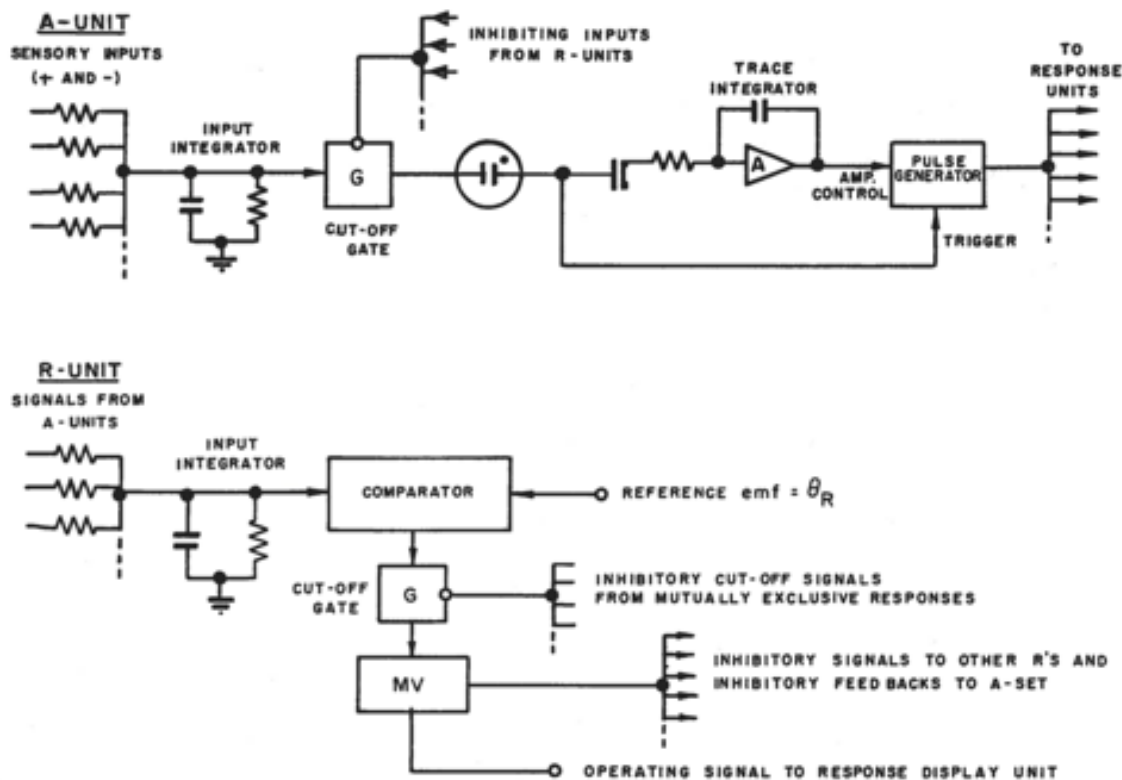
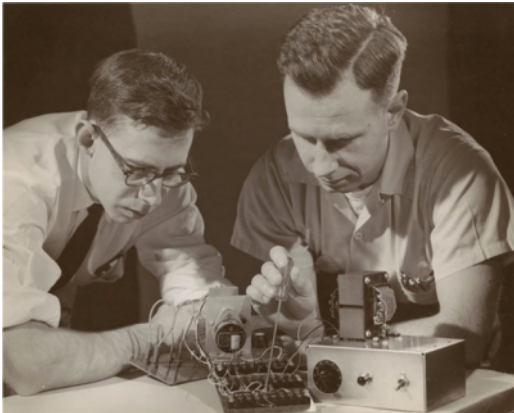


FIGURE 5
DESIGN OF TYPICAL UNITS

Perceptron: History

Imagine you are trying to build a new machine learning technique... your name is Frank Rosenblatt... and the year is 1957



The New Yorker, December 6, 1958 P. 44

Talk story about the perceptron, a new electronic brain which hasn't been built, but which has been successfully simulated on the I.B.M. 704. Talk with Dr. Frank Rosenblatt, of the Cornell Aeronautical Laboratory, who is one of the two men who developed the prodigy; the other man is Dr. Marshall C. Yovits, of the Office of Naval Research, in Washington. Dr. Rosenblatt defined the perceptron as the first non-biological object which will achieve an organization o its external environment in a meaningful way. It interacts with its environment, forming concepts that have not been made ready for it by a human agent. If a triangle is held up, the perceptron's eye picks up the image & conveys it along a random succession of lines to the response units, where the image is registered. It can tell the difference betw. a cat and a dog, although it wouldn't be able to tell whether the dog was to theleft or right of the cat. Right now it is of no practical use, Dr. Rosenblatt conceded, but he said that one day it might be useful to send one into outer space to take in impressions for us.



Linear Models for Classification

Key idea: Try to learn this hyperplane directly

Looking ahead:

- We'll see a number of commonly used Linear Classifiers
- These include:
 - Perceptron
 - Logistic Regression
 - Naïve Bayes (under certain conditions)
 - Support Vector Machines

Directly modeling the hyperplane would use a decision function:

$$h(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

for:

$$y \in \{-1, +1\}$$

Geometry

In-Class Exercise

Draw a picture of the region corresponding to:

$$w_1x_1 + w_2x_2 + b > 0$$

$$\text{where } w_1 = 2, w_2 = 3, b = 6$$

Draw the vector

$$\mathbf{w} = [w_1, w_2]$$

Answer Here:

