



NASSLLI 2018

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Machine Learning

Matt Gormley
June 23 - 24, 2018

Prerequisites

1. Be able to find Porter Hall 100
2. Later, be able to find GHC 6115 -- using a specific entrance door (near Cyert Hall) to enter between the hours of 12pm and 7pm.
3. Maybe speak a little bit of Python



MEAN HALL 5TH FLOOR
LOBBY CLOSED
USE
← DOWNTOWN HALL →
OR
← SCOTT HALL →
TO ACCESS MALL



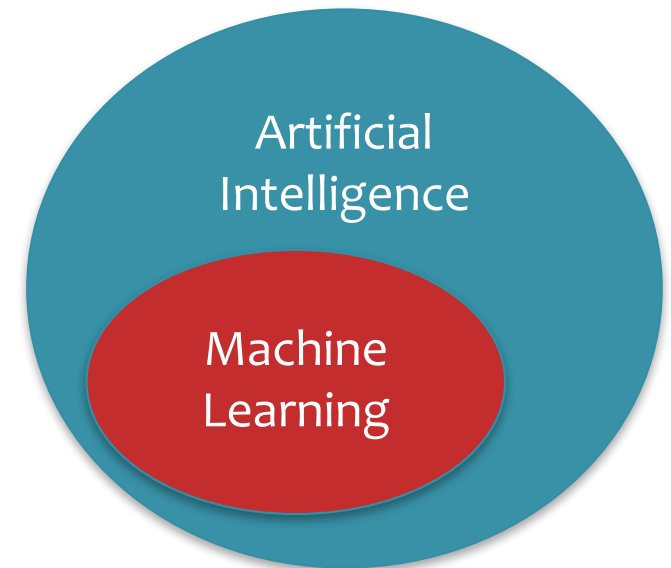
WHAT IS MACHINE LEARNING?

Artificial Intelligence

The basic goal of AI is to develop intelligent machines.

This consists of many sub-goals:

- Perception
- Reasoning
- Control / Motion / Manipulation
- Planning
- Communication
- Creativity
- Learning



What is Machine Learning?

The goal of this course is to provide you with a toolbox:

Machine Learning

Statistics

Probability

Computer Science

Optimization



Computer
Science

What is ML?

Domain of
Interest

Machine Learning

Optimization

Statistics

Probability

Calculus

Measure
Theory

Linear Algebra

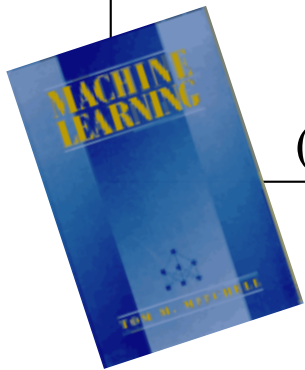
Speech Recognition

1. Learning to recognize spoken words

THEN

“...the SPHINX system (e.g. Lee 1989) learns speaker-specific strategies for recognizing the primitive sounds (phonemes) and words from the observed speech signal...neural network methods...hidden Markov models...”

(Mitchell, 1997)



NOW



Source: <https://www.stonetemple.com/great-knowledge-box-showdown/#VoiceStudyResults>

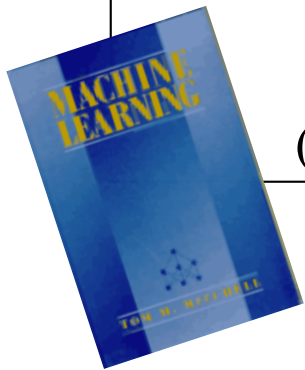
Robotics

2. Learning to drive an autonomous vehicle

THEN

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



NOW



waymo.com

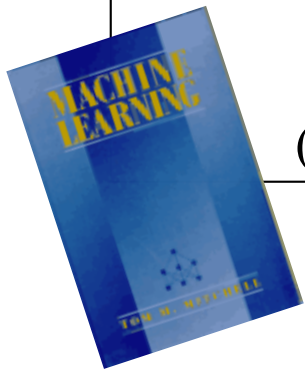
Robotics

2. Learning to drive an autonomous vehicle

THEN

“...the ALVINN system (Pomerleau 1989) has used its learned strategies to drive unassisted at 70 miles per hour for 90 miles on public highways among other cars...”

(Mitchell, 1997)



NOW



<https://www.geek.com/wp-content/uploads/2016/03/uber.jpg>

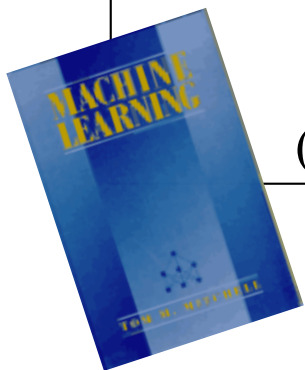
Games / Reasoning

3. Learning to beat the masters at board games

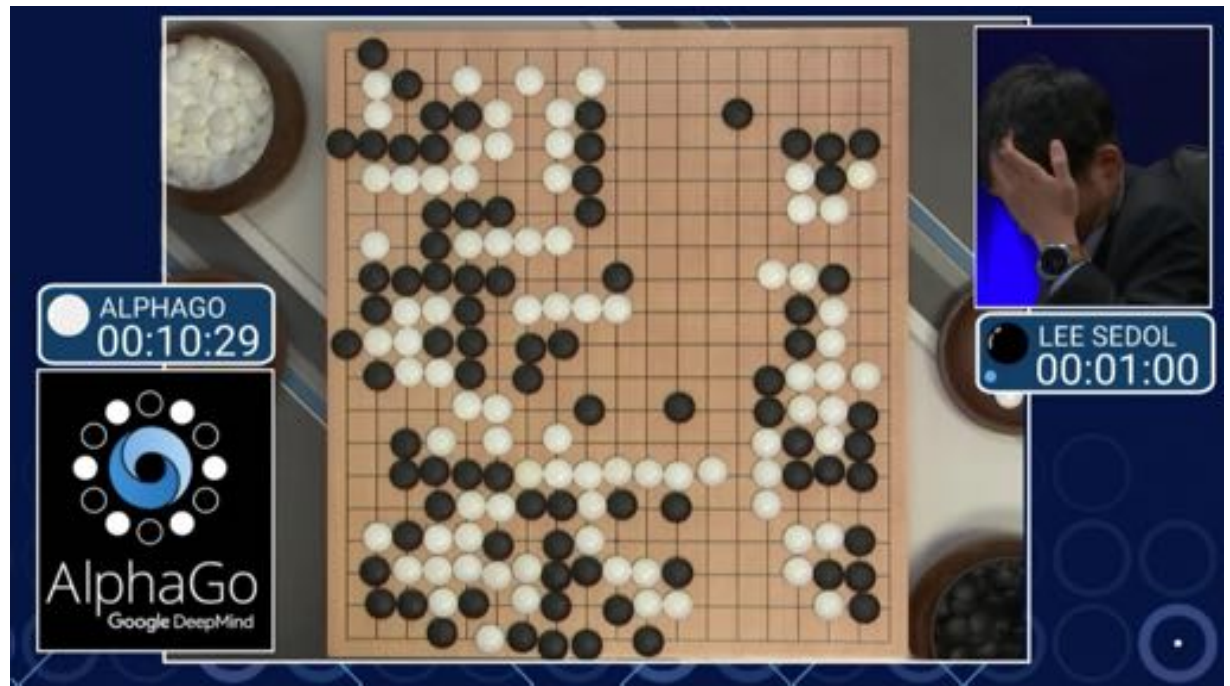
THEN

“...the world’s top computer program for backgammon, TD-GAMMON (Tesauro, 1992, 1995), learned its strategy by playing over one million practice games against itself...”

(Mitchell, 1997)



NOW

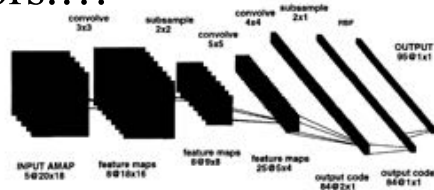


Computer Vision

4. Learning to recognize images

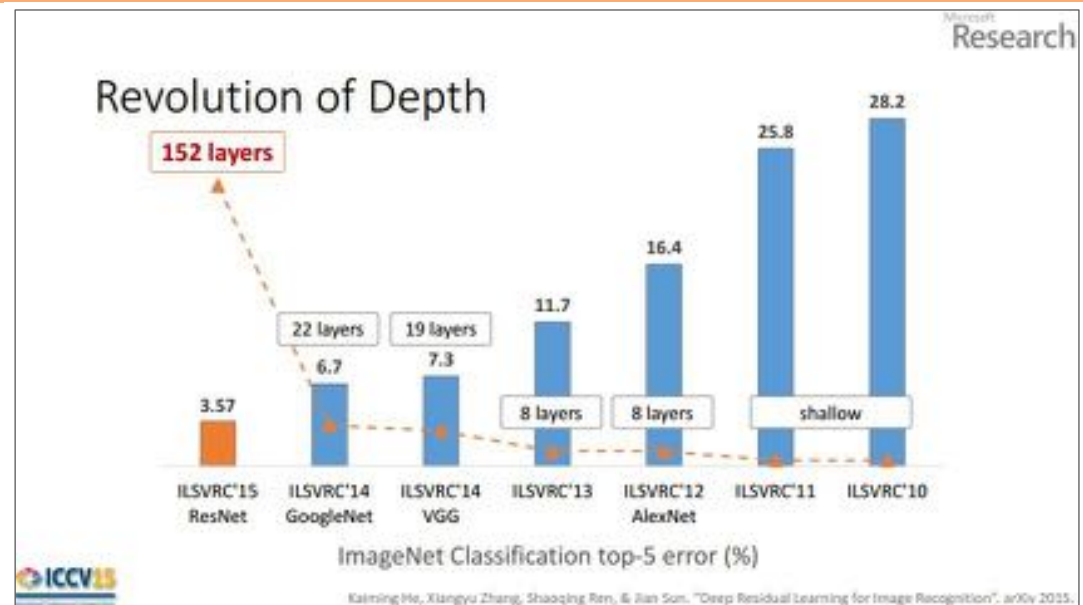
THEN

“...The recognizer is a convolution network that can be spatially replicated. From the network output, a hidden Markov model produces word scores. The entire system is globally trained to minimize word-level errors....”



(LeCun et al., 1995)

NOW



Learning Theory

• 5. In what cases and how well can we learn?

Sample Complexity Results

Definition 0.1. The **sample complexity** of a learning algorithm is the number of examples required to achieve arbitrarily small error (with respect to the optimal hypothesis) with high probability (i.e. close to 1).

Four Cases we care about...

	Realizable	Agnostic
Finite $ \mathcal{H} $	$N \geq \frac{1}{\epsilon} [\log(\mathcal{H}) + \log(\frac{1}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$.	$N \geq \frac{1}{2\epsilon^2} [\log(\mathcal{H}) + \log(\frac{2}{\delta})]$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) < \epsilon$.
Infinite $ \mathcal{H} $	$N = O(\frac{1}{\epsilon} [\text{VC}(\mathcal{H}) \log(\frac{1}{\epsilon}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ all $h \in \mathcal{H}$ with $R(h) \geq \epsilon$ have $\hat{R}(h) > 0$.	$N = O(\frac{1}{\epsilon^2} [\text{VC}(\mathcal{H}) + \log(\frac{1}{\delta})])$ labeled examples are sufficient so that with probability $(1 - \delta)$ for all $h \in \mathcal{H}$ we have that $ R(h) - \hat{R}(h) \leq \epsilon$.

Two Types of Error

① True Error (aka. expected risk) (aka. Generalization Error)

$$R(h) = \mathbb{P}_{x \sim p^*(x)} (c^*(x) \neq h(x)) \quad \leftarrow \text{always unknown.}$$

② Train Error (aka. empirical risk)

$$\begin{aligned} \hat{R}(h) &= \mathbb{P}_{x \sim S} (c^*(x) \neq h(x)) \quad \leftarrow S = \{x^{(1)}, \dots, x^{(N)}\} \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(c^*(x^{(i)}) \neq h(x^{(i)})) \quad \leftarrow \text{known, computable} \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y^{(i)} \neq h(x^{(i)})) \end{aligned}$$

PAC Learning

Q: Can we bound $R(h)$ in terms of $\hat{R}(h)$?
A: Yes!

PAC stands for Probably Approximately Correct

PAC learner yields hypothesis h , which is approximately correct $R(h) \approx 0$ with high probability $\Pr(R(h) \approx 0) \approx 1$

Def: PAC Criterion

$$\Pr(\forall h, |R(h) - \hat{R}(h)| \leq \epsilon) \geq 1 - \delta$$

1. How many examples do we need to learn?
2. How do we quantify our ability to generalize to unseen data?
3. Which algorithms are better suited to specific learning settings?

What is Machine Learning?

The goal of this course is to provide you with a toolbox:

Machine Learning

Statistics

Probability

Computer Science

Optimization

To solve all the problems above and more



Topics

- Foundations
 - Probability
 - MLE, MAP
 - Optimization
- Classifiers
 - KNN
 - Naïve Bayes
 - Logistic Regression
 - Perceptron
 - SVM
- Regression
 - Linear Regression
- Important Concepts
 - Kernels
 - Regularization and Overfitting
 - Experimental Design
- Unsupervised Learning
 - K-means / Lloyd's method
 - PCA
 - EM / GMMs
- Neural Networks
 - Feedforward Neural Nets
 - Basic architectures
 - Backpropagation
 - CNNs
- Graphical Models
 - Bayesian Networks
 - HMMs
 - Learning and Inference
- Learning Theory
 - Statistical Estimation (covered right before midterm)
 - PAC Learning
- Other Learning Paradigms
 - Matrix Factorization
 - Reinforcement Learning
 - Information Theory

ML Big Picture

Learning Paradigms:

What data is available and when? What form of prediction?

- supervised learning
- unsupervised learning
- semi-supervised learning
- reinforcement learning
- active learning
- imitation learning
- domain adaptation
- online learning
- density estimation
- recommender systems
- feature learning
- manifold learning
- dimensionality reduction
- ensemble learning
- distant supervision
- hyperparameter optimization

Theoretical Foundations:

What principles guide learning?

- ☐ probabilistic
- ☐ information theoretic
- ☐ evolutionary search
- ☐ ML as optimization

Problem Formulation:

What is the structure of our output prediction?

boolean	Binary Classification
categorical	Multiclass Classification
ordinal	Ordinal Classification
real	Regression
ordering	Ranking
multiple discrete	Structured Prediction
multiple continuous	(e.g. dynamical systems)
both discrete & cont.	(e.g. mixed graphical models)

Facets of Building ML Systems:

How to build systems that are robust, efficient, adaptive, effective?

1. Data prep
2. Model selection
3. Training (optimization / search)
4. Hyperparameter tuning on validation data
5. (Blind) Assessment on test data

Big Ideas in ML:

Which are the ideas driving development of the field?

- inductive bias
- generalization / overfitting
- bias-variance decomposition
- generative vs. discriminative
- deep nets, graphical models
- PAC learning
- distant rewards

Application Areas

Key challenges?

NLP, Speech, Computer Vision, Robotics, Medicine, Search

Machine Learning & Ethics

What ethical responsibilities do we have as machine learning experts?

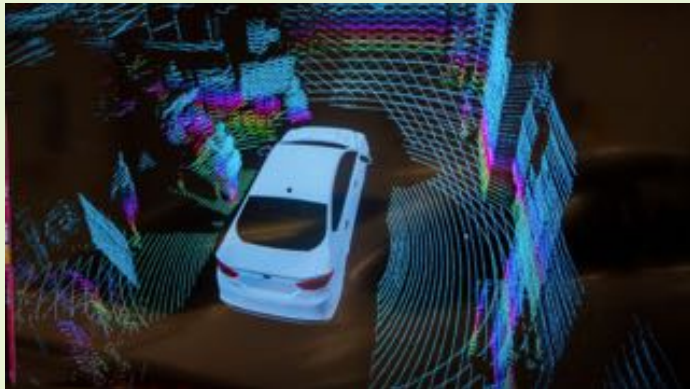
Some topics that we won't cover are probably deserve an entire course

If our search results for news are optimized for ad revenue, might they reflect gender / racial / socio-economic biases?



<http://bing.com/>

<http://arstechnica.com/>



How do autonomous vehicles make decisions when all of the outcomes are likely to be negative?

Should restrictions be placed on intelligent agents that are capable of interacting with the world?



<http://vizdoom.cs.put.edu.pl/>

DEFINING LEARNING PROBLEMS

Well-Posed Learning Problems

Three components $\langle T, P, E \rangle$:

1. Task, T
2. Performance measure, P
3. Experience, E

Definition of learning:

A computer program **learns** if its performance at tasks in T , as measured by P , improves with experience E .

Example Learning Problems

3. Learning to beat the masters at **chess**

1. Task, T :
2. Performance measure, P :
3. Experience, E :

Example Learning Problems

4. Learning to **respond to voice commands (Siri)**

1. Task, T :
2. Performance measure, P :
3. Experience, E :

Capturing the Knowledge of Experts



Solution #1: Expert Systems

- Over 20 years ago, we had rule based systems
- Ask the expert to
 1. Obtain a PhD in Linguistics
 2. Introspect about the structure of their native language
 3. Write down the rules they devise

Give me directions to Starbucks

If: "give me directions to X"
Then: `directions(here, nearest(X))`

How do I get to Starbucks?

If: "how do i get to X"
Then: `directions(here, nearest(X))`

Where is the nearest Starbucks?

If: "where is the nearest X"
Then: `directions(here, nearest(X))`

Capturing the Knowledge of Experts



Solution #1: Expert Systems

- Over 20 years ago, we had rule based systems
- Ask the expert to
 1. Obtain a PhD in Linguistics
 2. Introspect about the structure of their native language
 3. Write down the rules they devise

I need directions to Starbucks

If: "I need directions to X"
Then: `directions(here, nearest(X))`

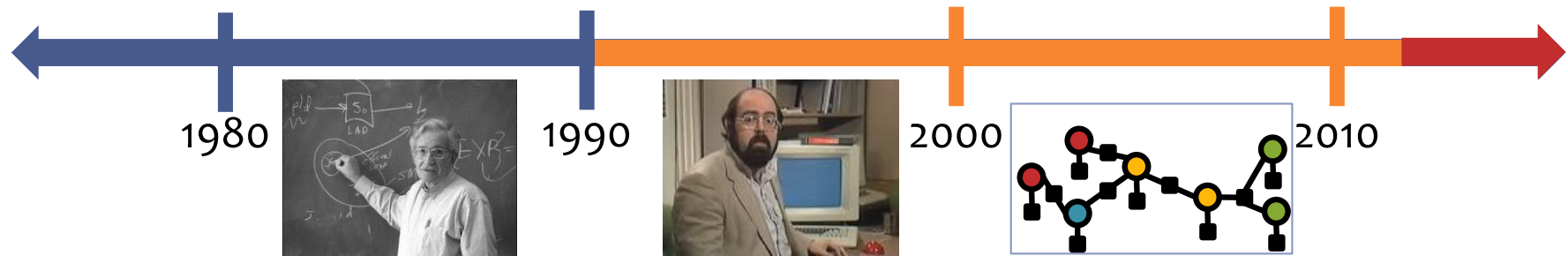
Starbucks directions

If: "X directions"
Then: `directions(here, nearest(X))`

Is there a Starbucks nearby?

If: "Is there an X nearby"
Then: `directions(here, nearest(X))`

Capturing the Knowledge of Experts



Solution #2: Annotate Data and Learn

- Experts:
 - **Very good at** answering questions about specific cases
 - **Not very good at** telling **HOW** they do it
- 1990s: So why not just have them tell you what they do on **SPECIFIC CASES** and then let **MACHINE LEARNING** tell you how to come to the same decisions that they did

Capturing the Knowledge of Experts



Solution #2: Annotate Data and Learn

1. Collect raw sentences $\{x_1, \dots, x_n\}$
2. Experts annotate their meaning $\{y_1, \dots, y_n\}$

x_1 : How do I get to Starbucks?

y_1 : `directions(here,
nearest(Starbucks))`

x_2 : Show me the closest Starbucks

y_2 : `map(nearest(Starbucks))`

x_3 : Send a text to John that I'll be late

y_3 : `txtmsg(John, I'll be late)`

x_4 : Set an alarm for seven in the morning

y_4 : `setalarm(7:00AM)`

Example Learning Problems

4. Learning to **respond to voice commands (Siri)**

1. Task, T :

predicting action from speech

2. Performance measure, P :

percent of correct actions taken in user pilot study

3. Experience, E :

examples of (speech, action) pairs

Problem Formulation

- Often, the same task can be formulated in more than one way:
- Ex: Loan applications
 - creditworthiness/score (regression)
 - probability of default (density estimation)
 - loan decision (classification)

Problem Formulation:

What is the structure of our output prediction?

boolean	Binary Classification	}
categorical	Multiclass Classification	
ordinal	Ordinal Classification	
real	Regression	
ordering	Ranking	
multiple discrete	Structured Prediction	
multiple continuous	(e.g. dynamical systems)	
both discrete & cont.	(e.g. mixed graphical models)	

Well-posed Learning Problems

In-Class Exercise

1. Select a **task**, T
2. Identify **performance measure**, P
3. Identify **experience**, E
4. Report ideas back to rest of class

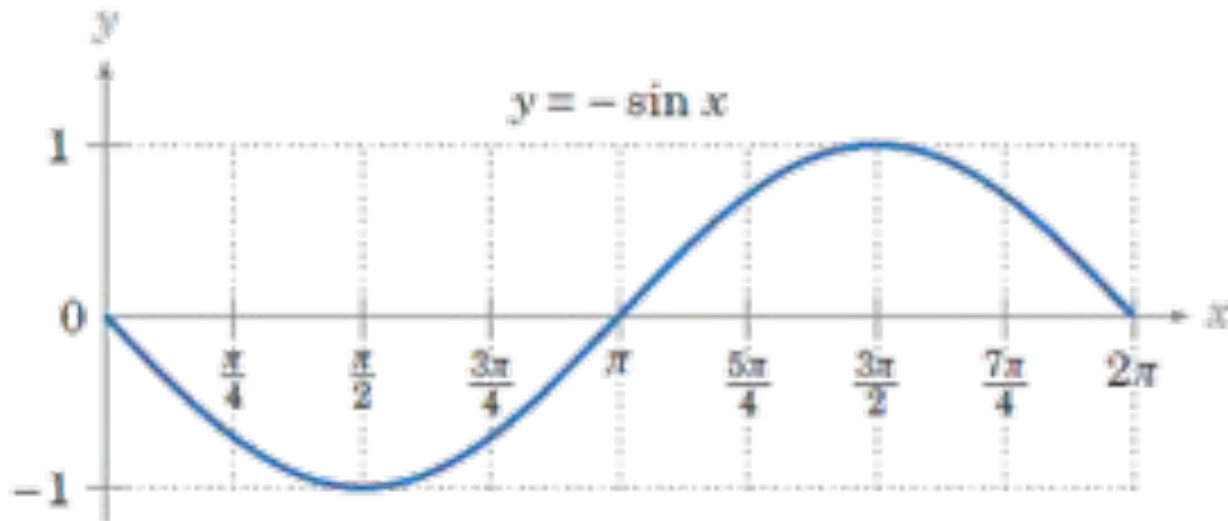
Example Tasks

- Identify objects in an image
- Translate from one human language to another
- Recognize speech
- Assess risk (e.g. in loan application)
- Make decisions (e.g. in loan application)
- Assess potential (e.g. in admission decisions)
- Categorize a complex situation (e.g. medical diagnosis)
- Predict outcome (e.g. medical prognosis, stock prices, inflation, temperature)
- Predict events (default on loans, quitting school, war)
- Plan ahead under perfect knowledge (chess)
- Plan ahead under partial knowledge (Poker, Bridge)

FUNCTION APPROXIMATION

Function Approximation

Quiz: Implement a simple function which returns $\sin(x)$.



A few constraints are imposed:

1. You can't call any other trigonometric functions
2. You *can* call an existing implementation of $\sin(x)$ a few times (e.g. 100) to test your solution
3. You only need to evaluate it for x in $[0, 2\pi]$

ML as Function Approximation

Chalkboard

– ML as Function Approximation

- Problem setting
- Input space
- Output space
- Unknown target function
- Hypothesis space
- Training examples

CLASSIFICATION

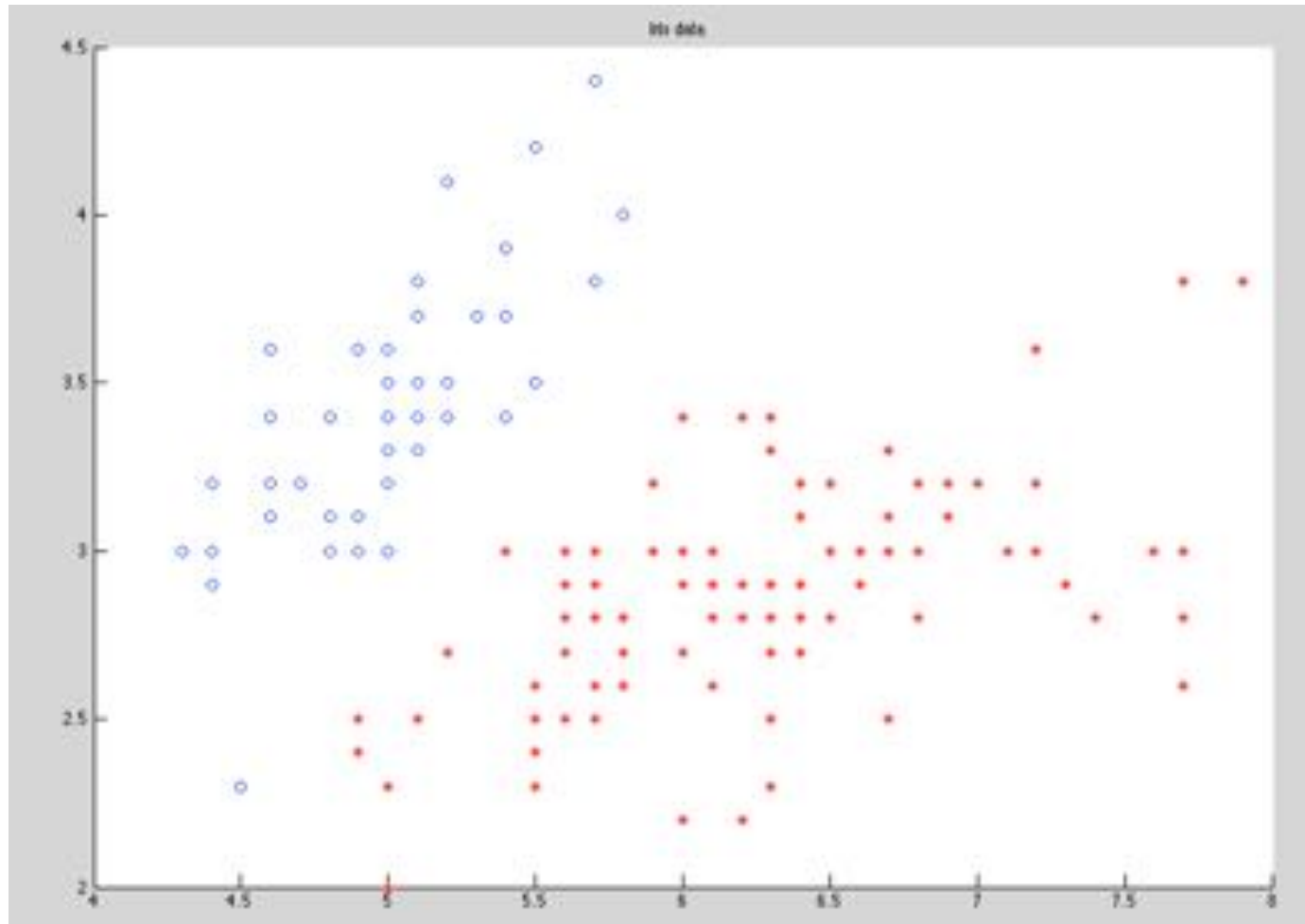


Fisher Iris Dataset

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



Classification

Chalkboard:

- Binary classification
- 2D examples
- Decision rules / hypotheses

K-NEAREST NEIGHBORS

k-Nearest Neighbors

Chalkboard:

- KNN for binary classification
- Distance functions
- Efficiency of KNN
- Inductive bias of KNN
- KNN Properties

KNN ON FISHER IRIS DATA

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

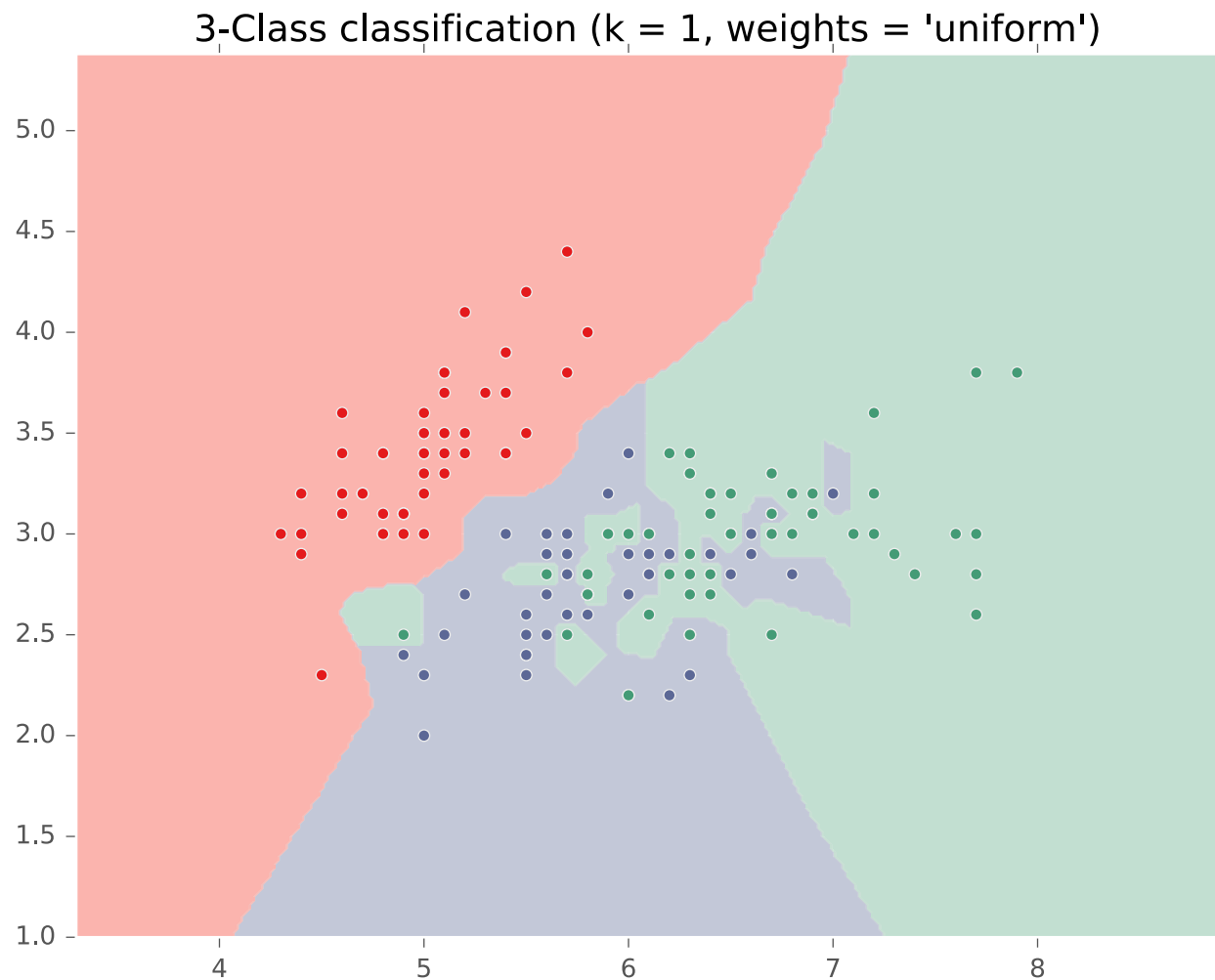


KNN on Fisher Iris Data



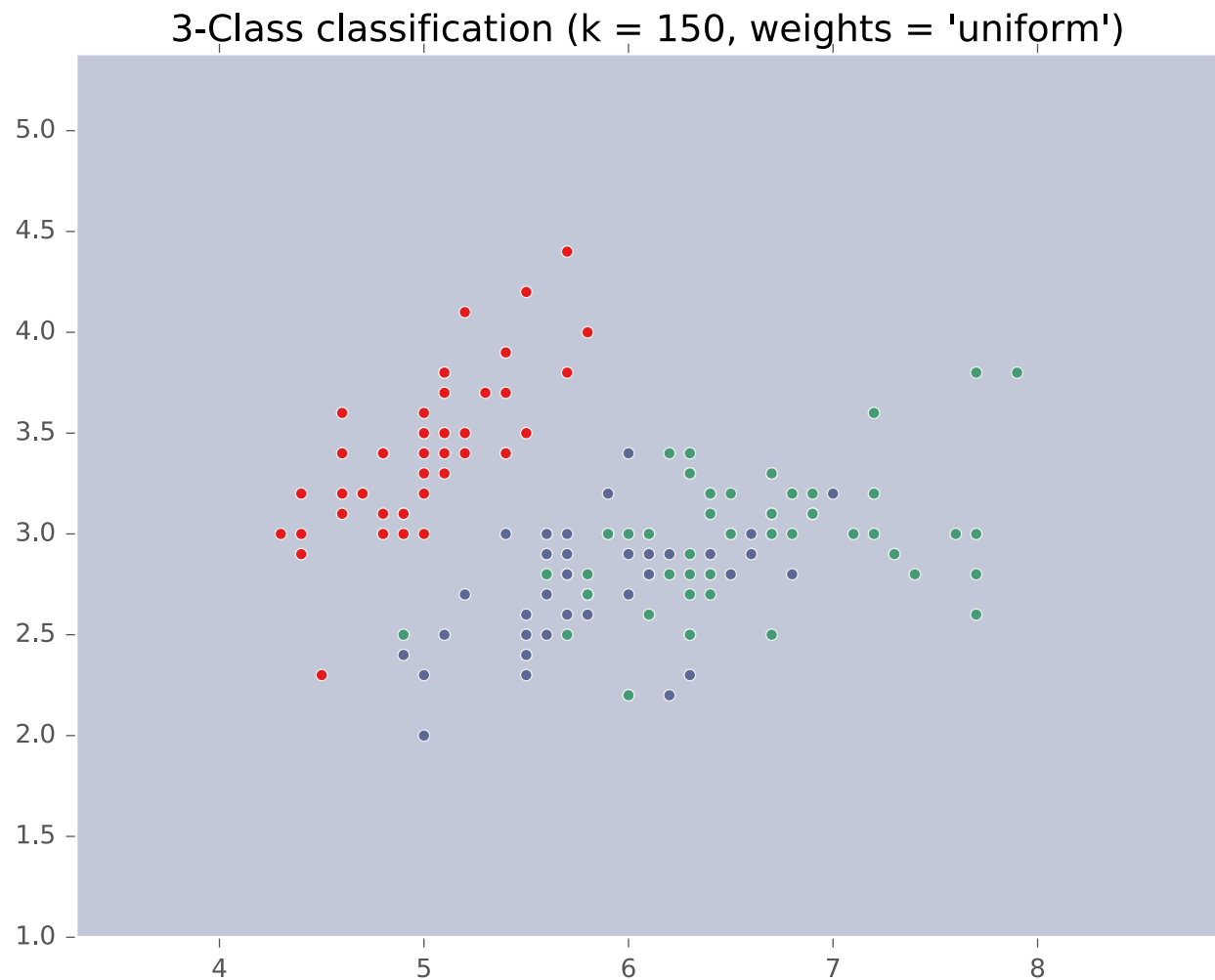
KNN on Fisher Iris Data

Special Case: Nearest Neighbor

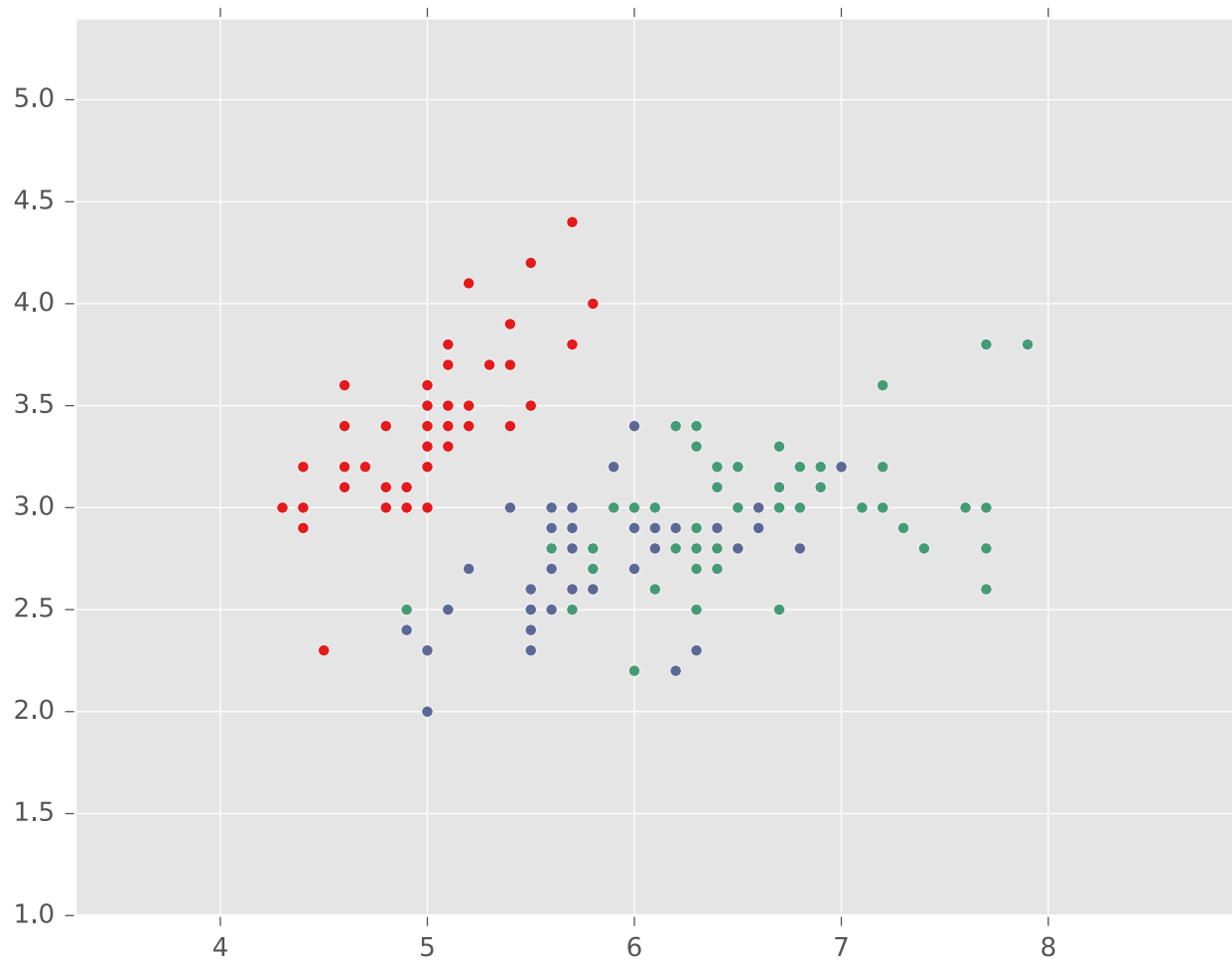


KNN on Fisher Iris Data

Special Case: Majority Vote

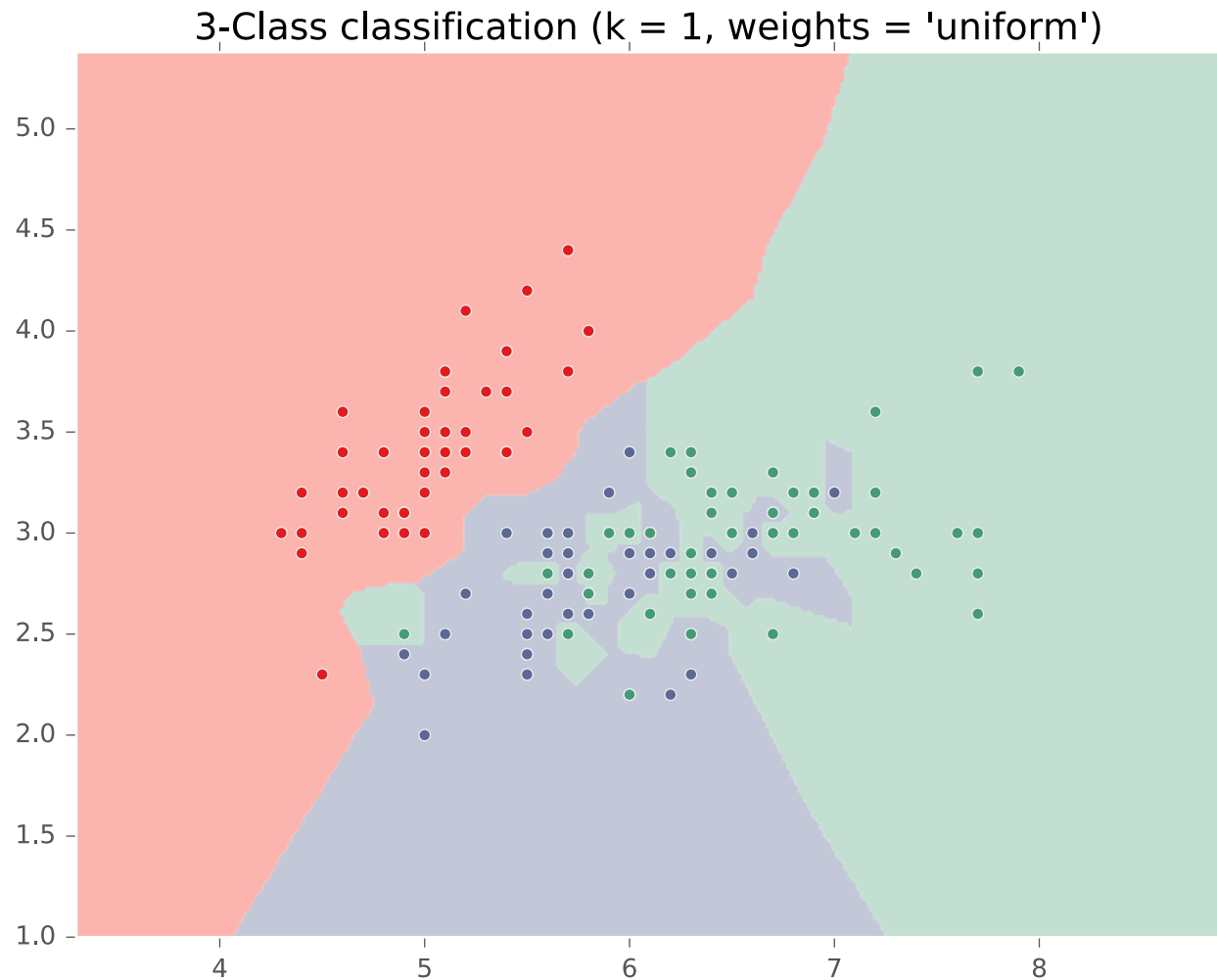


KNN on Fisher Iris Data

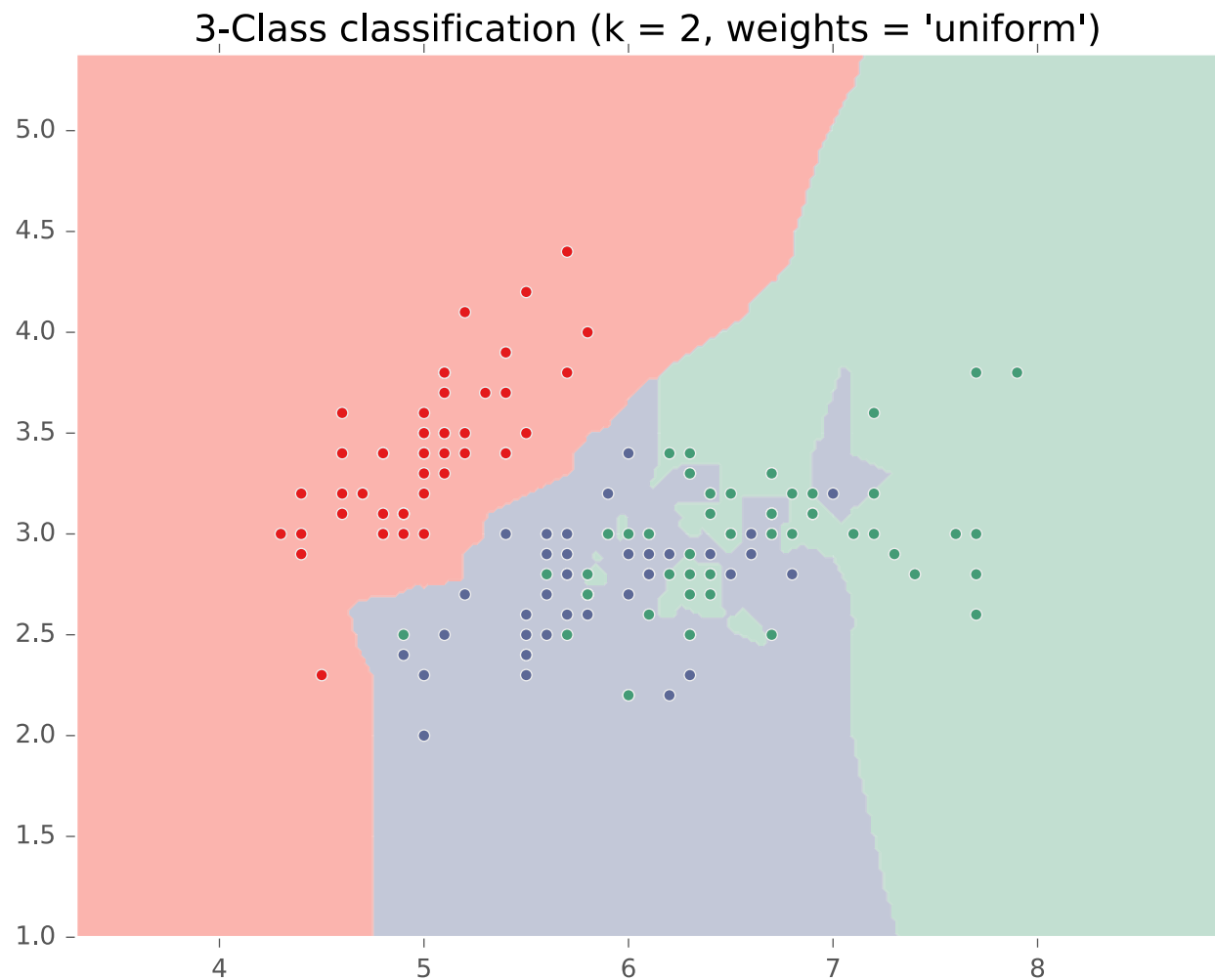


KNN on Fisher Iris Data

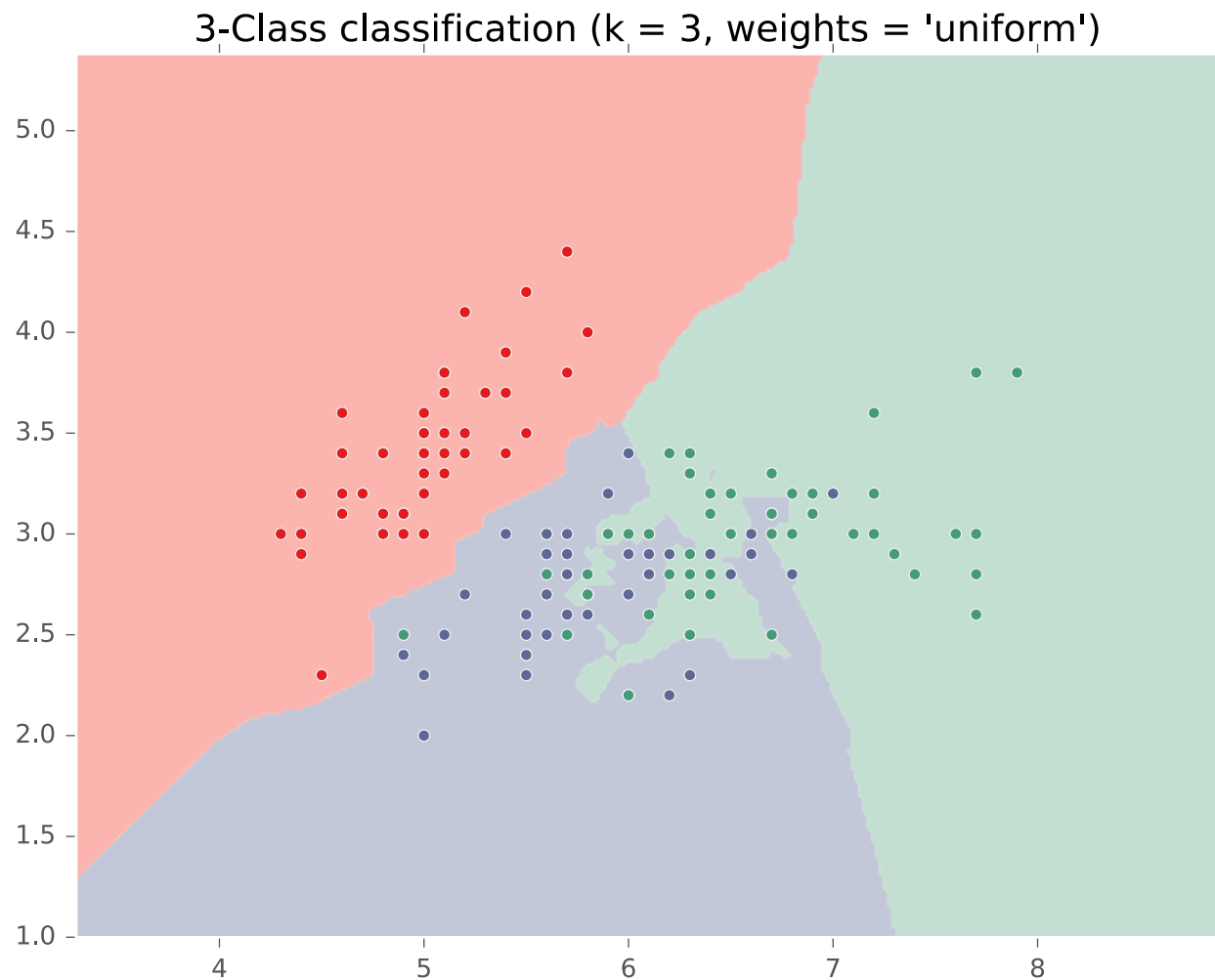
Special Case: Nearest Neighbor



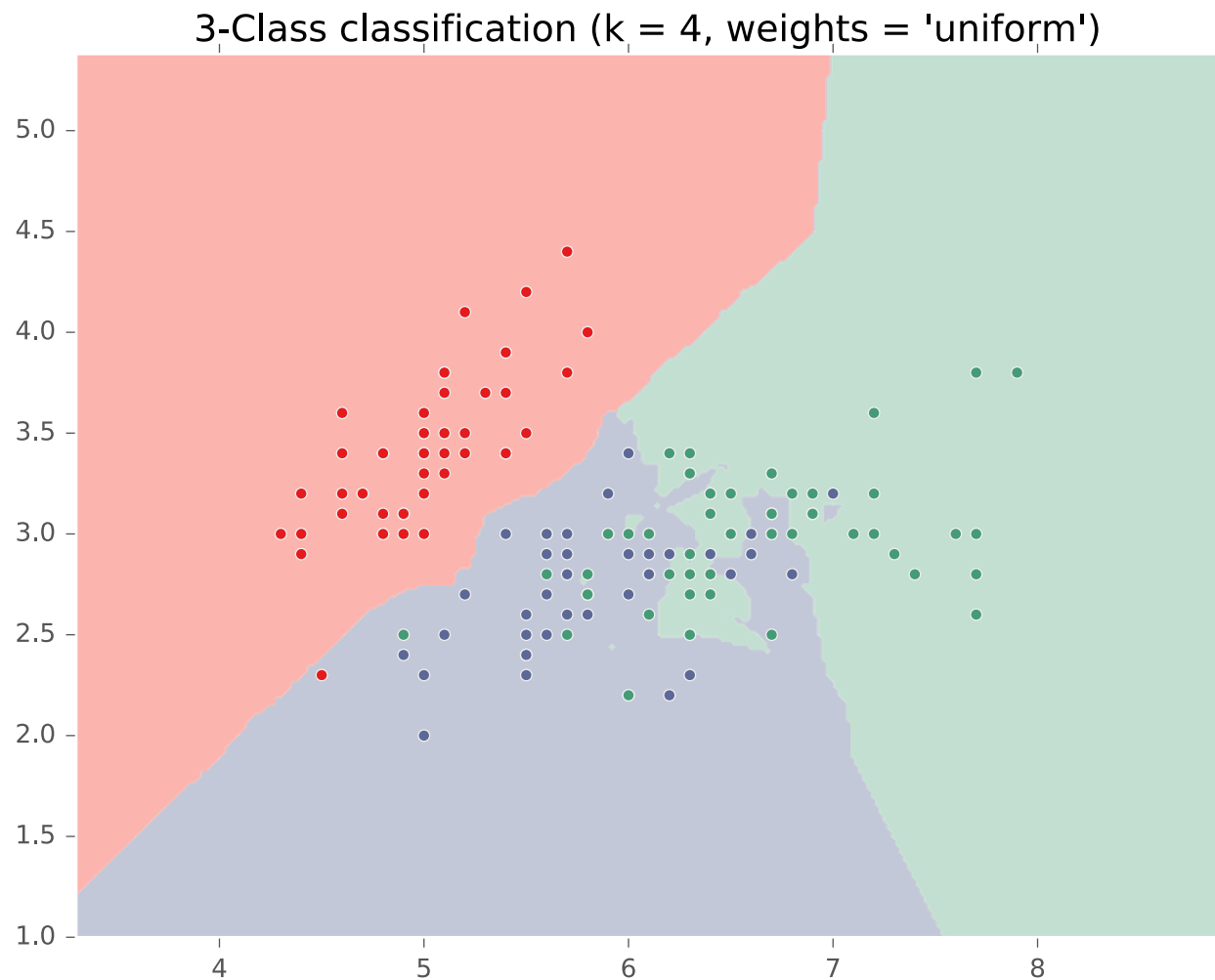
KNN on Fisher Iris Data



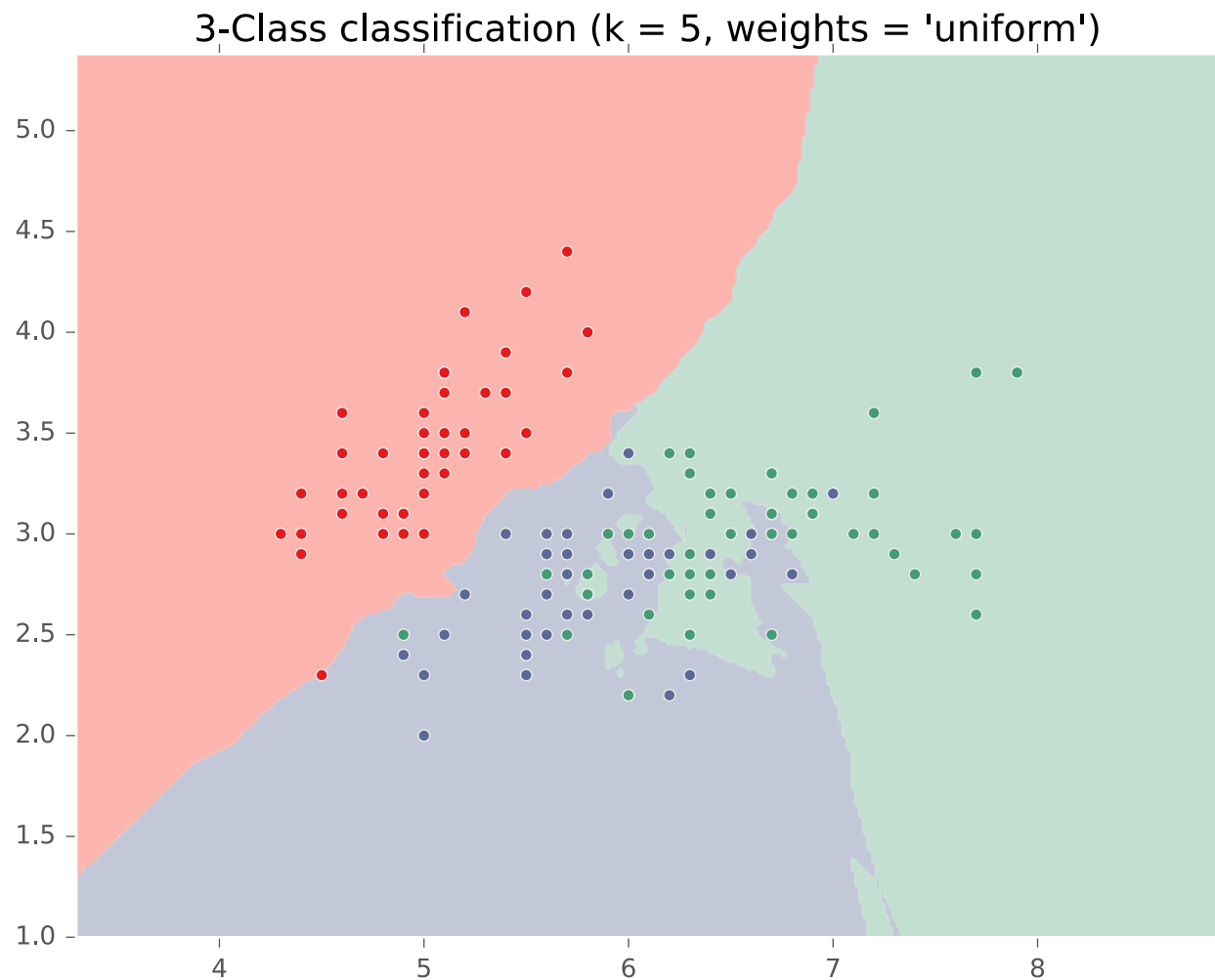
KNN on Fisher Iris Data



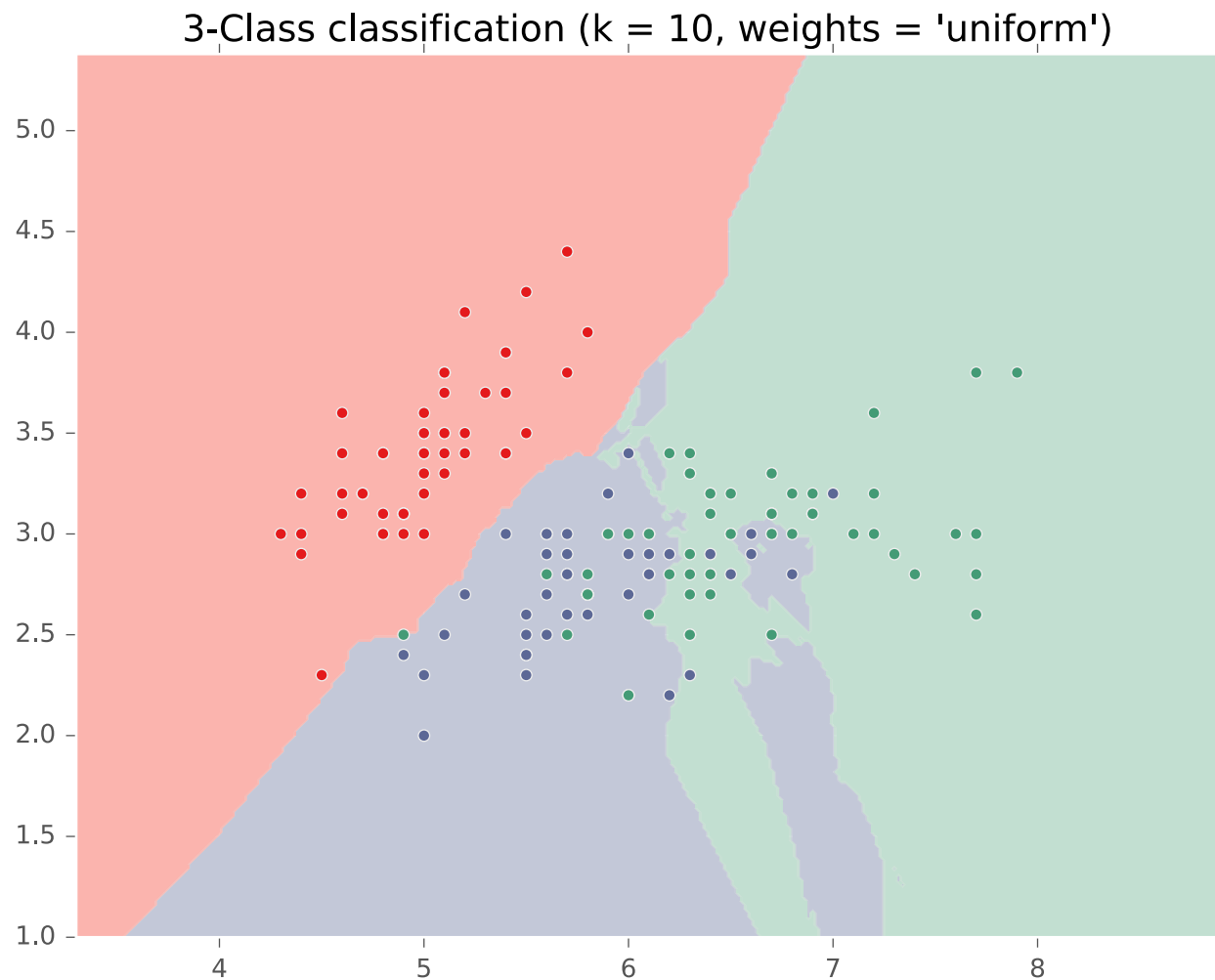
KNN on Fisher Iris Data



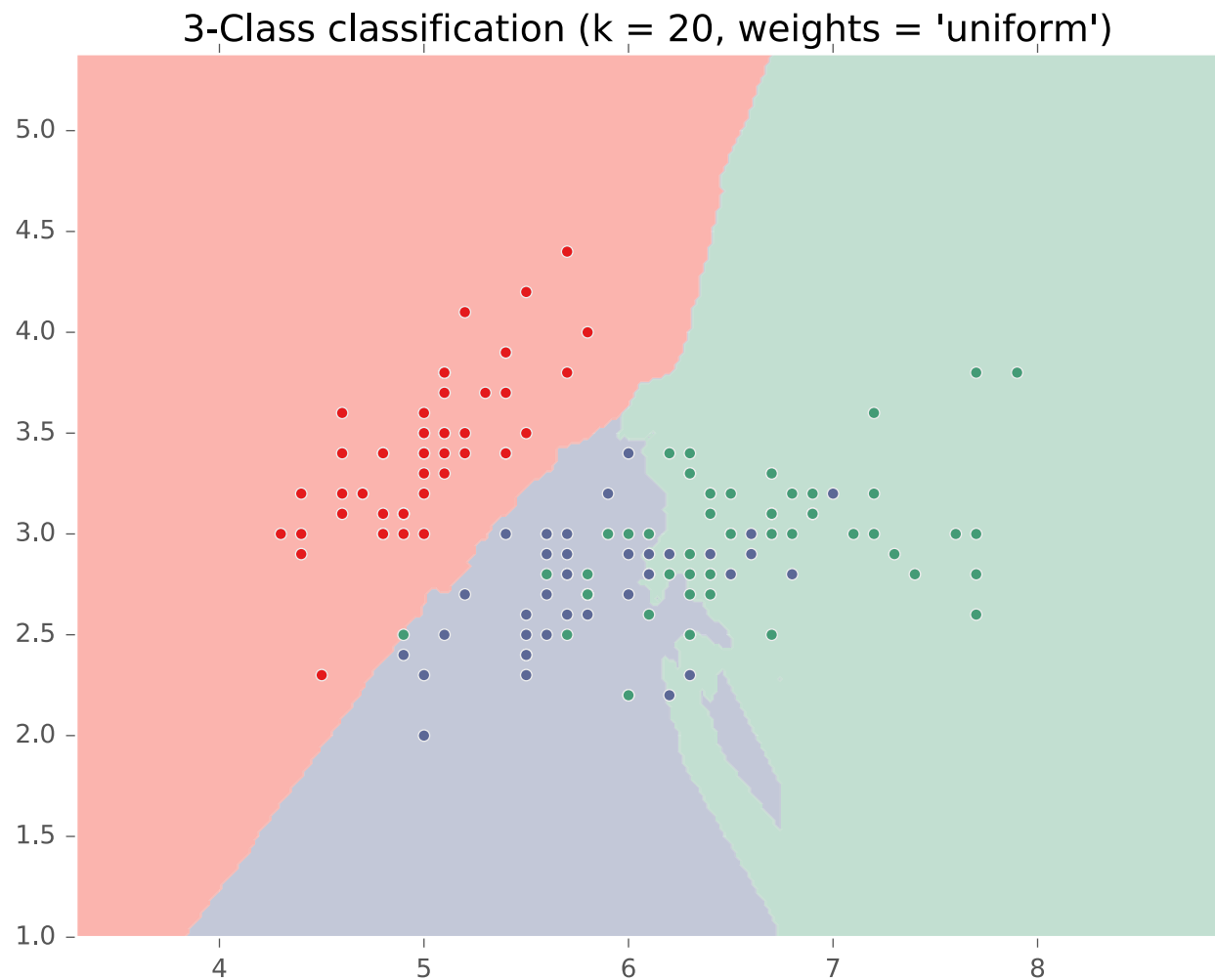
KNN on Fisher Iris Data



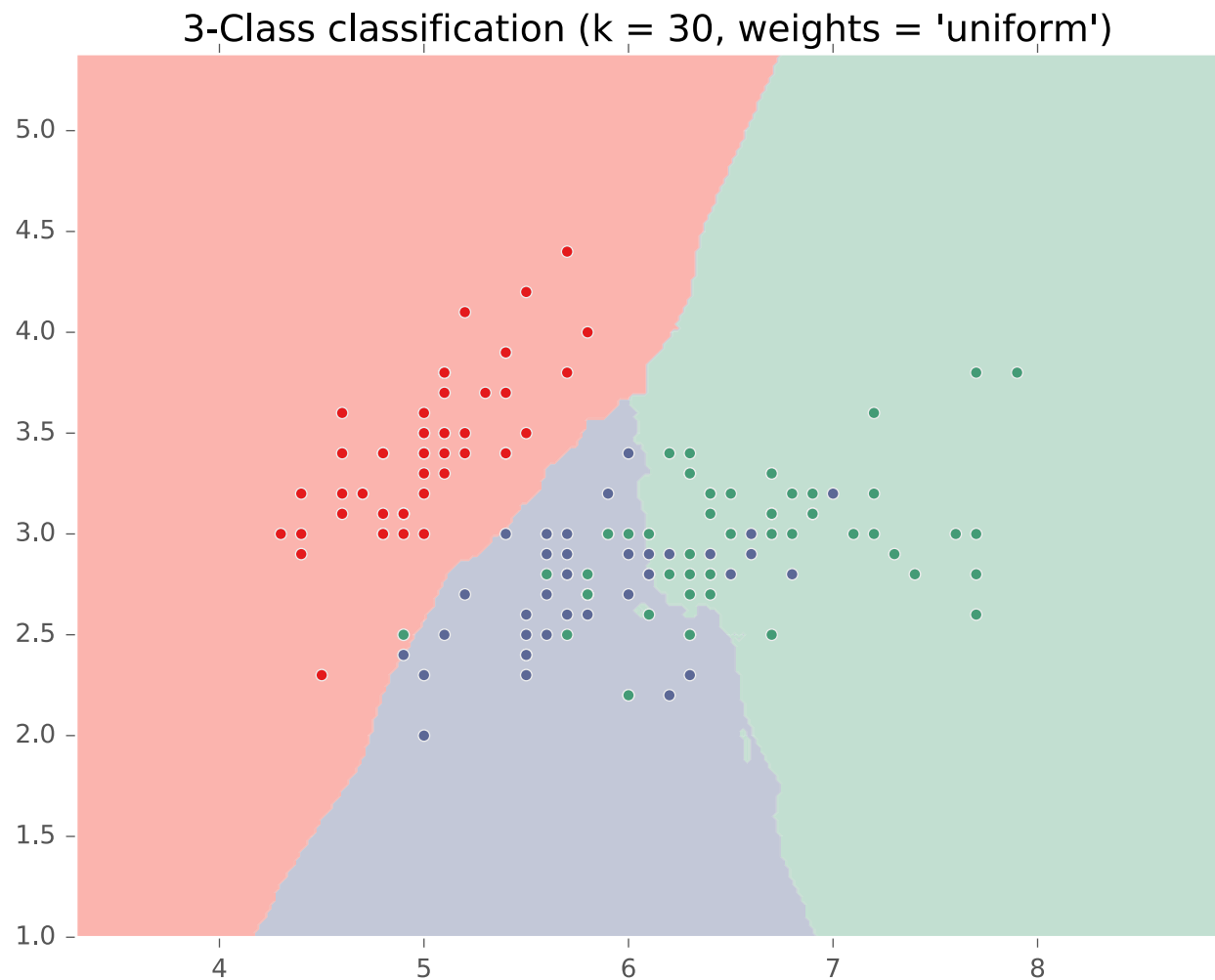
KNN on Fisher Iris Data



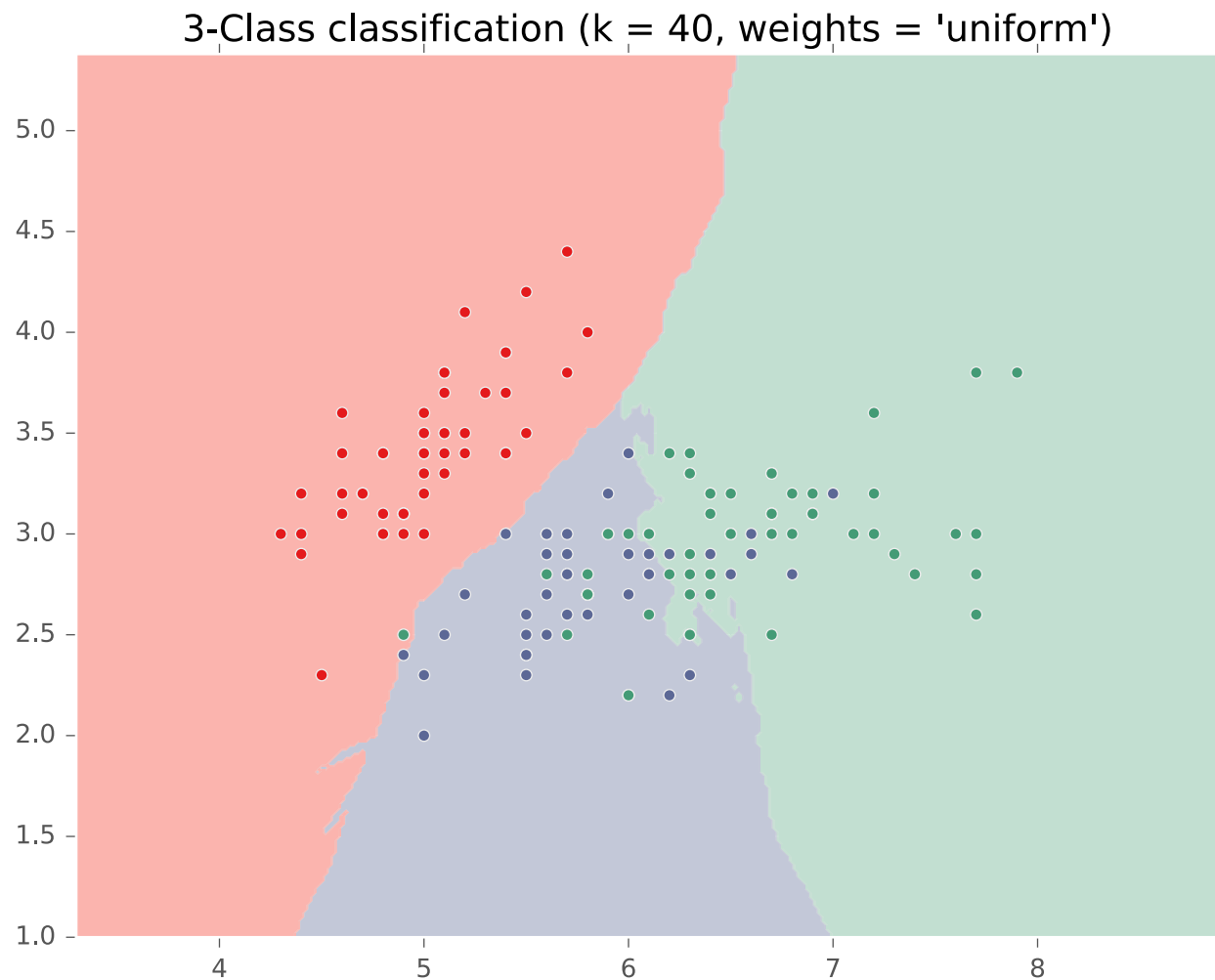
KNN on Fisher Iris Data



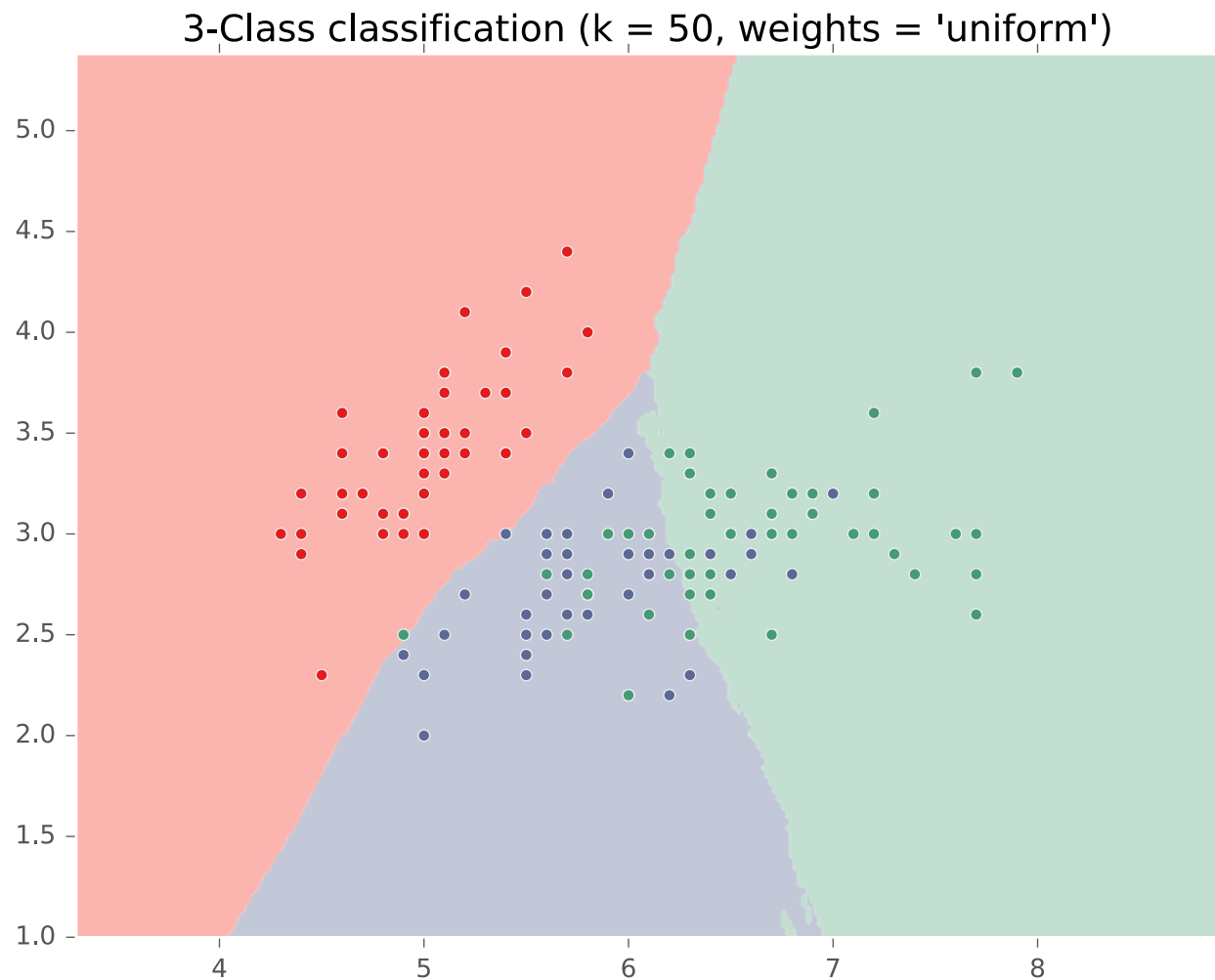
KNN on Fisher Iris Data



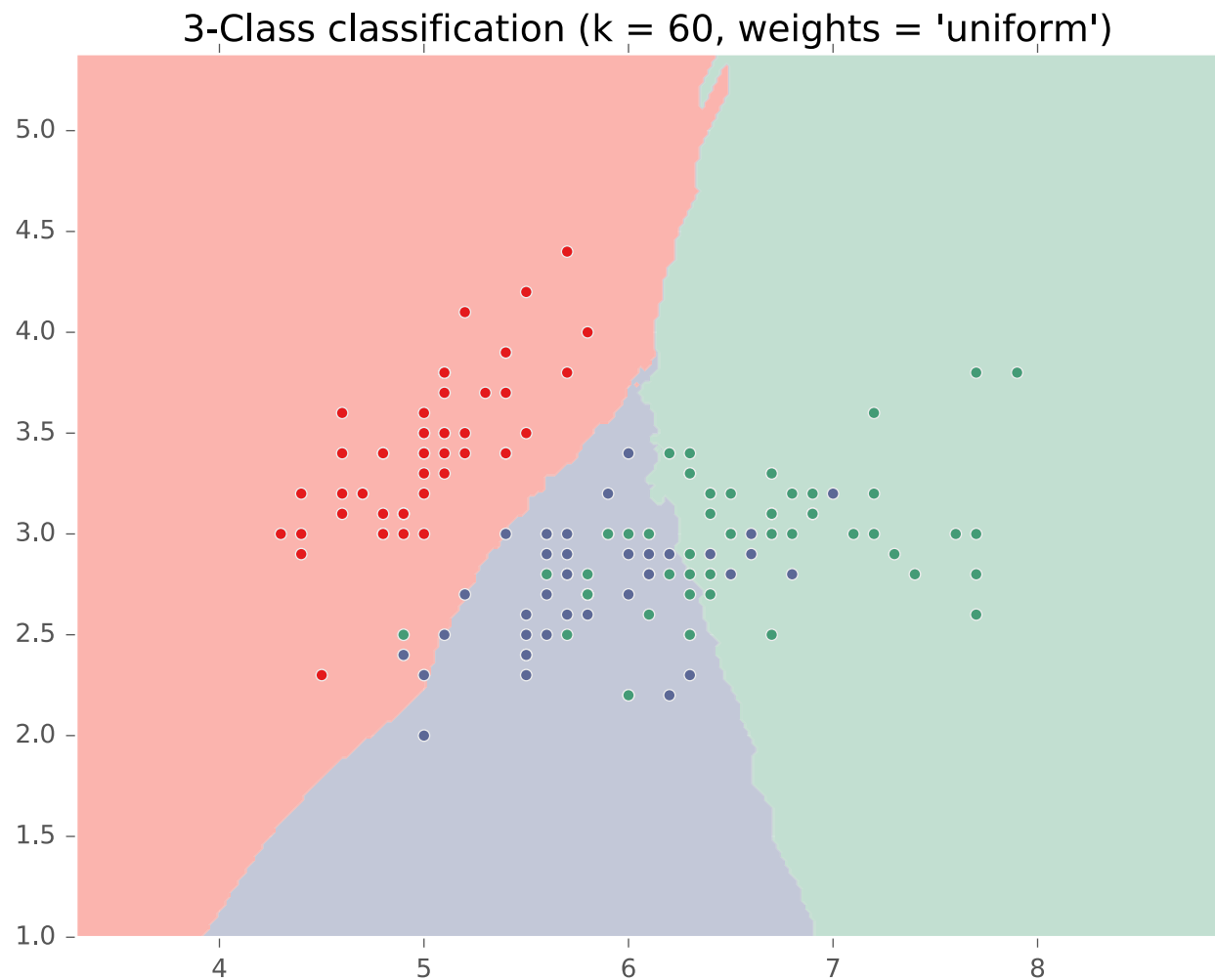
KNN on Fisher Iris Data



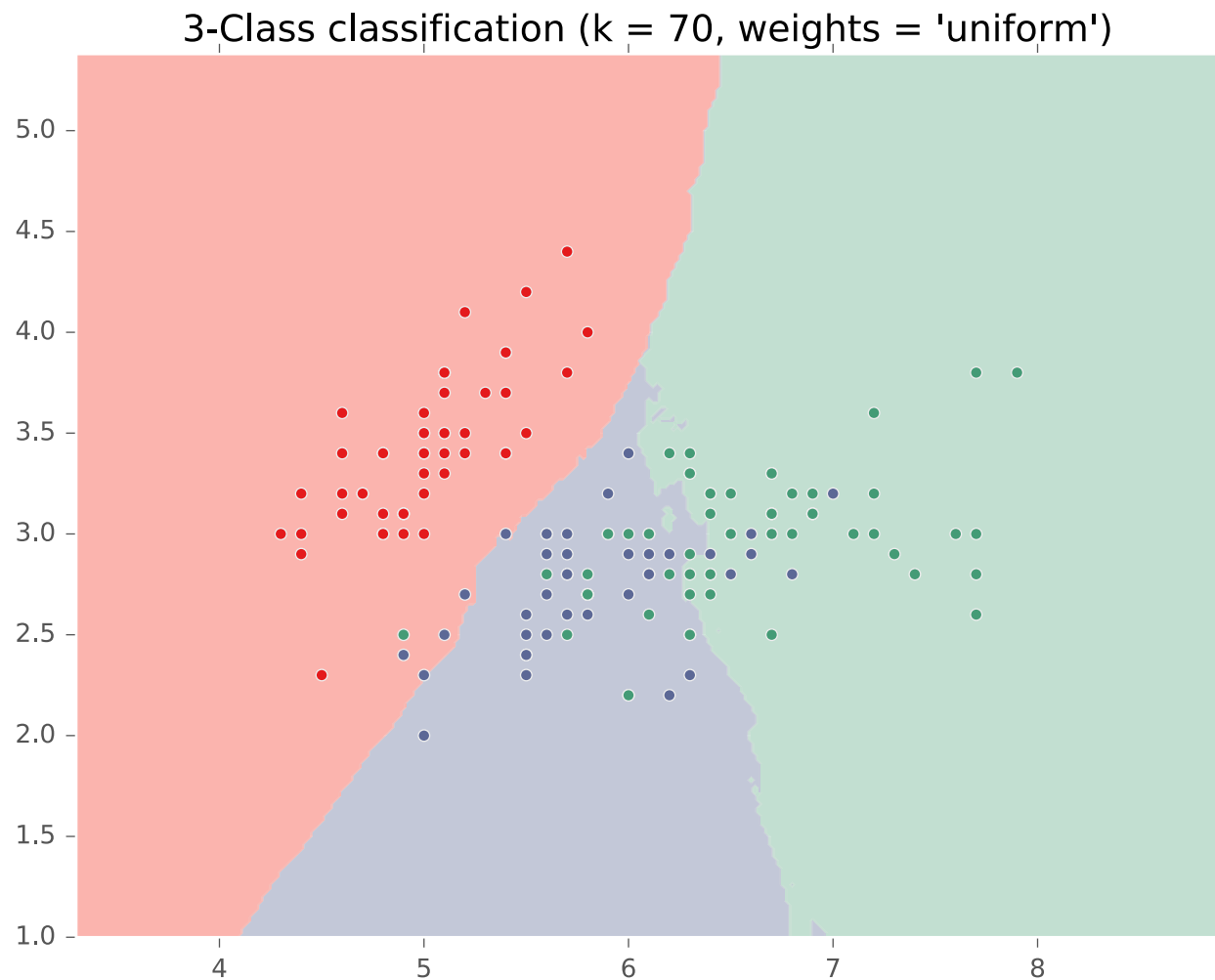
KNN on Fisher Iris Data



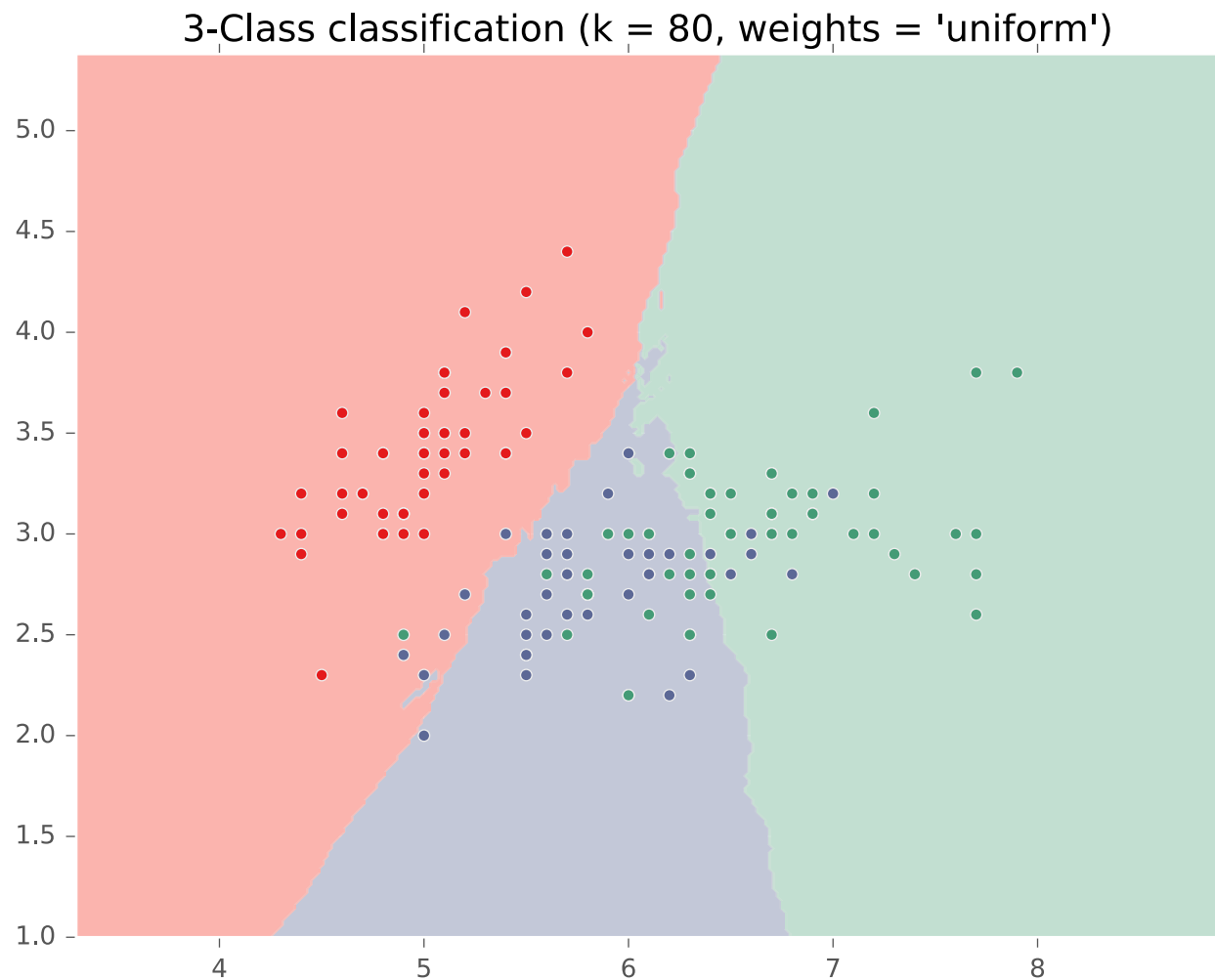
KNN on Fisher Iris Data



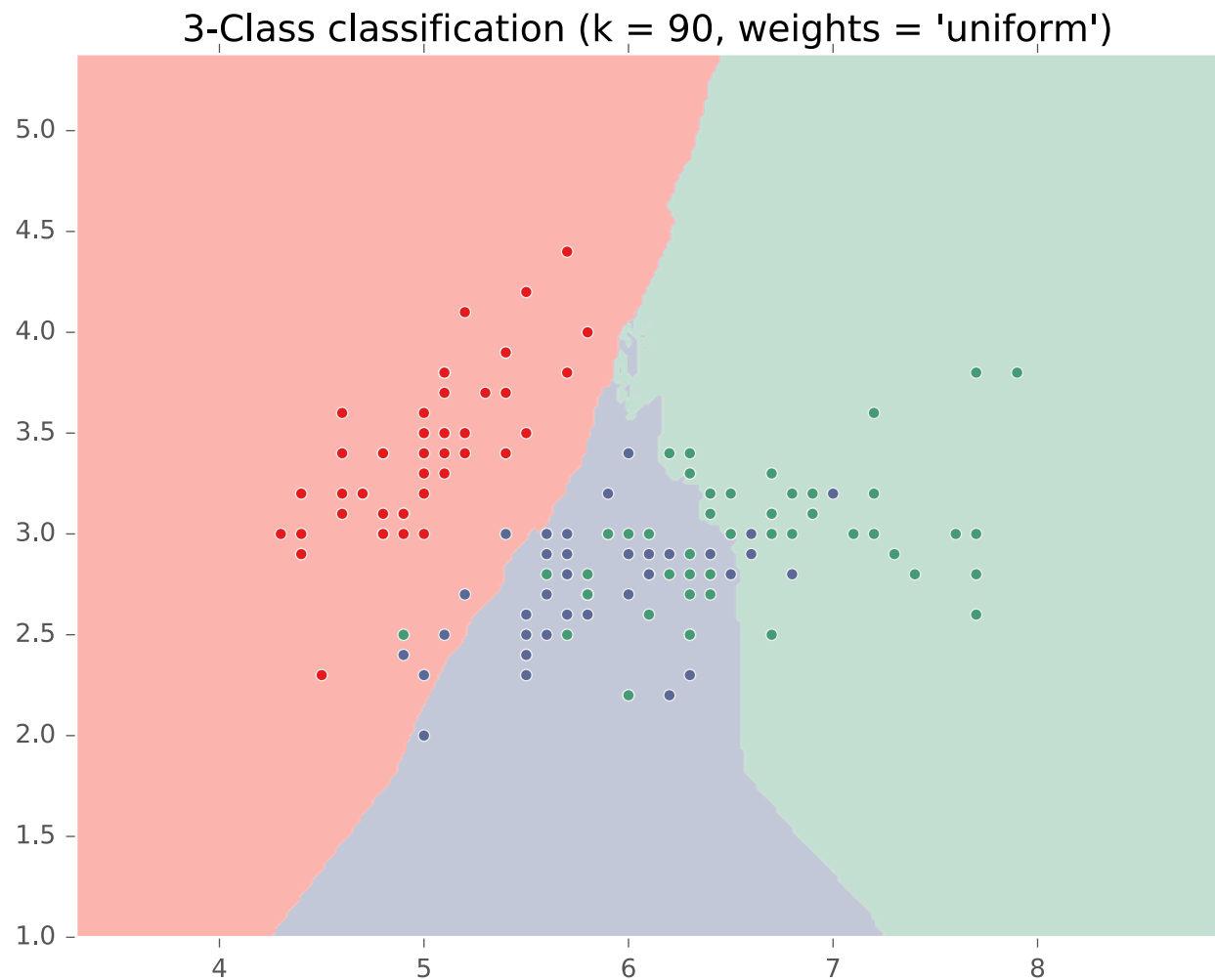
KNN on Fisher Iris Data



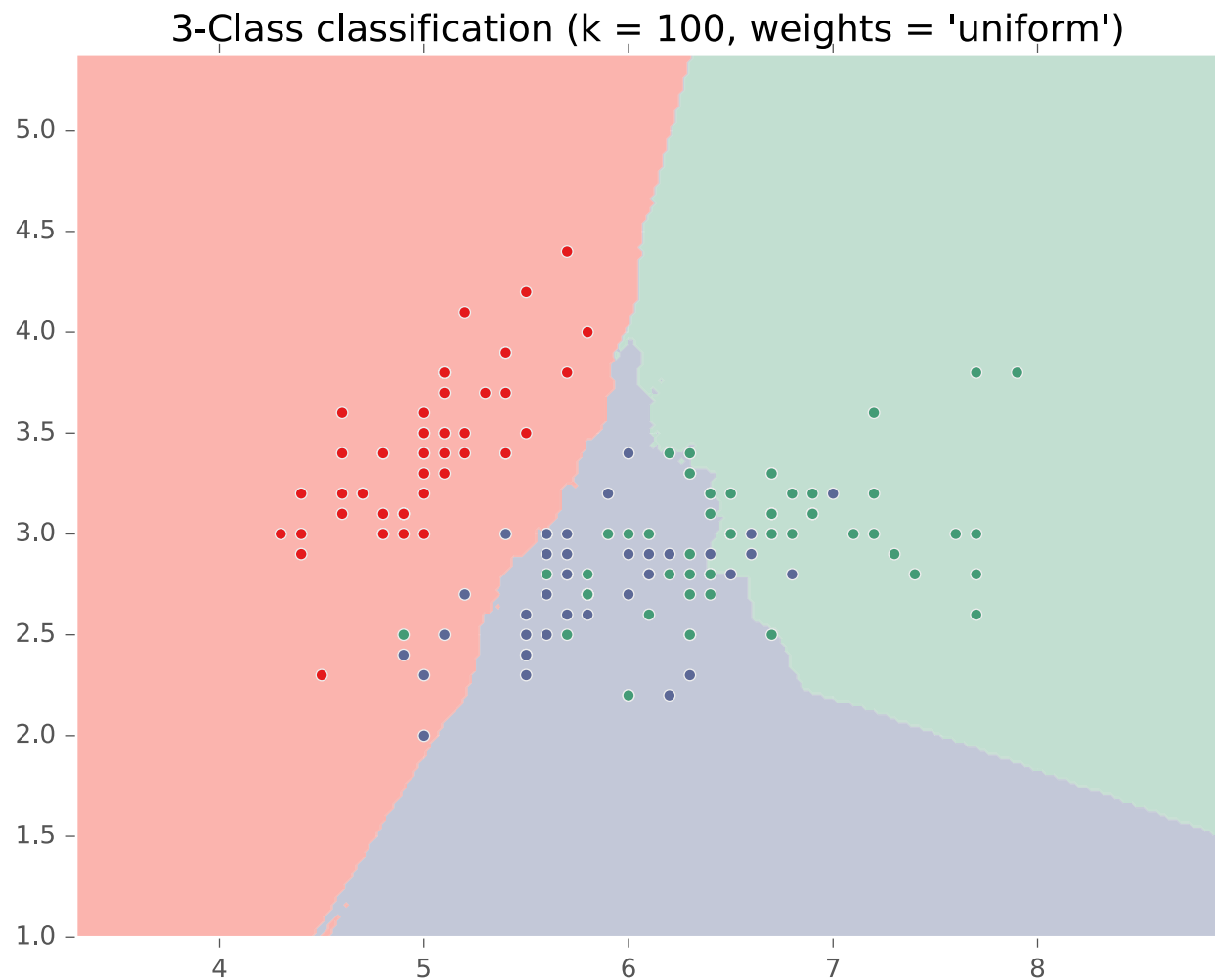
KNN on Fisher Iris Data



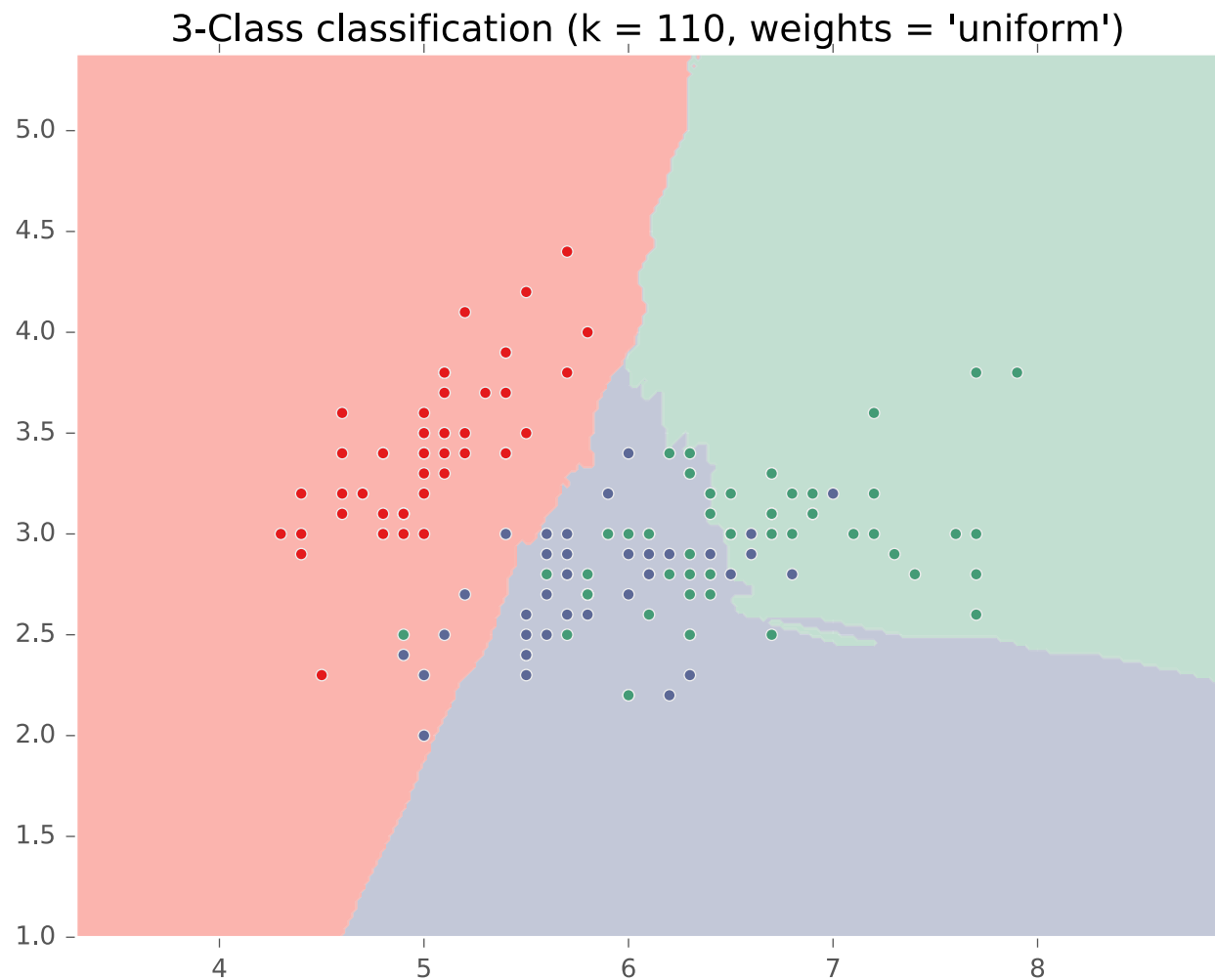
KNN on Fisher Iris Data



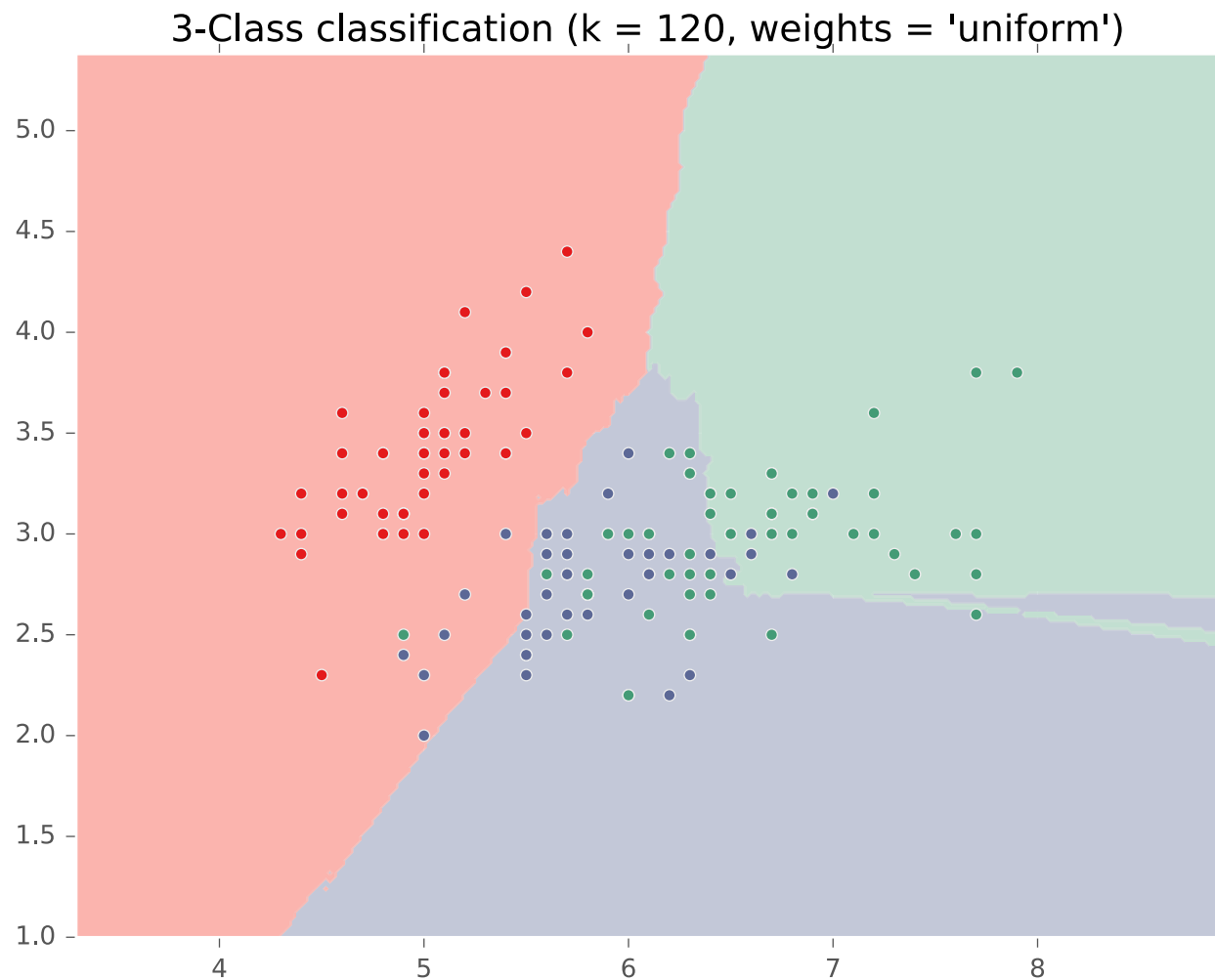
KNN on Fisher Iris Data



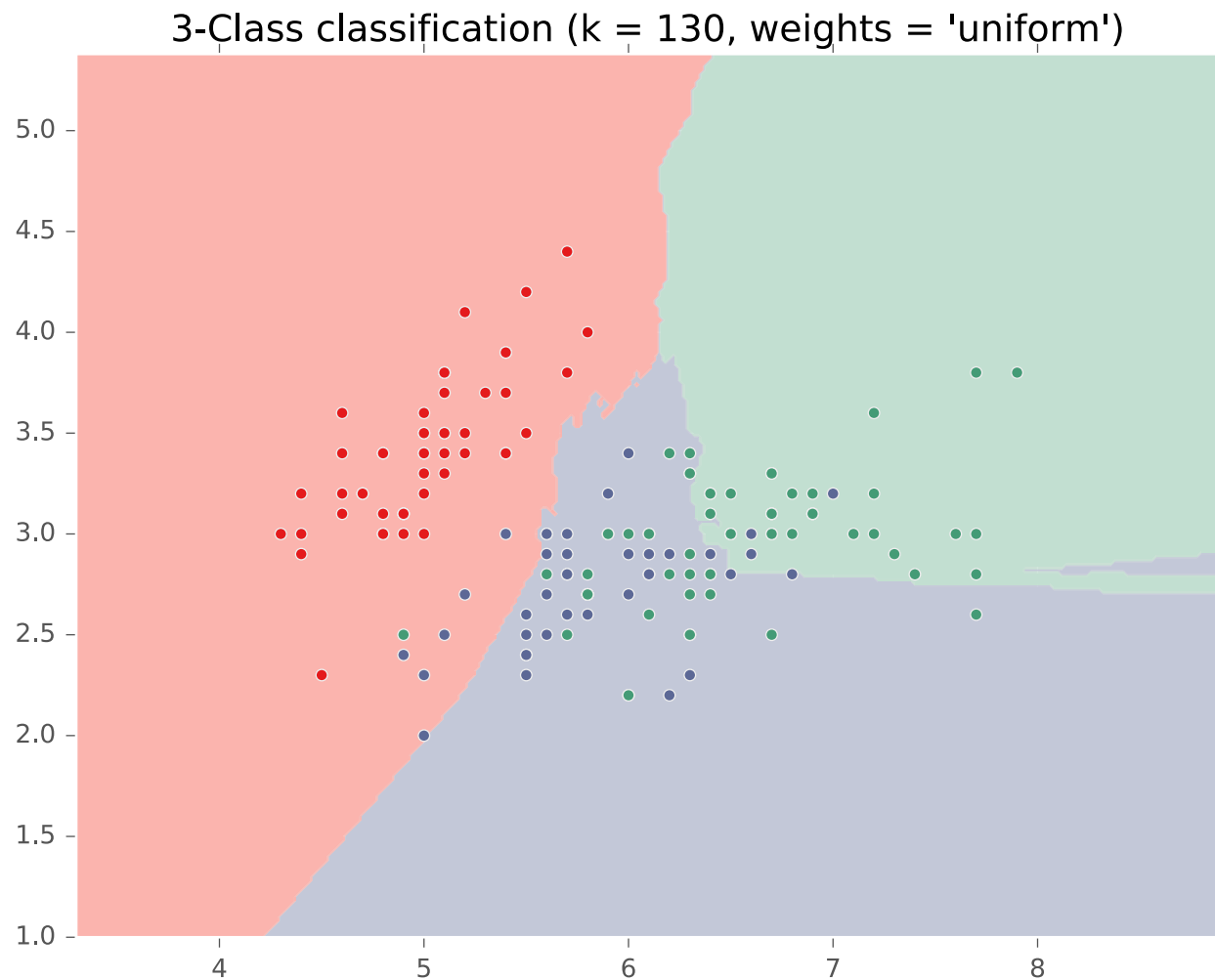
KNN on Fisher Iris Data



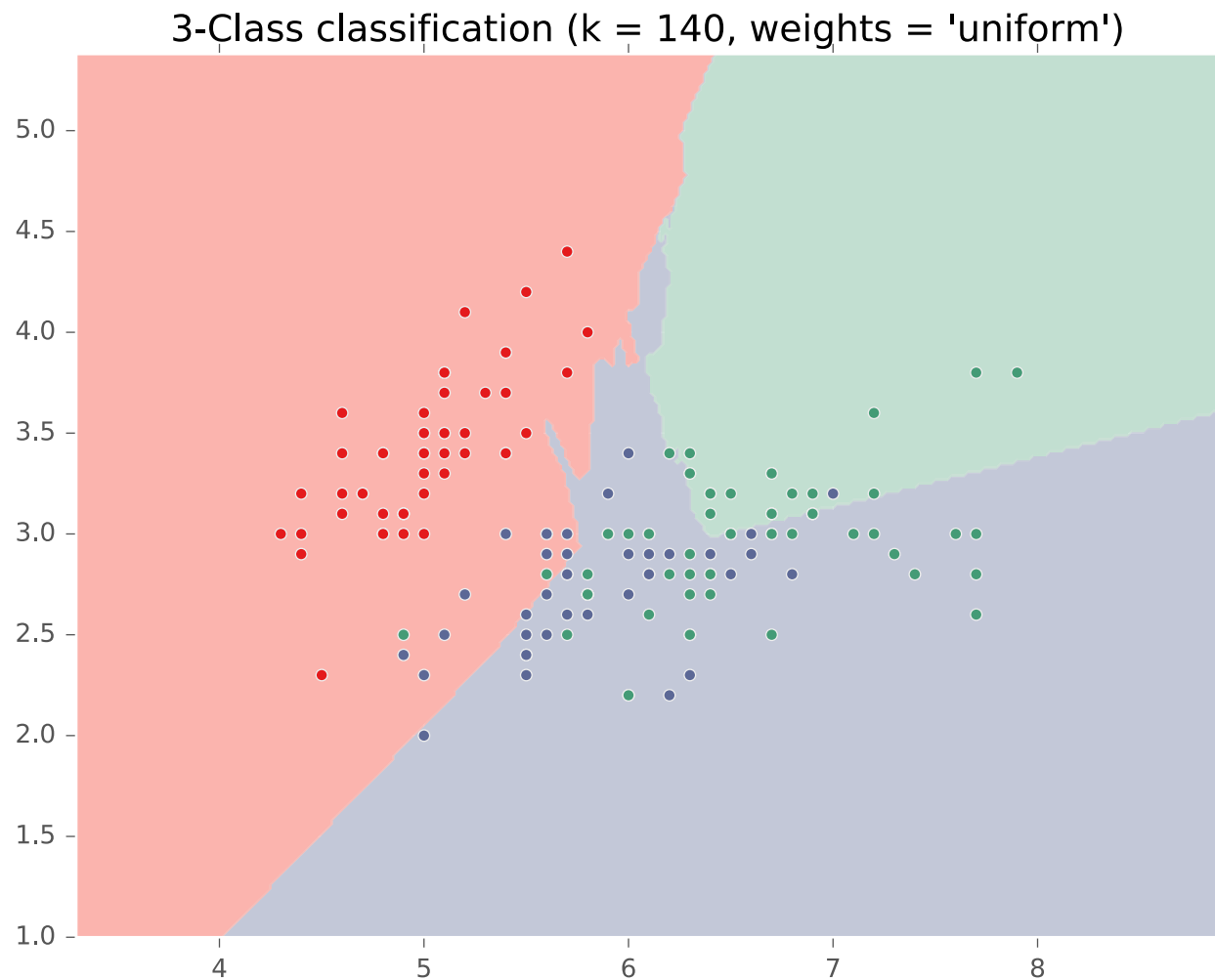
KNN on Fisher Iris Data



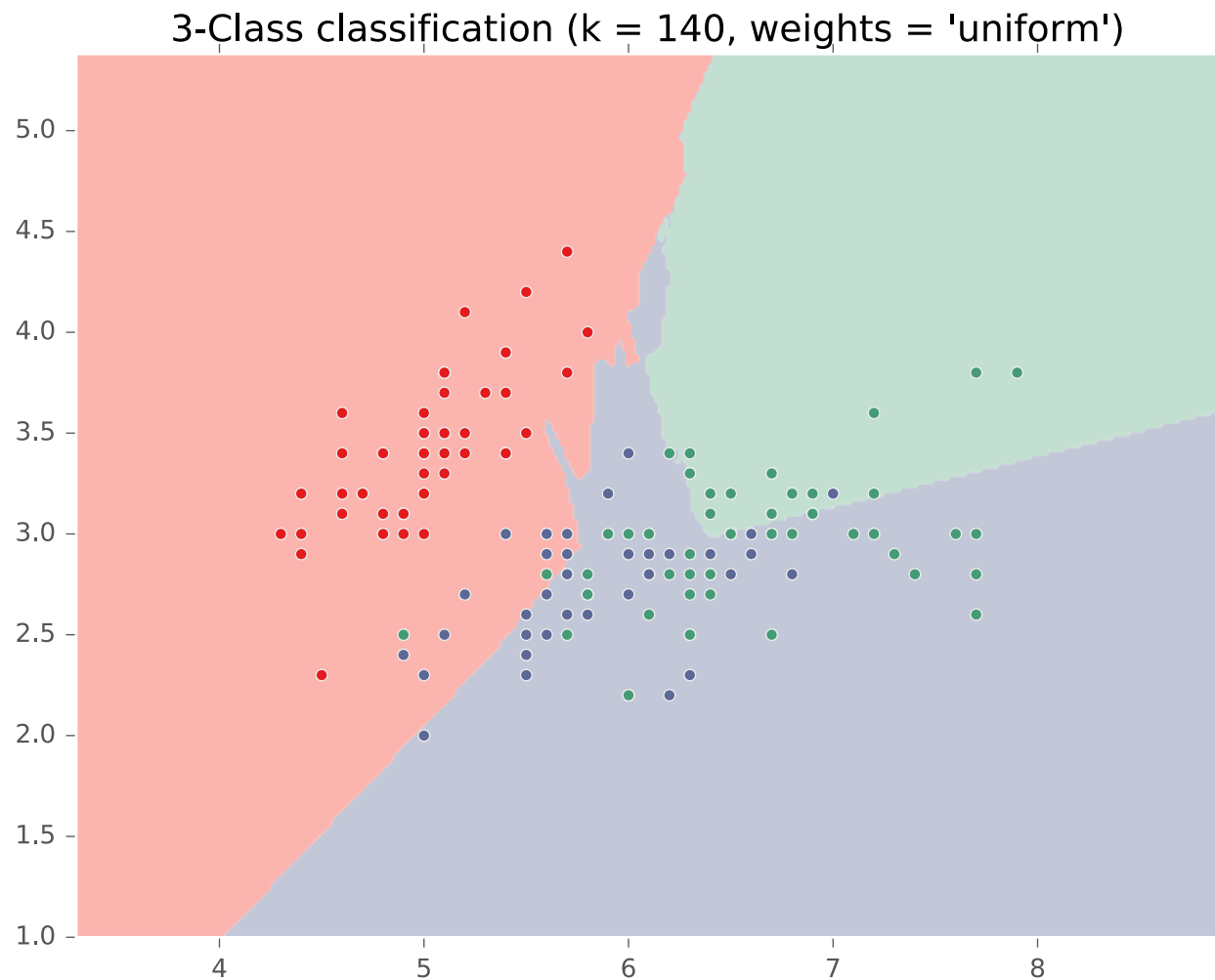
KNN on Fisher Iris Data



KNN on Fisher Iris Data

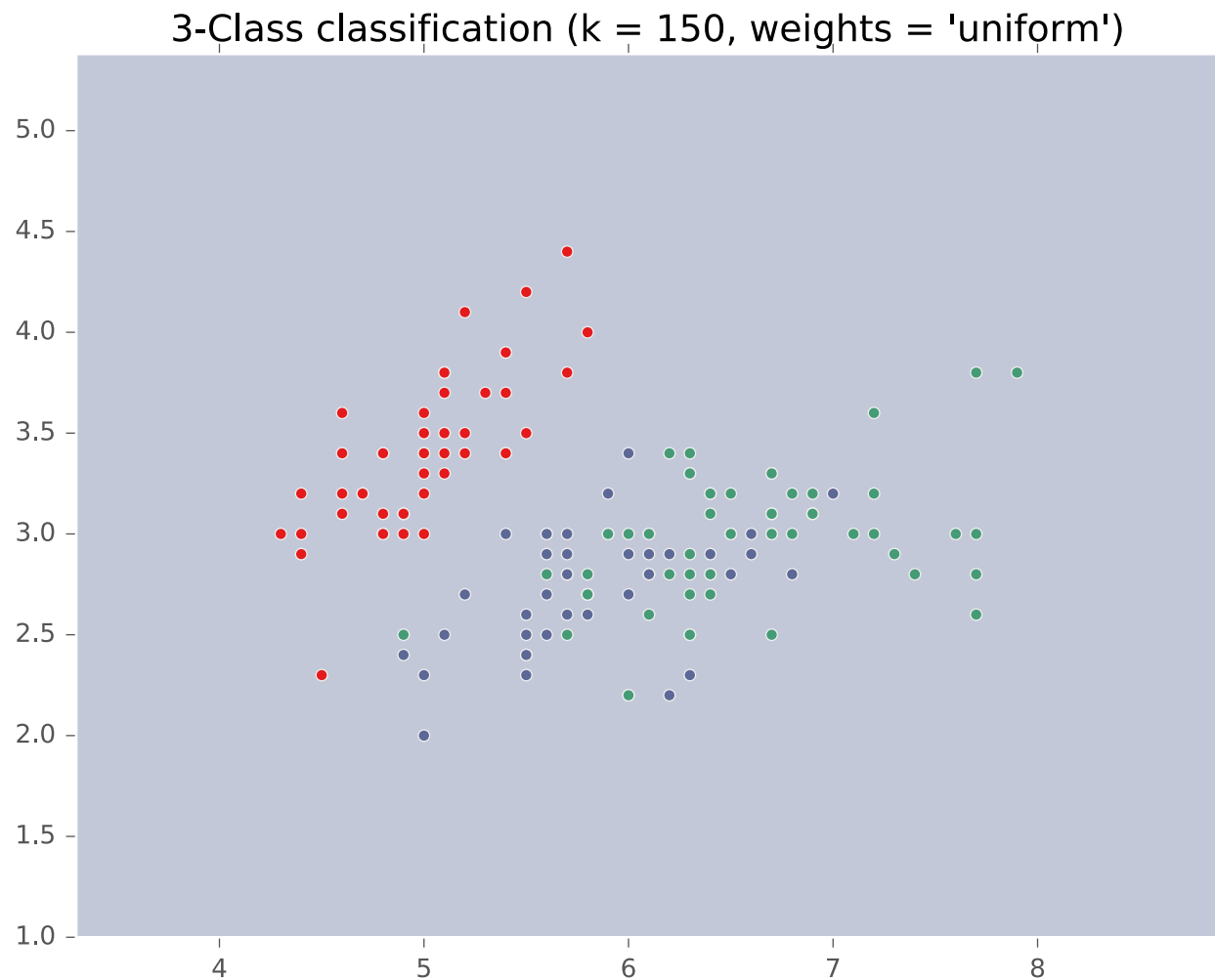


KNN on Fisher Iris Data



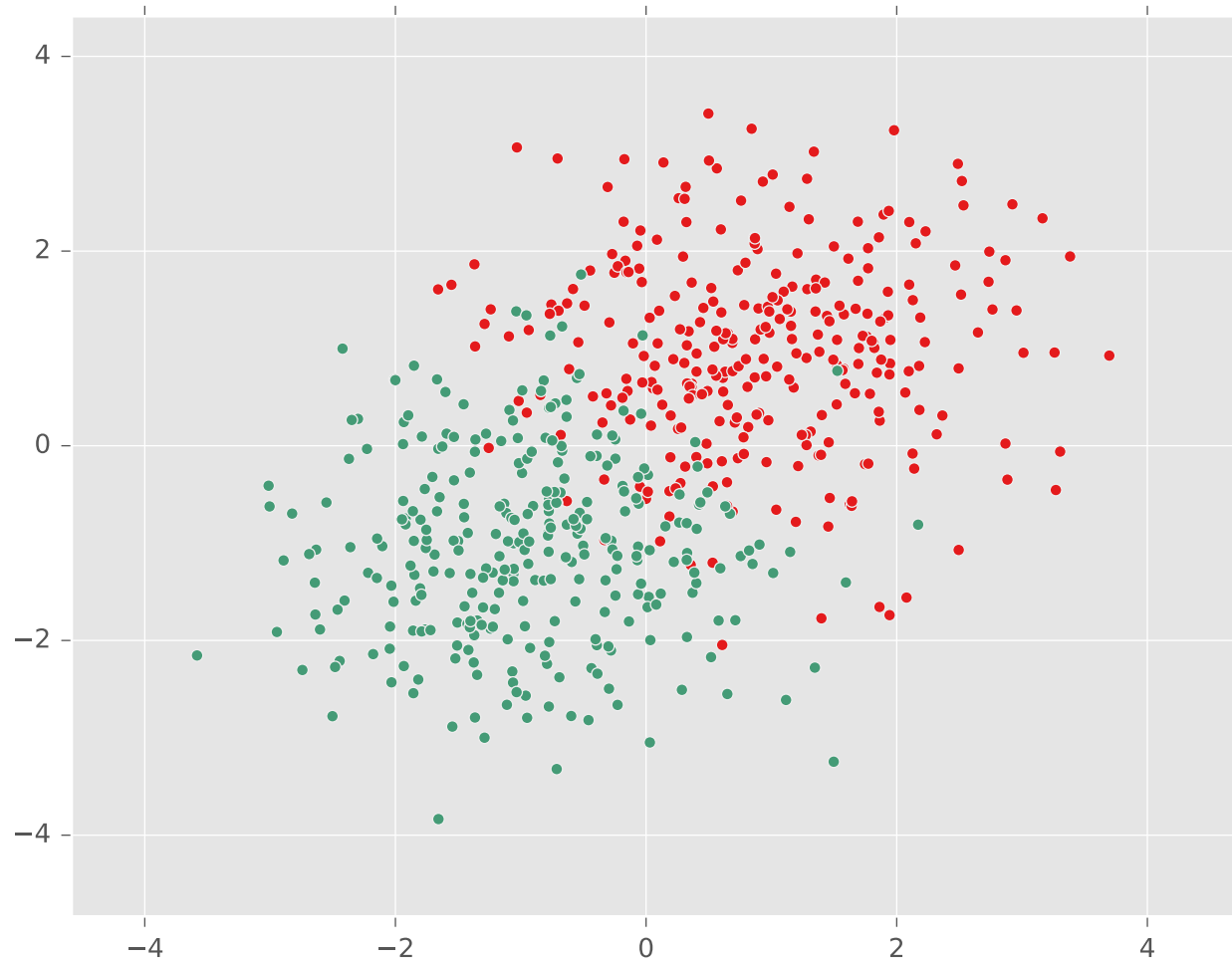
KNN on Fisher Iris Data

Special Case: Majority Vote

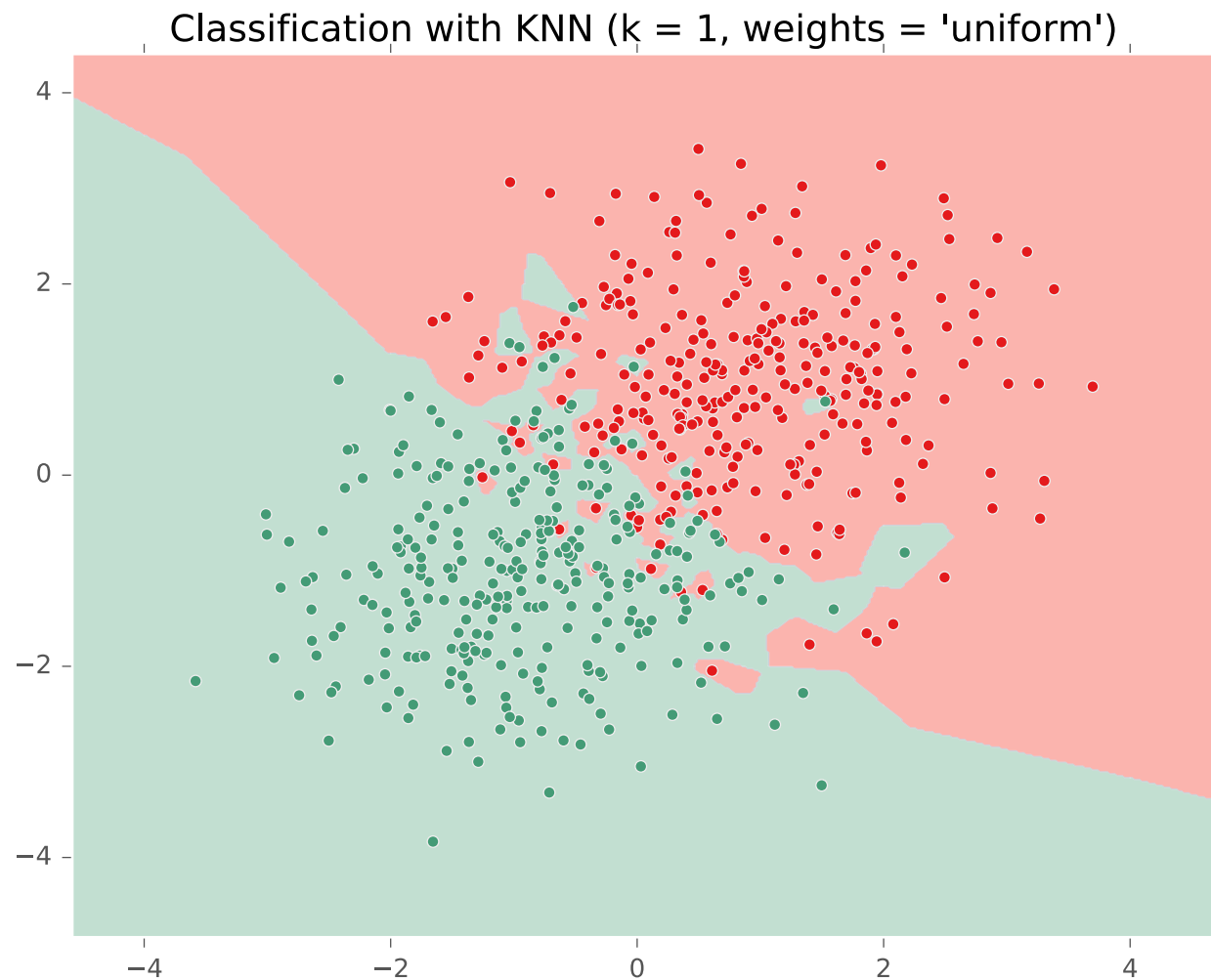


KNN ON GAUSSIAN DATA

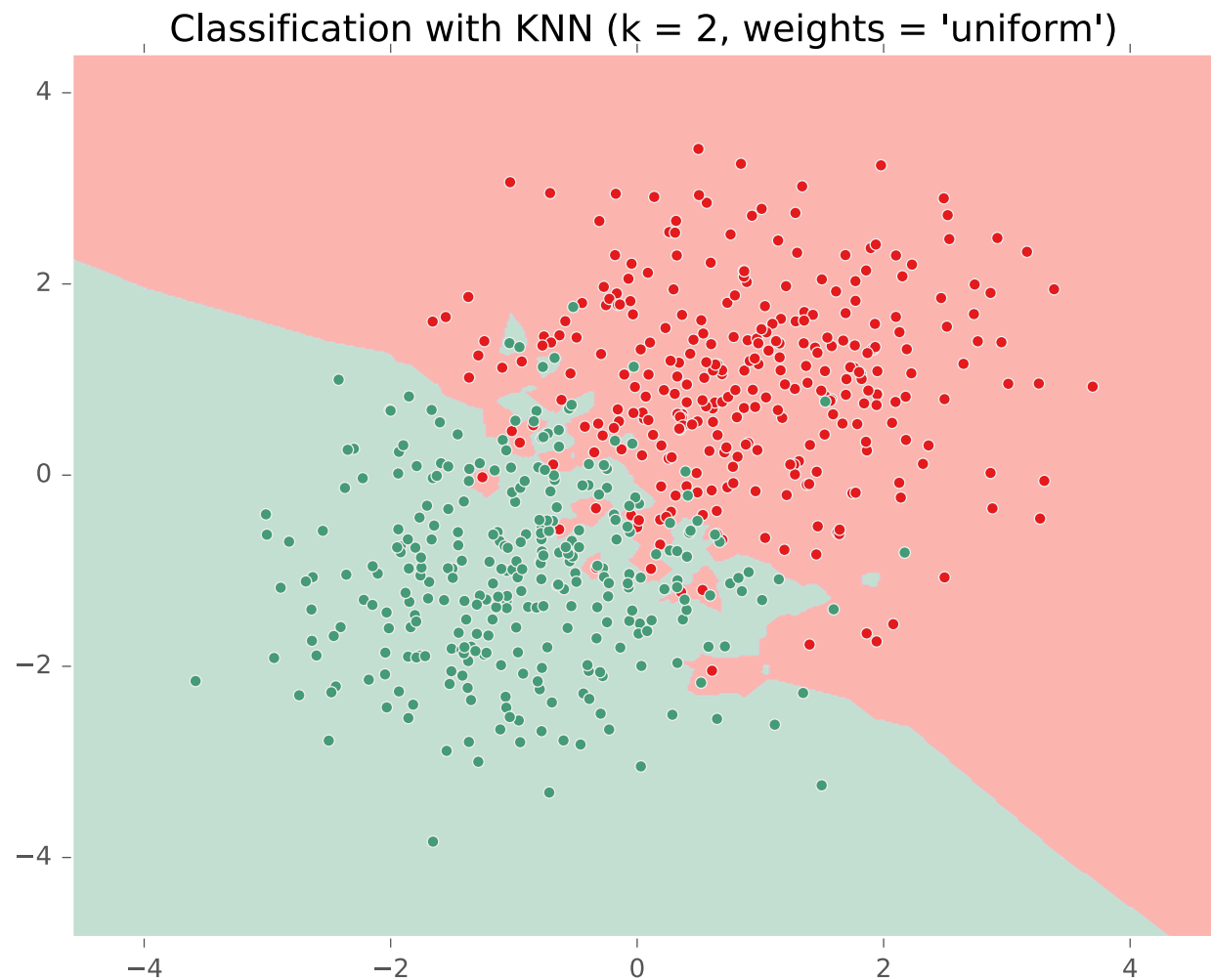
KNN on Gaussian Data



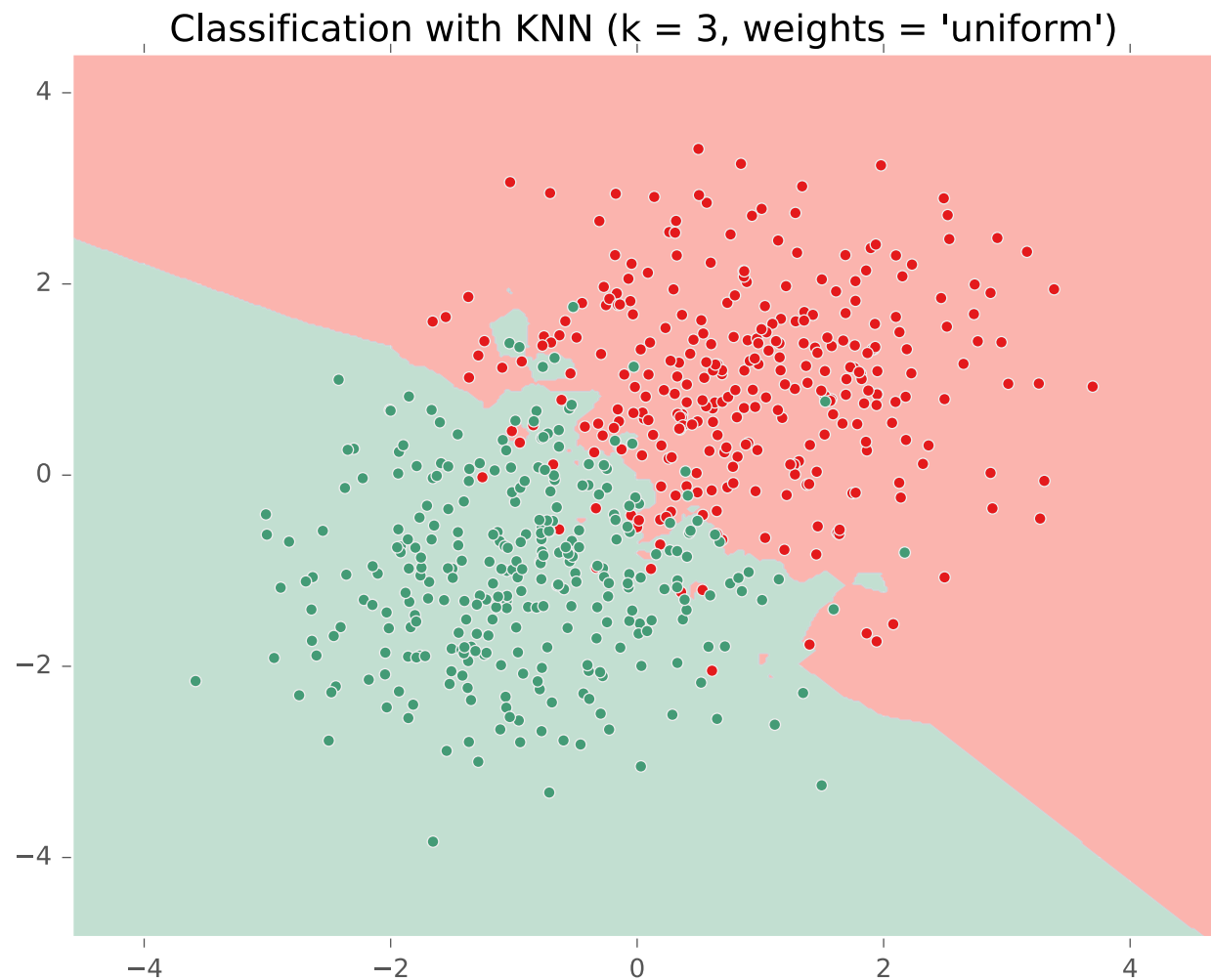
KNN on Gaussian Data



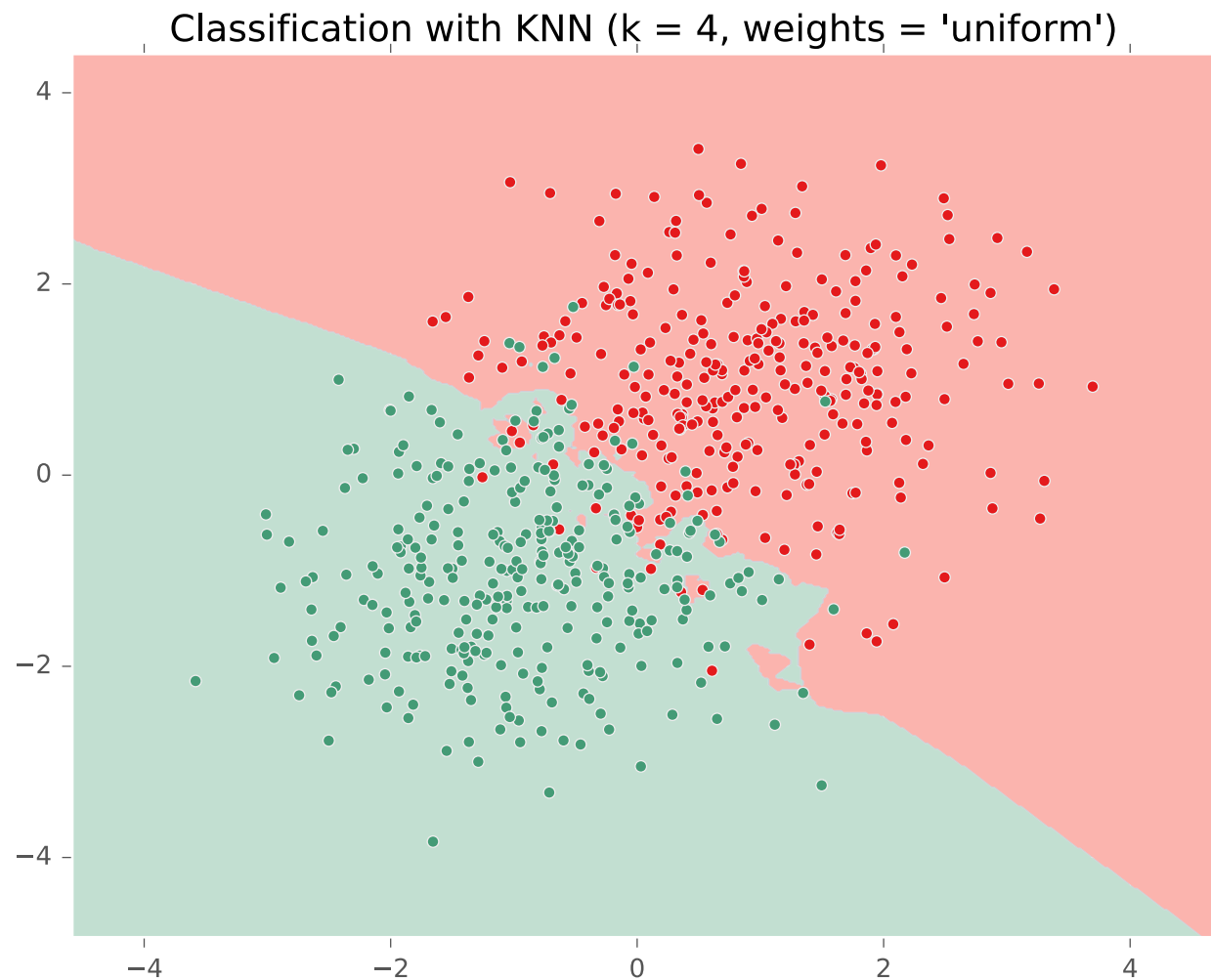
KNN on Gaussian Data



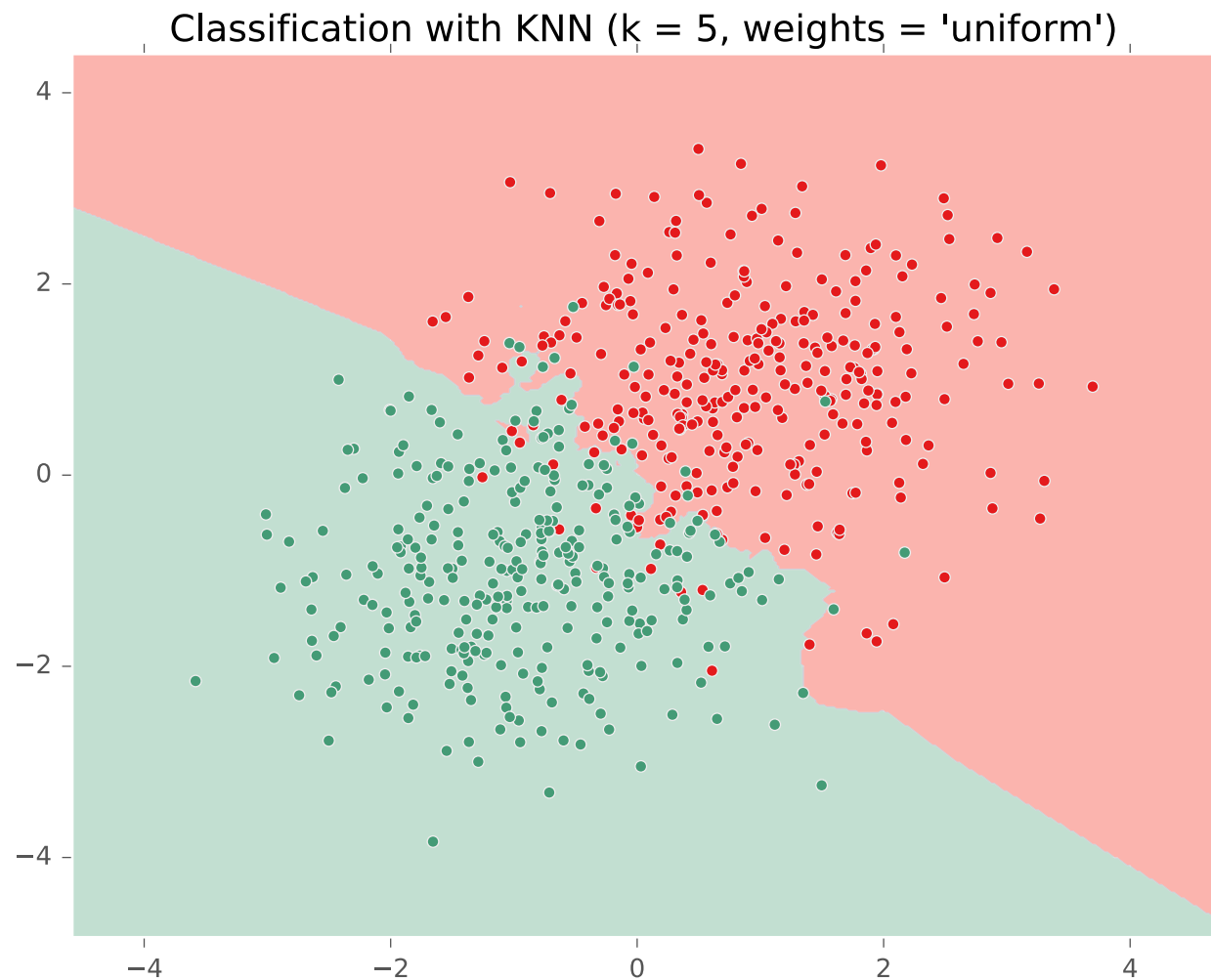
KNN on Gaussian Data



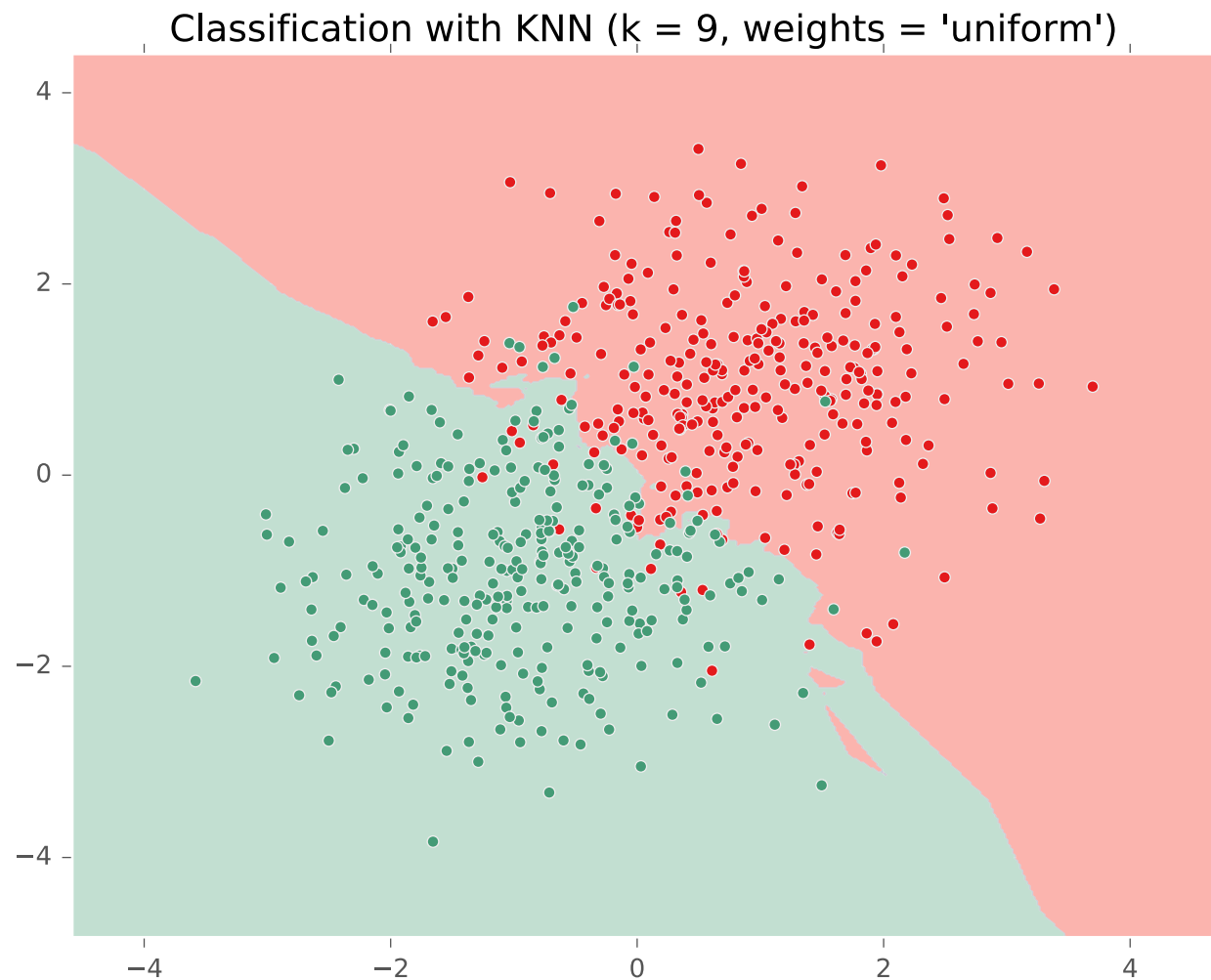
KNN on Gaussian Data



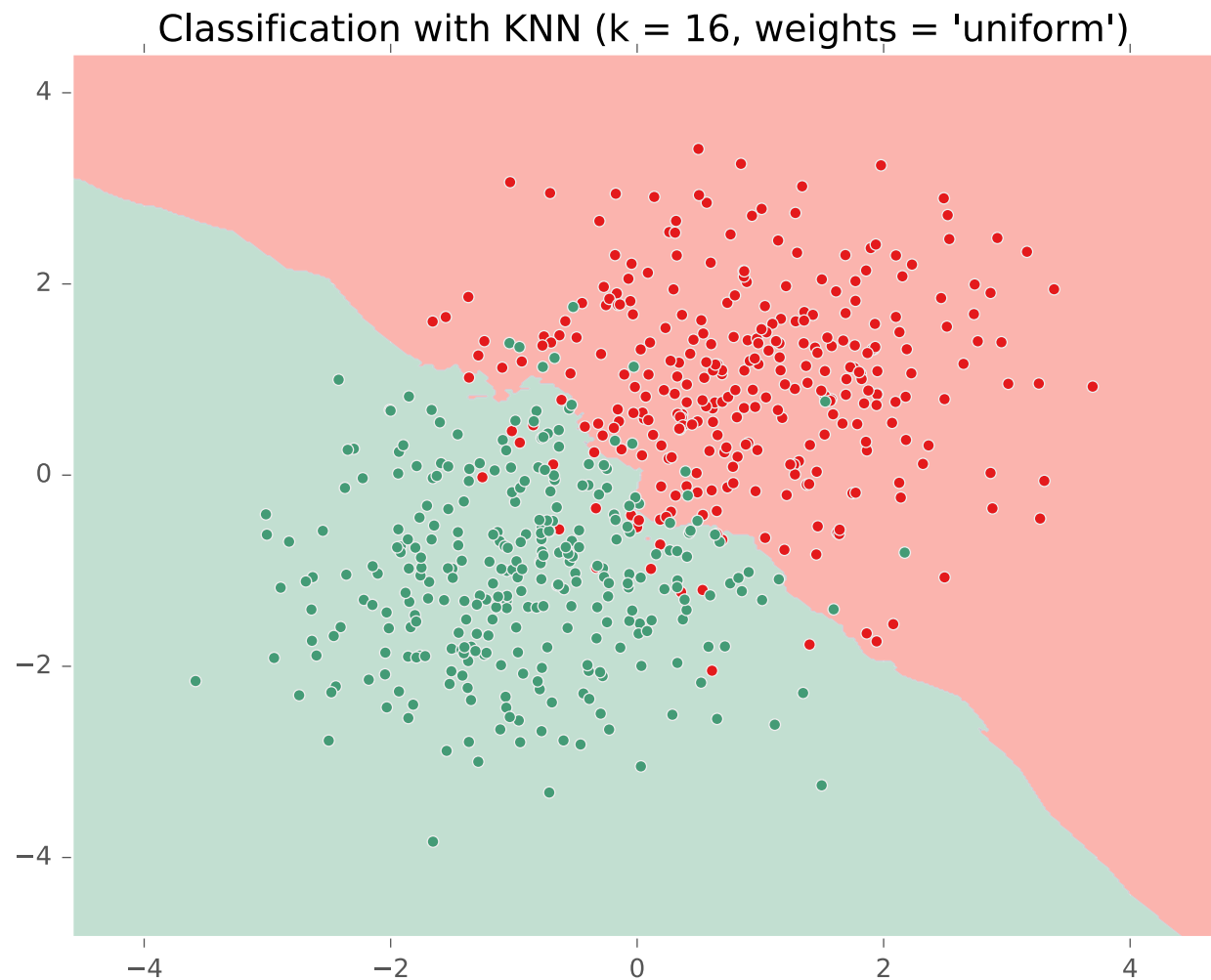
KNN on Gaussian Data



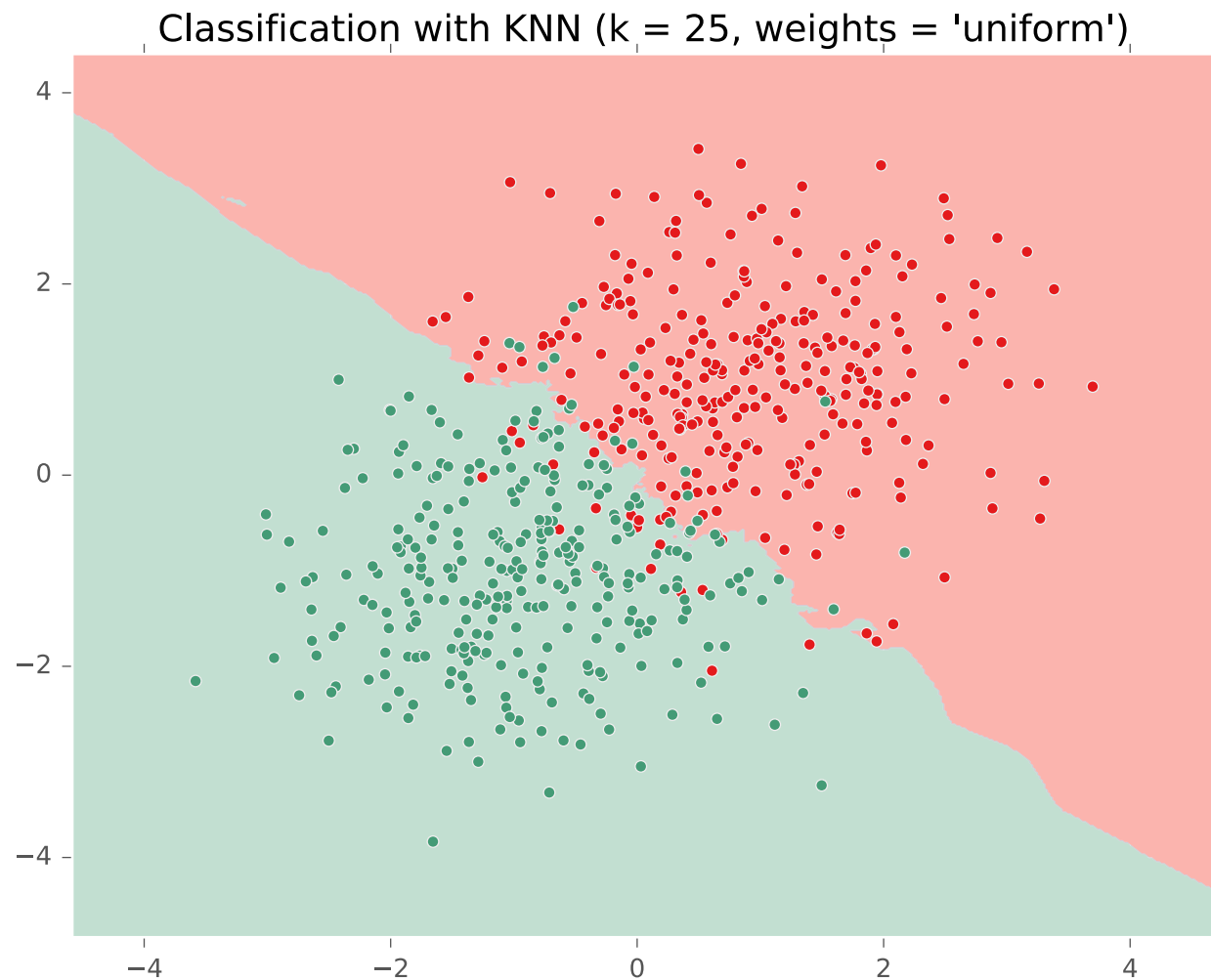
KNN on Gaussian Data



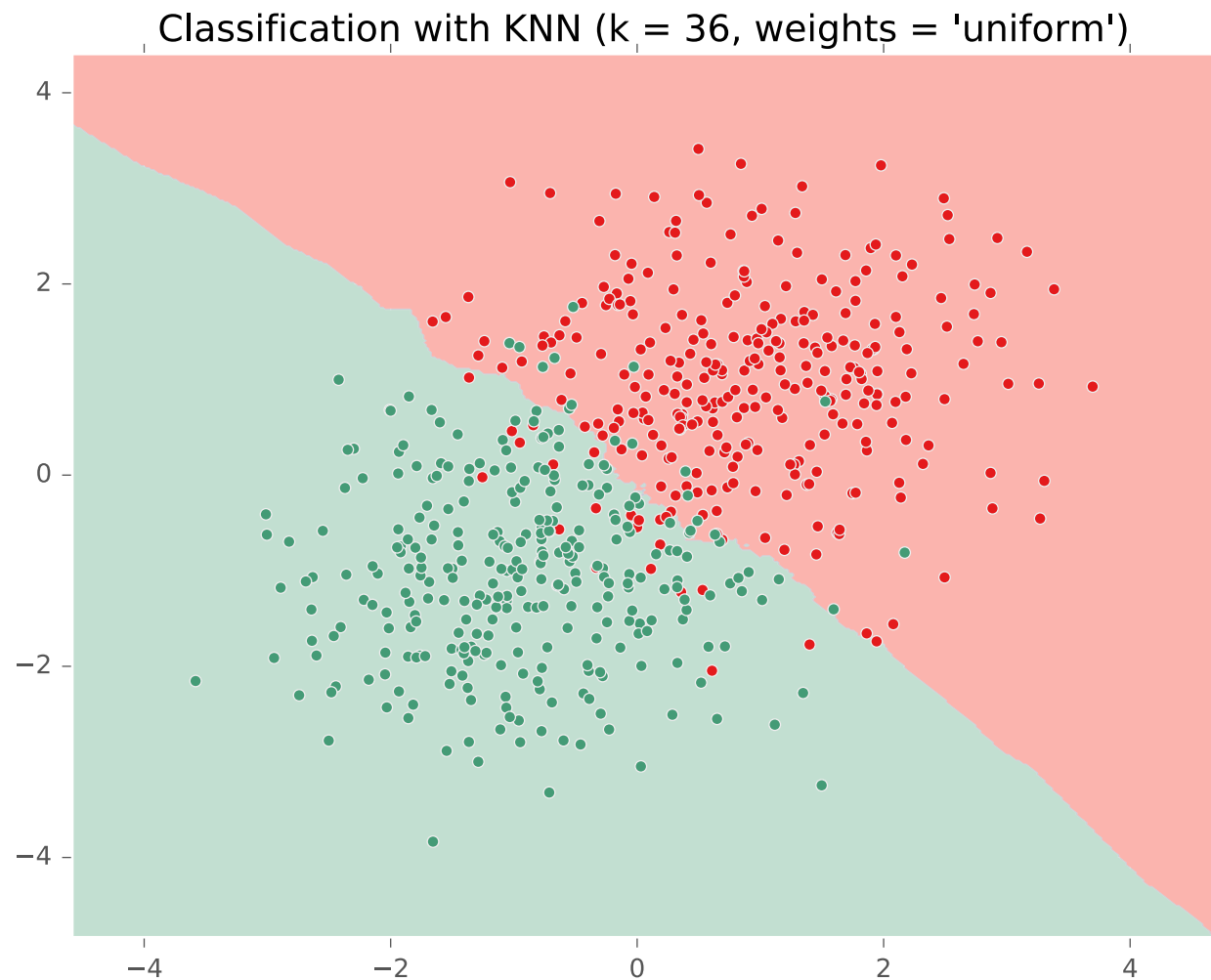
KNN on Gaussian Data



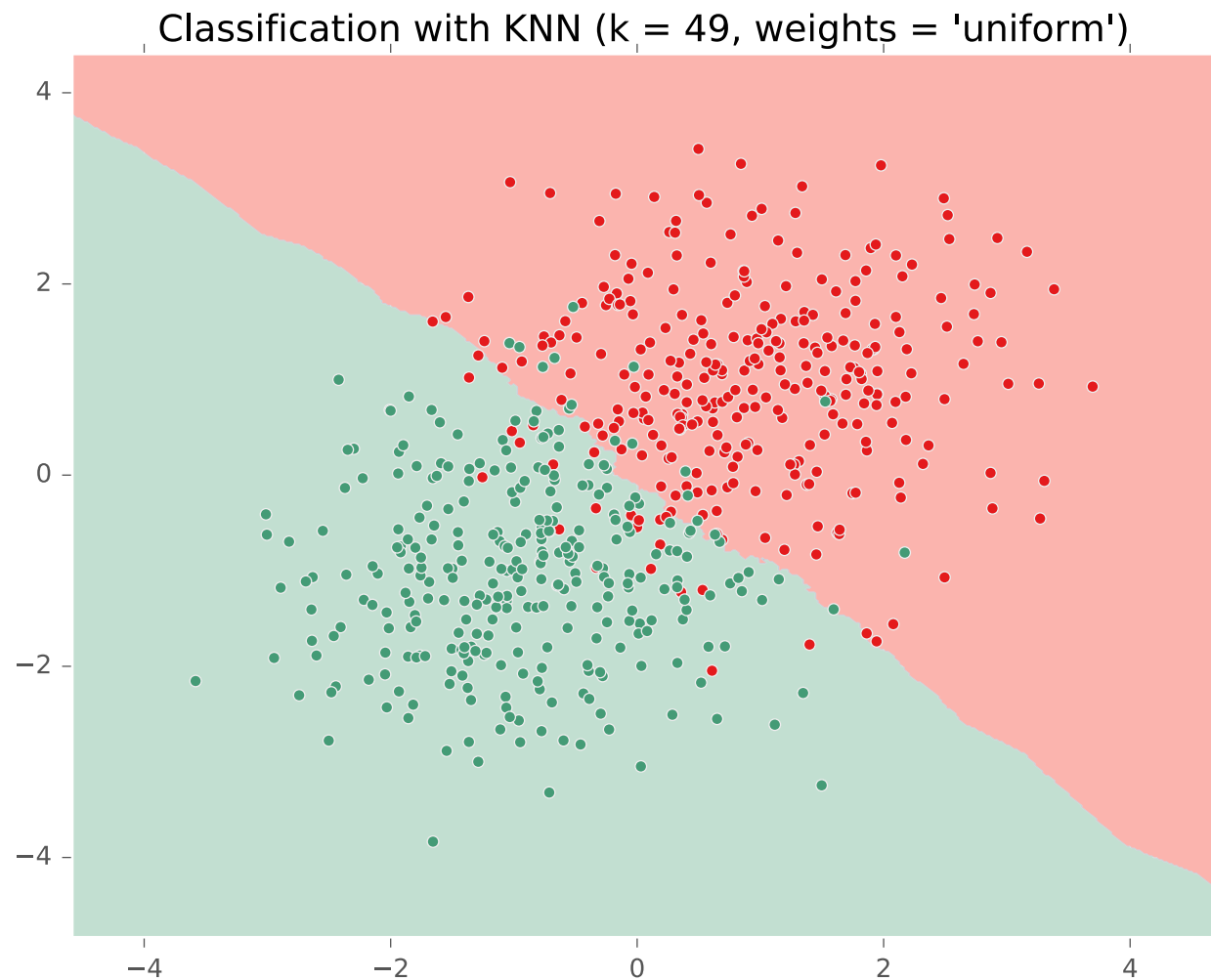
KNN on Gaussian Data



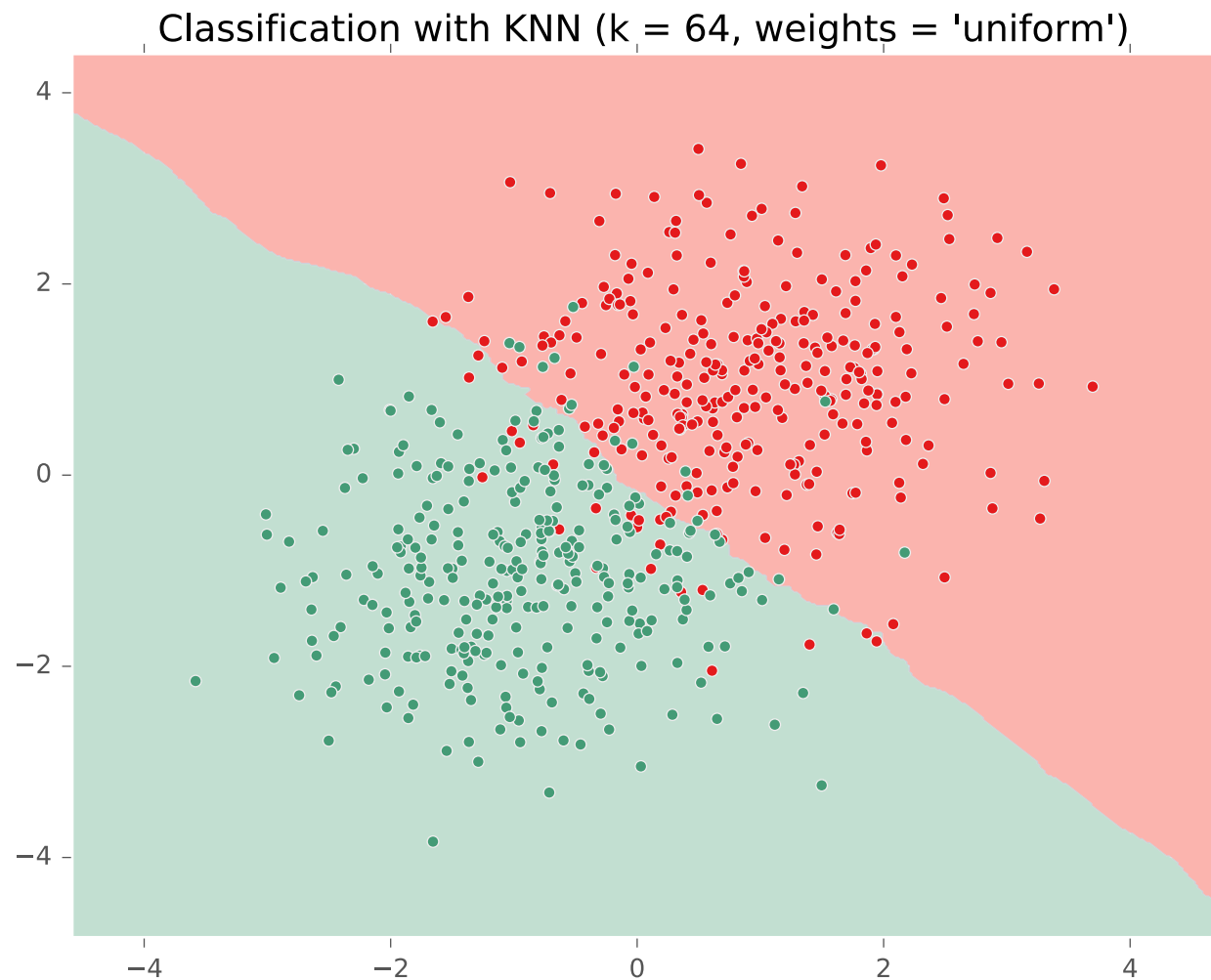
KNN on Gaussian Data



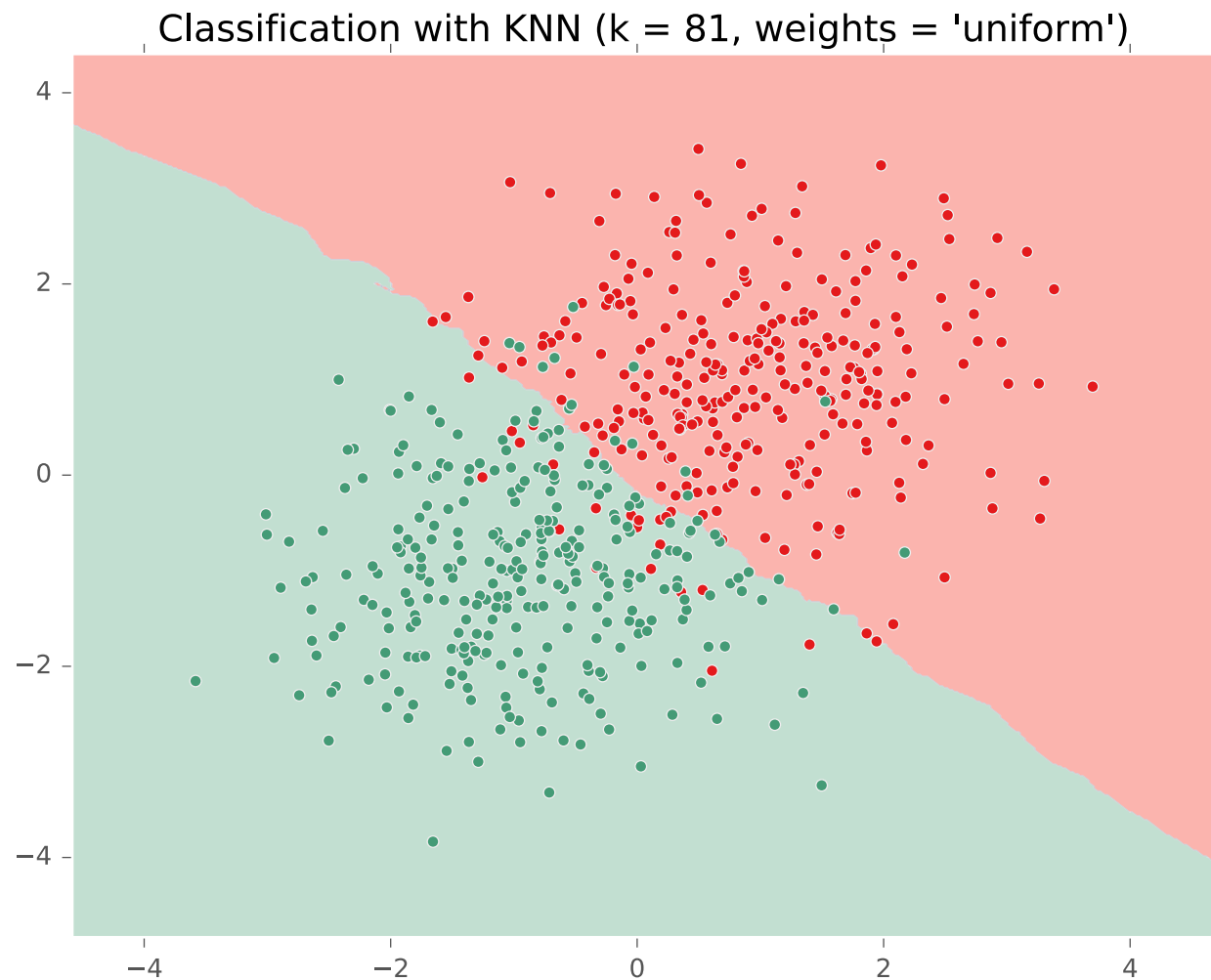
KNN on Gaussian Data



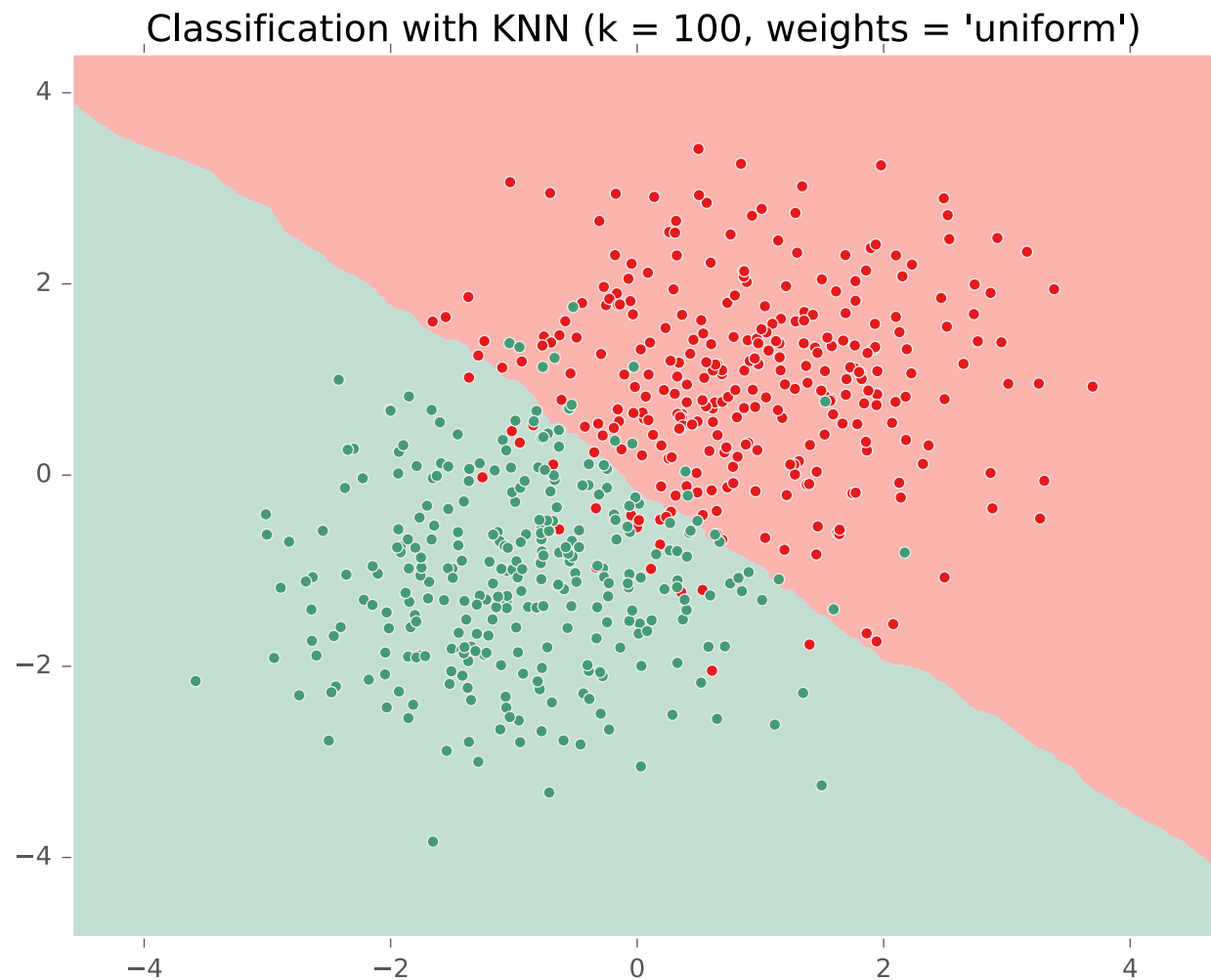
KNN on Gaussian Data



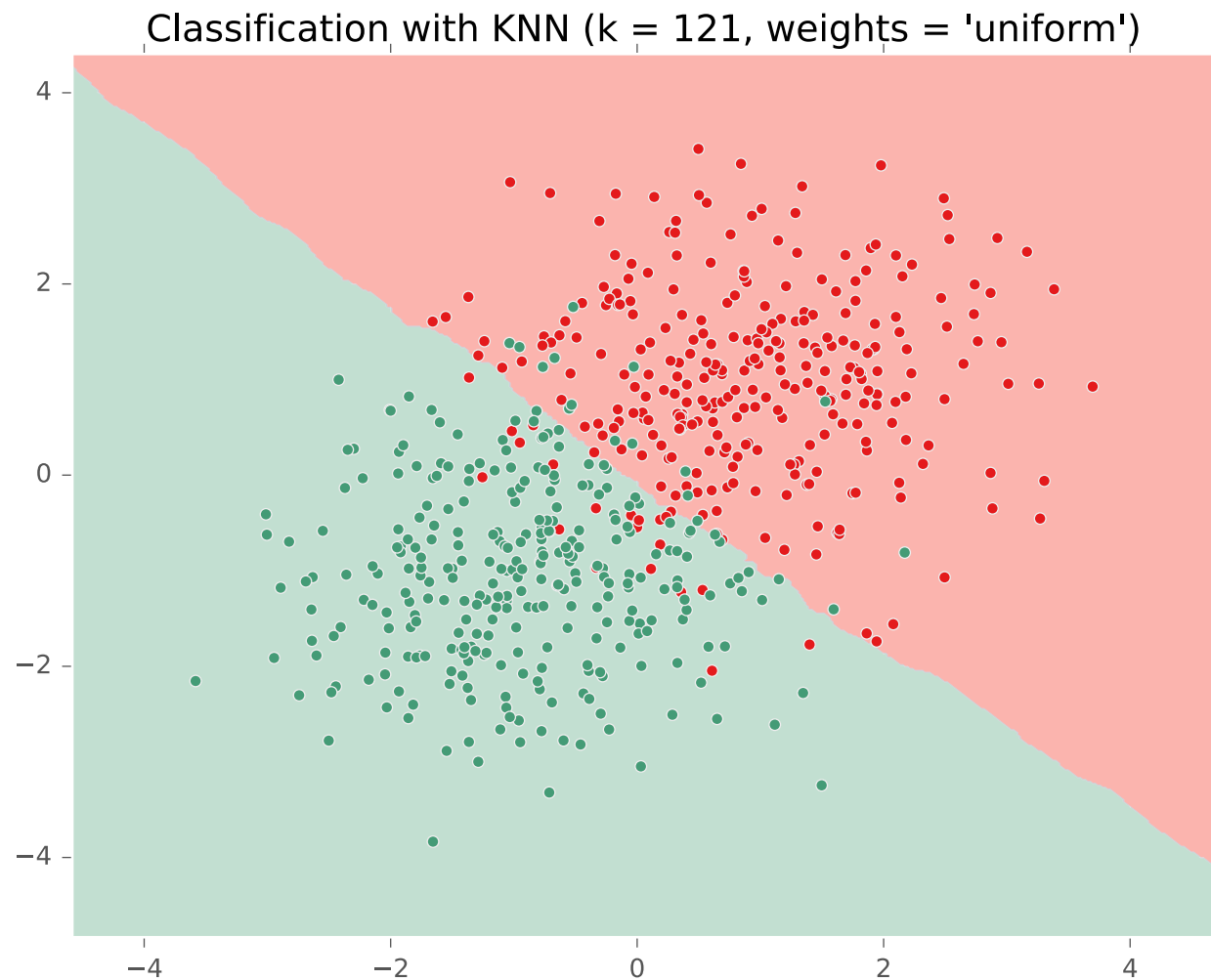
KNN on Gaussian Data



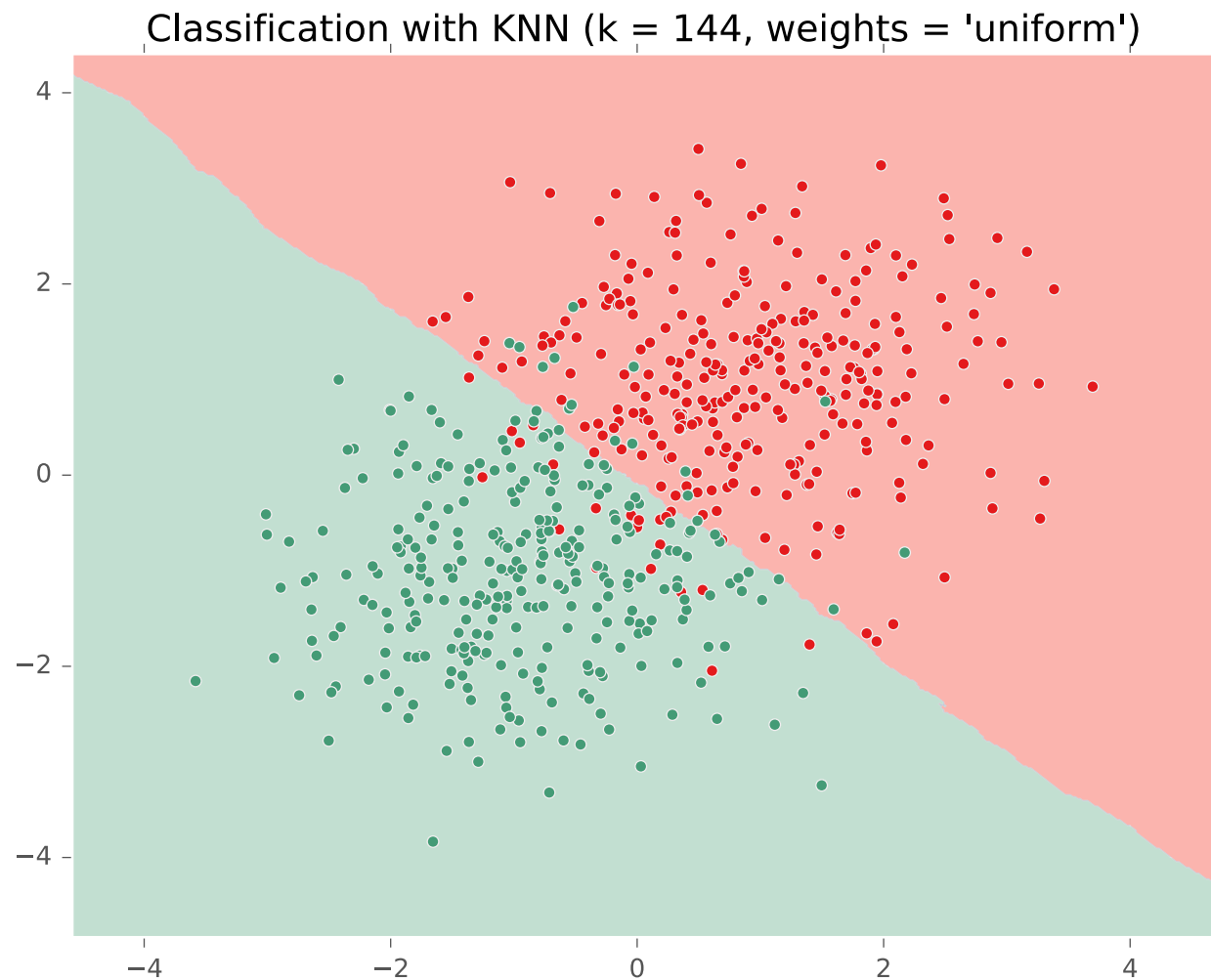
KNN on Gaussian Data



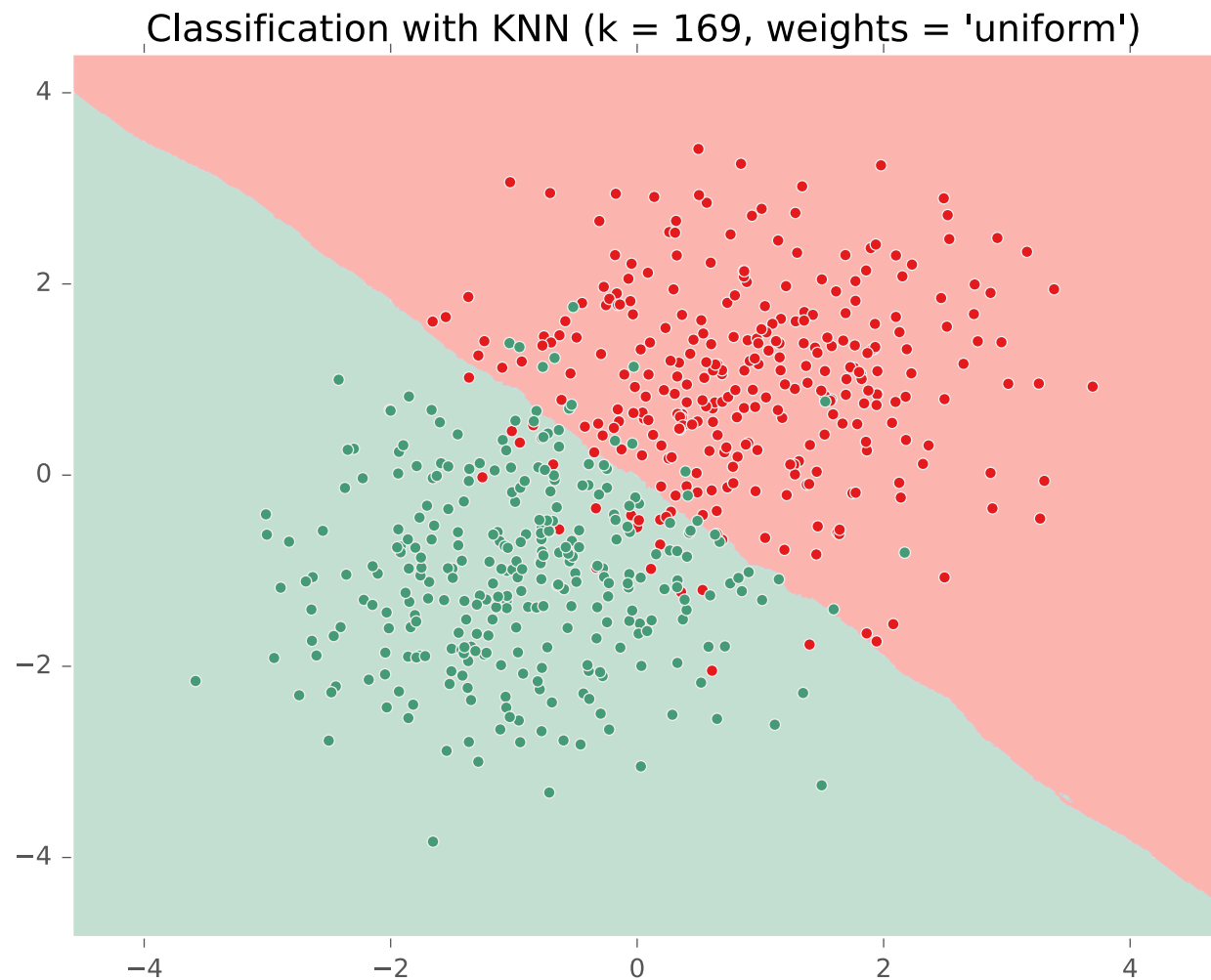
KNN on Gaussian Data



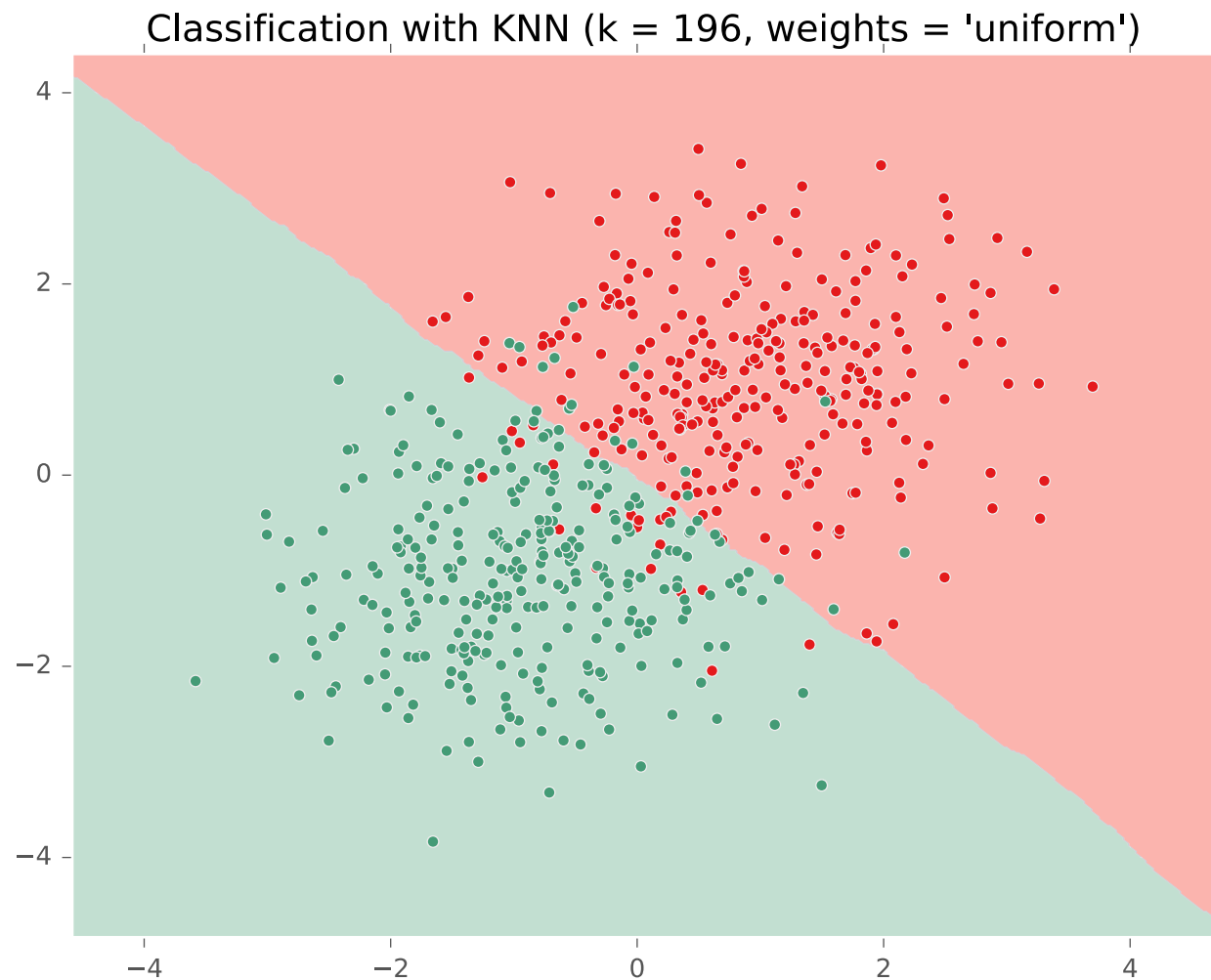
KNN on Gaussian Data



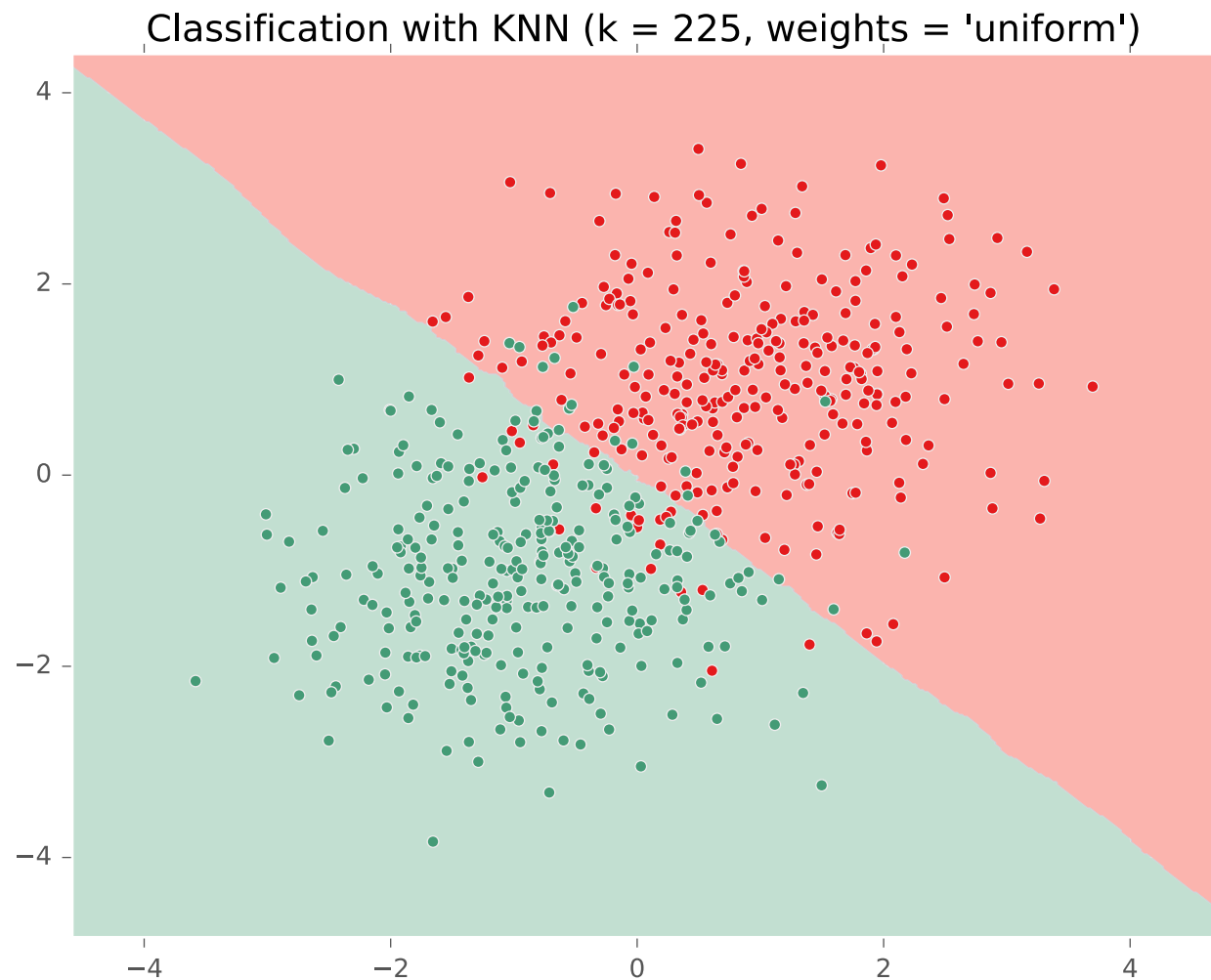
KNN on Gaussian Data



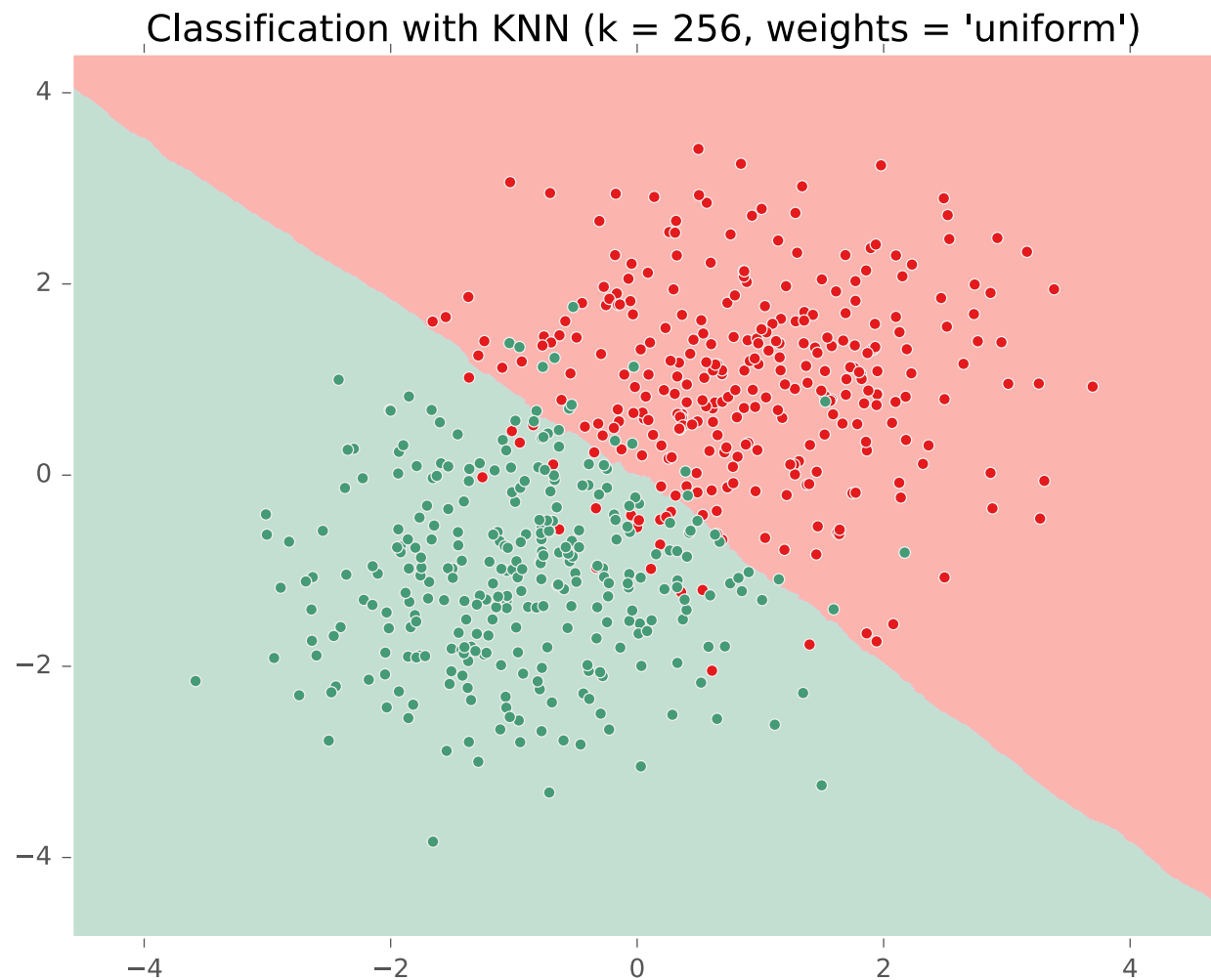
KNN on Gaussian Data



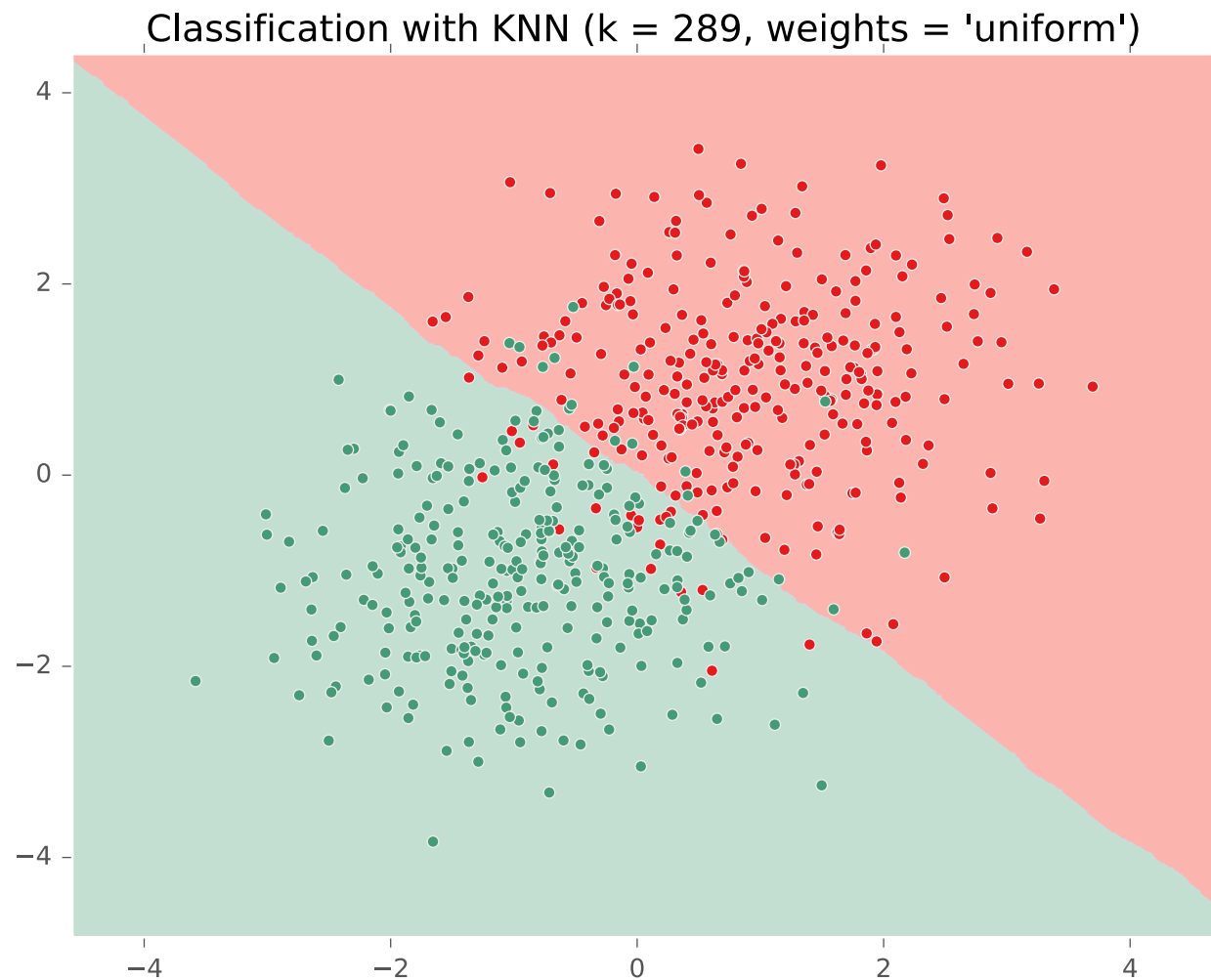
KNN on Gaussian Data



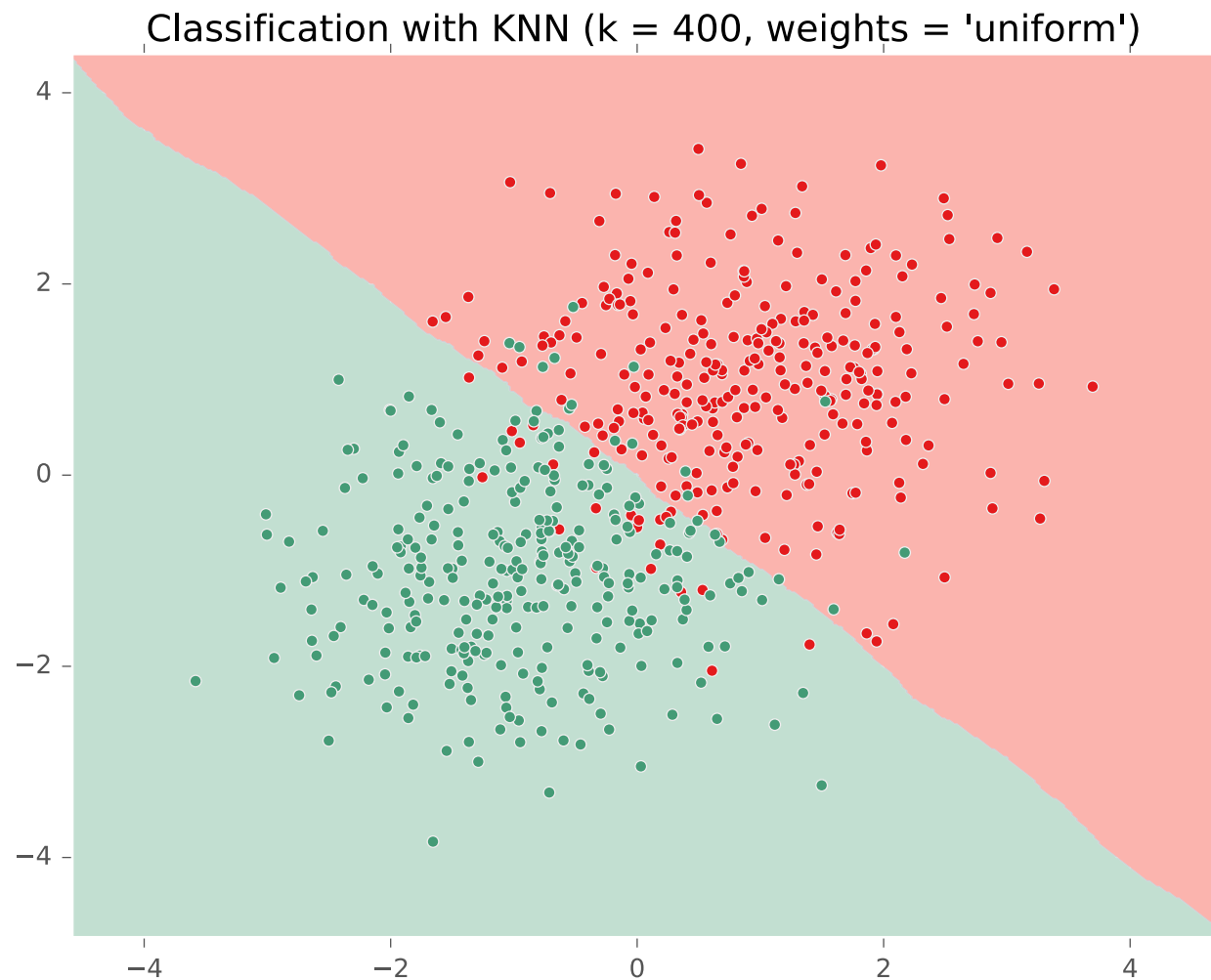
KNN on Gaussian Data



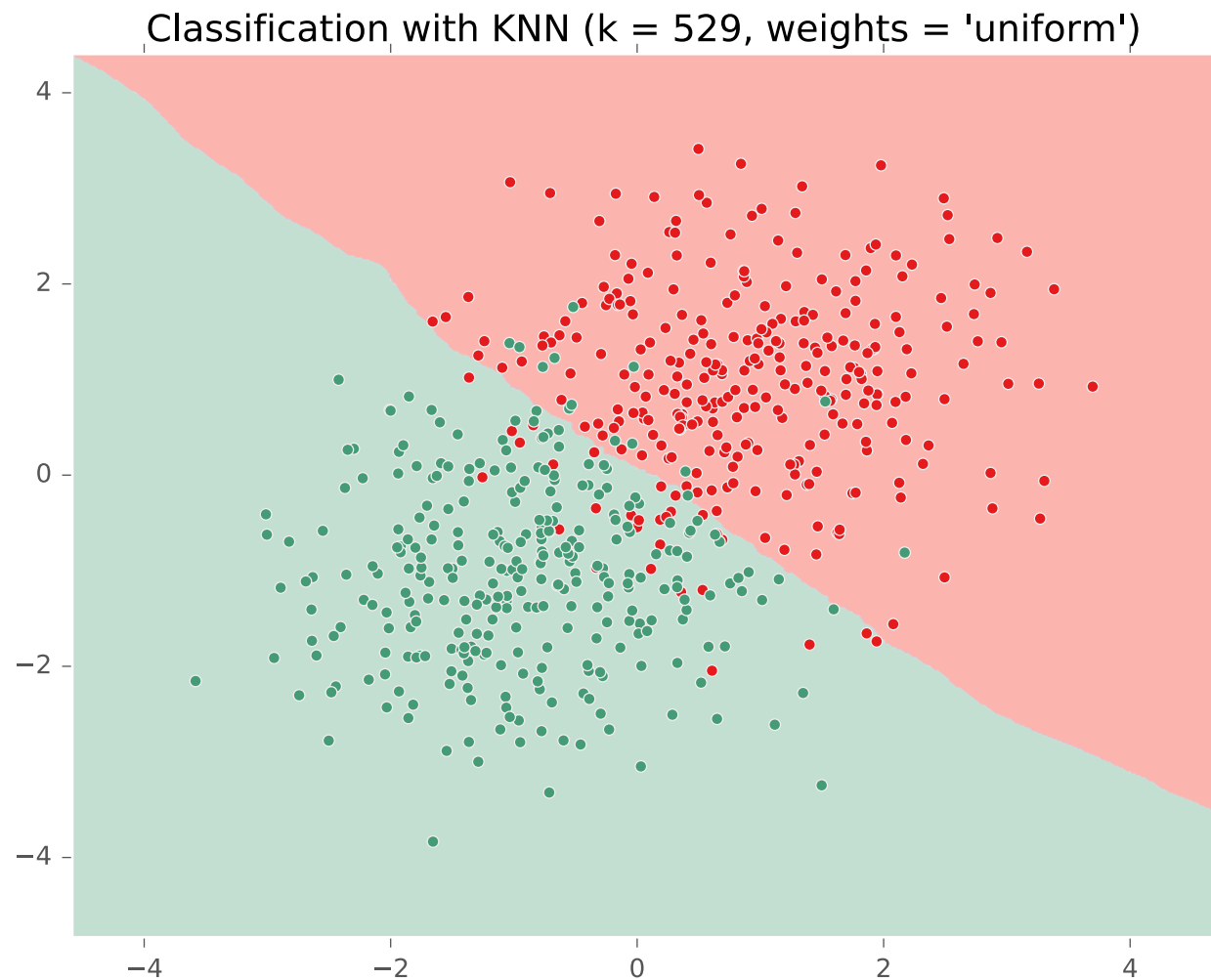
KNN on Gaussian Data



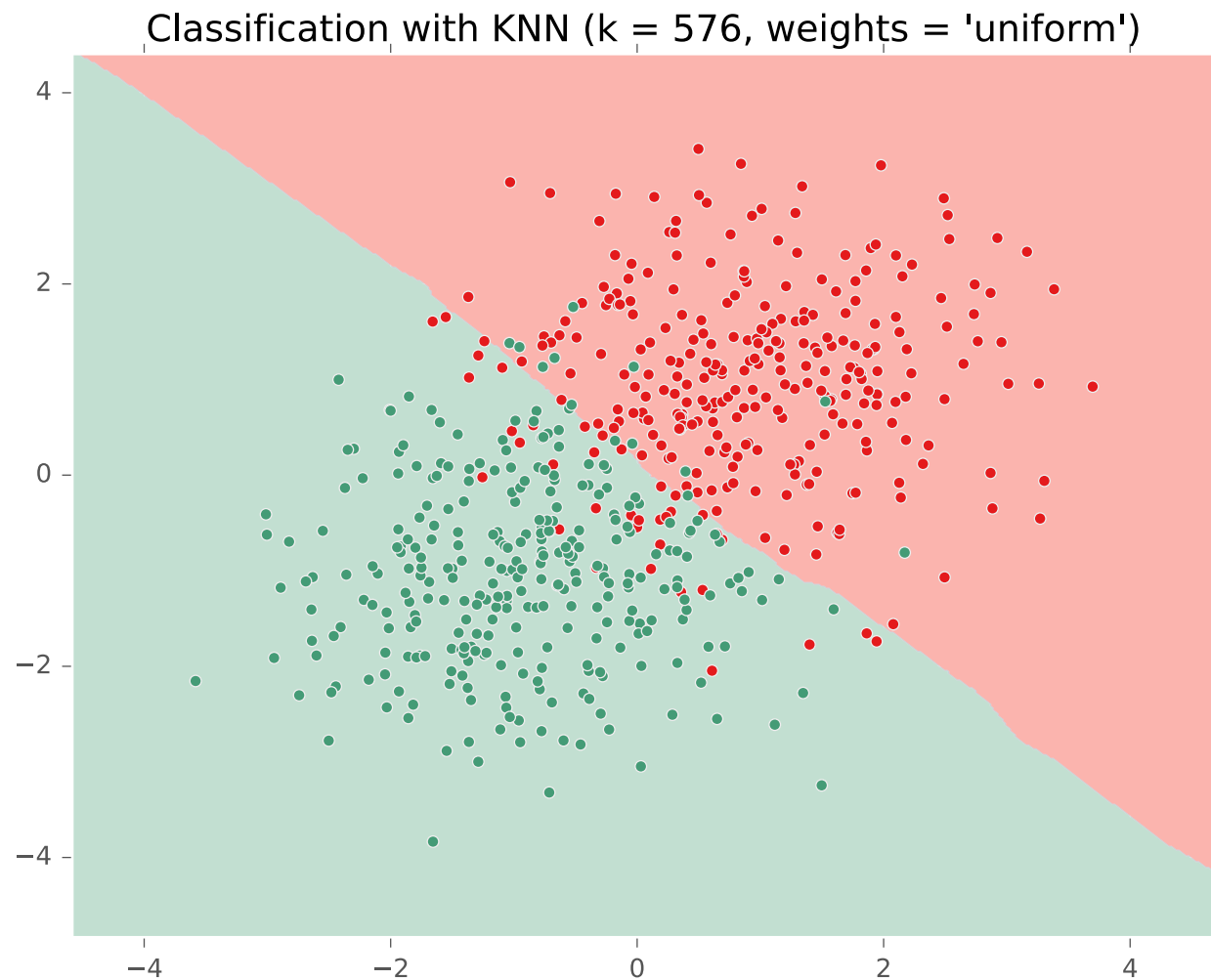
KNN on Gaussian Data



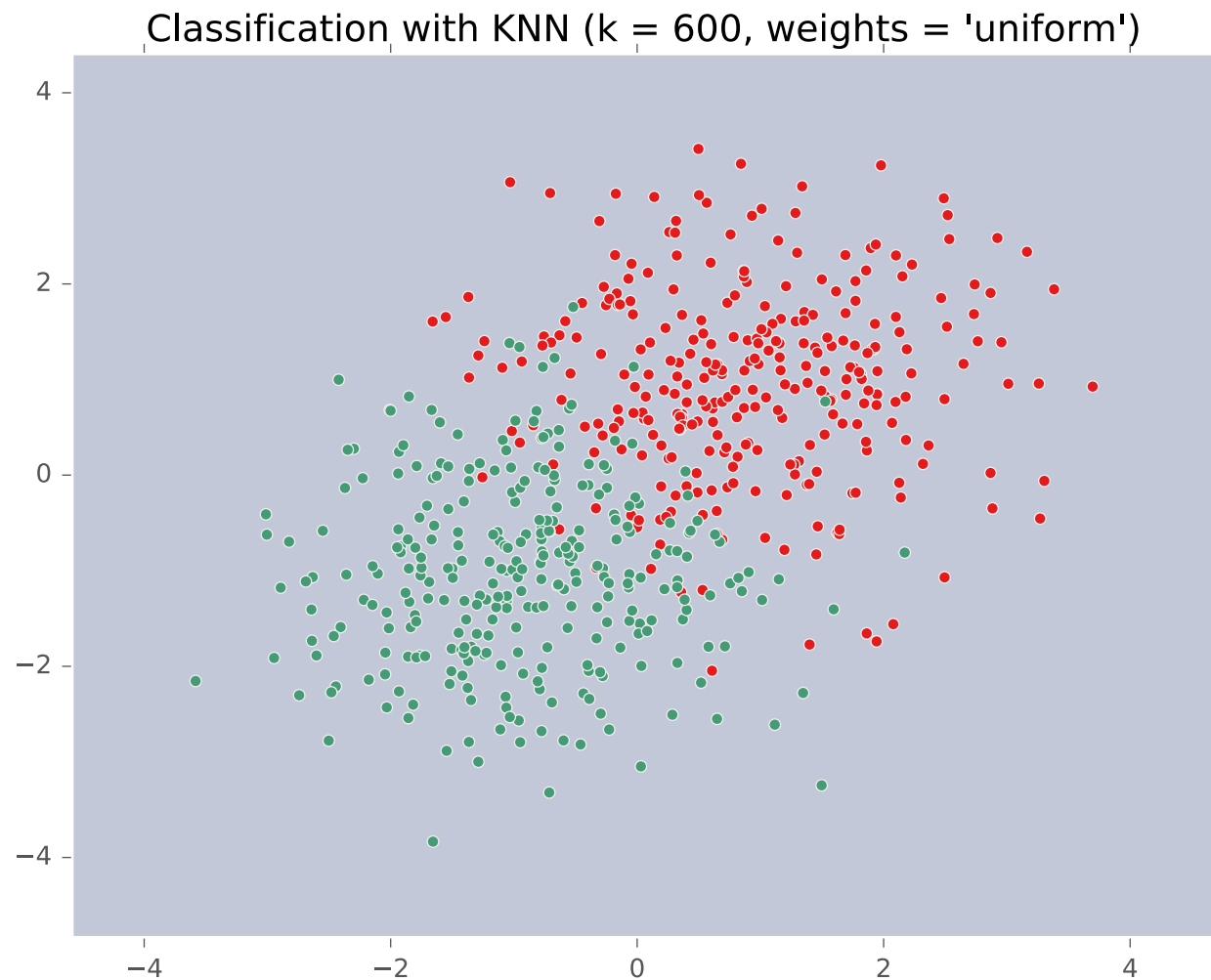
KNN on Gaussian Data



KNN on Gaussian Data



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

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

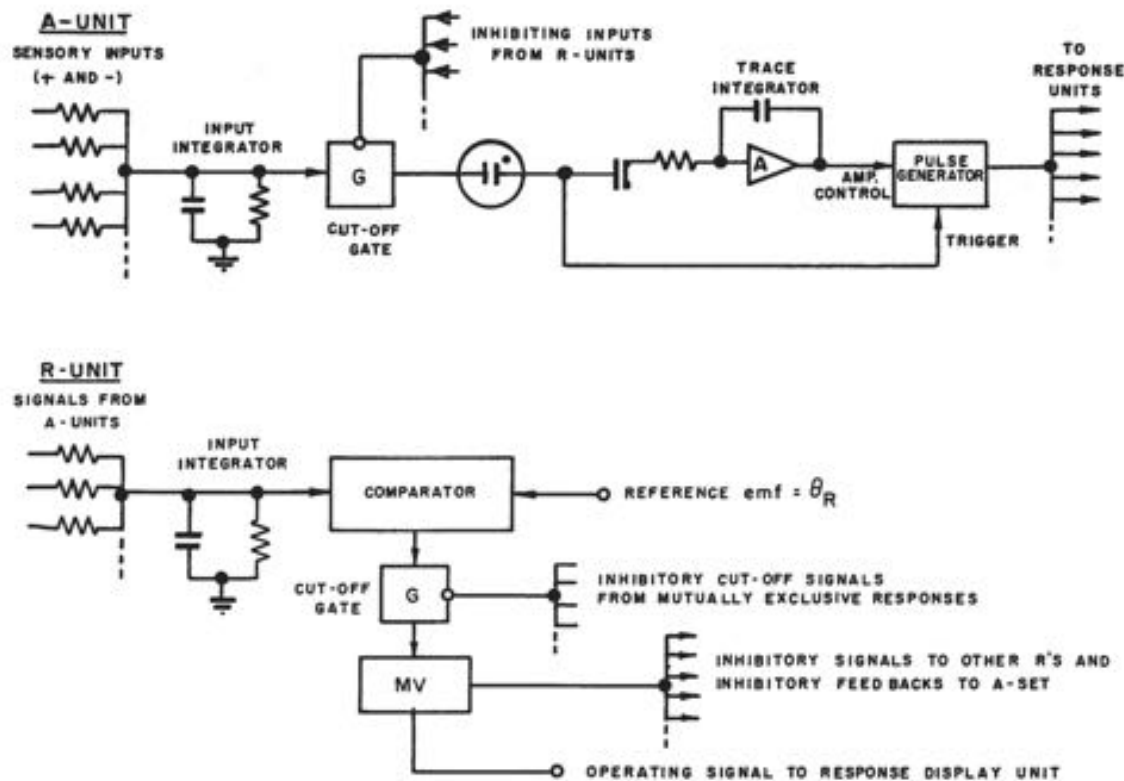


FIGURE 5
DESIGN OF TYPICAL UNITS

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:

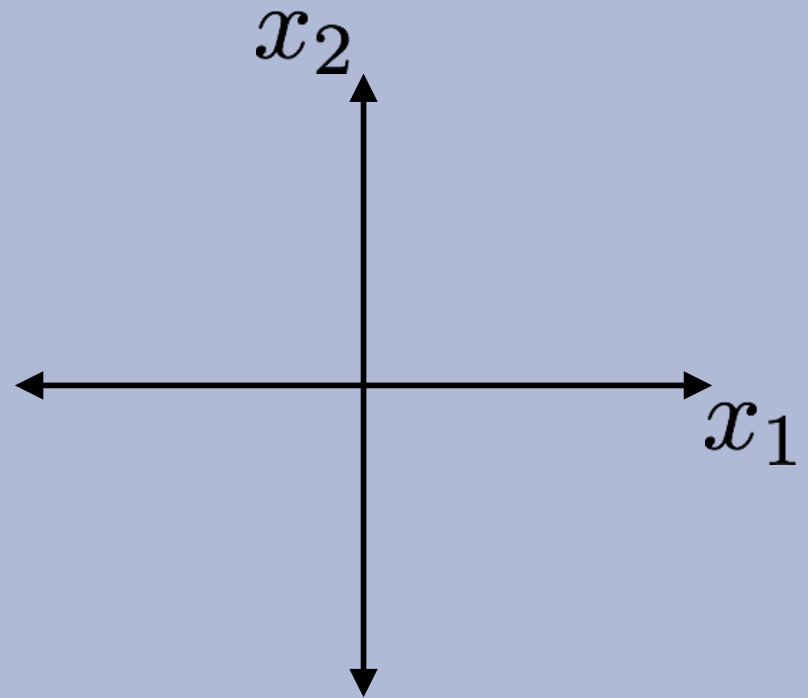
$$w_1x_1 + w_2x_2 + b > 0$$

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

Draw the vector

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

Answer Here:



Visualizing Dot-Products

Chalkboard:

- vector in 2D
- line in 2D
- adding a bias term
- definition of orthogonality
- vector projection
- hyperplane definition
- half-space definitions

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:

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

Online vs. Batch Learning

Batch Learning

Learn from all the examples at once

Online Learning

Gradually learn as each example is received

Online Learning

Examples

1. **Stock market** prediction (what will the value of Alphabet Inc. be tomorrow?)
2. **Email** classification (distribution of both spam and regular mail changes over time, but the target function stays fixed - last year's spam still looks like spam)
3. **Recommendation** systems. Examples: recommending movies; predicting whether a user will be interested in a new news article
4. **Ad placement** in a new market

Online Learning

For $i = 1, 2, 3, \dots$:

- **Receive** an unlabeled instance $\mathbf{x}^{(i)}$
- **Predict** $y' = h_{\theta}(\mathbf{x}^{(i)})$
- **Receive** true label $y^{(i)}$
- **Suffer loss** if a mistake was made, $y' \neq y^{(i)}$
- **Update** parameters θ

Goal:

- **Minimize** the number of **mistakes**

Perceptron

Chalkboard:

- (Online) Perceptron Algorithm

Perceptron Algorithm: Example

Example: $(-1,2) -$ ✗

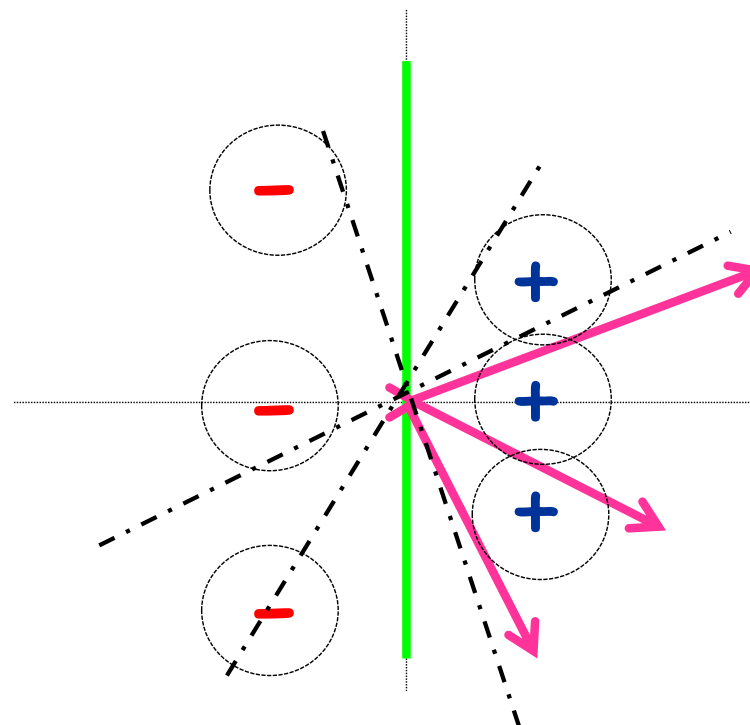
$(1,0) +$ ✓

$(1,1) +$ ✗

$(-1,0) -$ ✓

$(-1,-2) -$ ✗

$(1,-1) +$ ✓



Perceptron Algorithm: (without the bias term)

- Set $t=1$, start with all-zeroes weight vector w_1 .
- Given example x , predict positive iff $w_t \cdot x \geq 0$.
- On a mistake, update as follows:
 - Mistake on positive, update $w_{t+1} \leftarrow w_t + x$
 - Mistake on negative, update $w_{t+1} \leftarrow w_t - x$

$$w_1 = (0,0)$$

$$w_2 = w_1 - (-1,2) = (1,-2)$$

$$w_3 = w_2 + (1,1) = (2,-1)$$

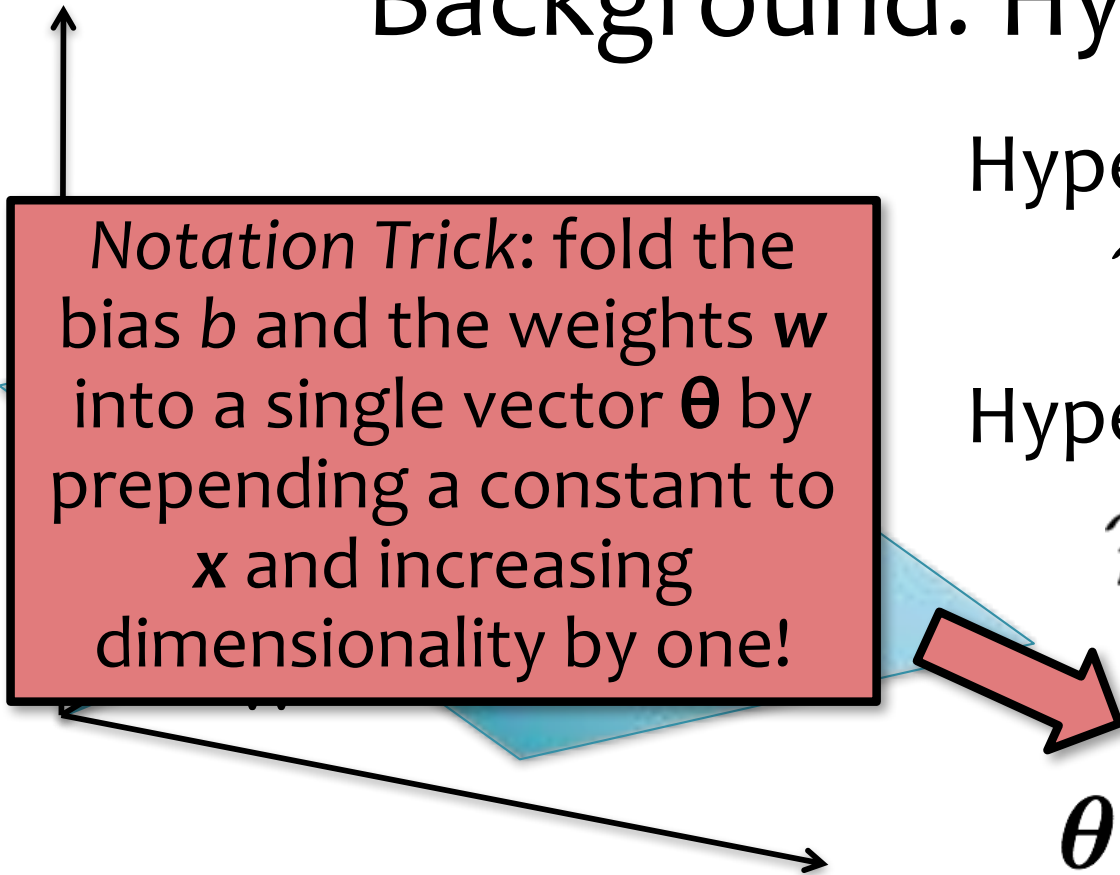
$$w_4 = w_3 - (-1,-2) = (3,1)$$

Perceptron

Chalkboard:

- Why do we need a bias term?
- (Batch) Perceptron Algorithm
- Inductive Bias of Perceptron
- Limitations of Linear Models

Background: Hyperplanes



Notation Trick: fold the bias b and the weights \mathbf{w} into a single vector $\boldsymbol{\theta}$ by prepending a constant to \mathbf{x} and increasing dimensionality by one!

Hyperplane (Definition 1):

$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} = b\}$$

Hyperplane (Definition 2):

$$\mathcal{H} = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} = 0$$

$$\text{and } x_0 = 1\}$$

$$\boldsymbol{\theta} = [b, w_1, \dots, w_M]^T$$

Half-spaces:

$$\mathcal{H}^+ = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} > 0 \text{ and } x_0 = 1\}$$

$$\mathcal{H}^- = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} < 0 \text{ and } x_0 = 1\}$$

(Online) Perceptron Algorithm

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$$

$$\text{where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{+1, -1\}$$

Prediction: Output determined by hyperplane.

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \text{sign}(\boldsymbol{\theta}^T \mathbf{x})$$

$$\text{sign}(a) = \begin{cases} 1, & \text{if } a \geq 0 \\ -1, & \text{otherwise} \end{cases}$$

$$\text{Assume } \boldsymbol{\theta} = [b, w_1, \dots, w_M]^T \text{ and } x_0 = 1$$

Learning: Iterative procedure:

- **initialize** parameters to vector of all zeroes
- **while** not converged
 - **receive** next example $(\mathbf{x}^{(i)}, y^{(i)})$
 - **predict** $y' = h(\mathbf{x}^{(i)})$
 - **if** positive mistake: **add** $\mathbf{x}^{(i)}$ to parameters
 - **if** negative mistake: **subtract** $\mathbf{x}^{(i)}$ from parameters

(Online) Perceptron Algorithm

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots$$

where $\mathbf{x} \in \mathbb{R}^M$ and $y \in \{+1, -1\}$

Prediction: Output determined by

$$\hat{y} = h_{\theta}(\mathbf{x}) = \text{sign}(\theta^T \mathbf{x})$$

Assume $\theta = [b, w_1, \dots, w_M]$

Learning:

Algorithm 1 Perceptron Learning Algorithm

```
1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}$ )
2:    $\theta \leftarrow \mathbf{0}$ 
3:   for  $i \in \{1, 2, \dots\}$  do
4:      $\hat{y} \leftarrow \text{sign}(\theta^T \mathbf{x}^{(i)})$ 
5:     if  $\hat{y} \neq y^{(i)}$  then
6:        $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$ 
7:   return  $\theta$ 
```

Implementation Trick: same behavior as our “add on positive mistake and subtract on negative mistake” version, because $y^{(i)}$ takes care of the sign

- ▷ Initialize parameters
- ▷ For each example
- ▷ Predict
- ▷ If mistake
- ▷ Update parameters

(Batch) Perceptron Algorithm

Learning for Perceptron also works if we have a fixed training dataset, D . We call this the “batch” setting in contrast to the “online” setting that we’ve discussed so far.

Algorithm 1 Perceptron Learning Algorithm (Batch)

```
1: procedure PERCEPTRON( $\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ )  
2:    $\theta \leftarrow \mathbf{0}$  ▷ Initialize parameters  
3:   while not converged do  
4:     for  $i \in \{1, 2, \dots, N\}$  do ▷ For each example  
5:        $\hat{y} \leftarrow \text{sign}(\theta^T \mathbf{x}^{(i)})$  ▷ Predict  
6:       if  $\hat{y} \neq y^{(i)}$  then ▷ If mistake  
7:          $\theta \leftarrow \theta + y^{(i)} \mathbf{x}^{(i)}$  ▷ Update parameters  
8:   return  $\theta$ 
```

(Batch) Perceptron Algorithm

Learning for Perceptron also works if we have a fixed training dataset, D . We call this the “batch” setting in contrast to the “online” setting that we’ve discussed so far.

Discussion:

The Batch Perceptron Algorithm can be derived in two ways.

1. By extending the online Perceptron algorithm to the batch setting (as mentioned above)
2. By applying **Stochastic Gradient Descent (SGD)** to minimize a so-called **Hinge Loss** on a linear separator

Extensions of Perceptron

- **Voted Perceptron**
 - generalizes better than (standard) perceptron
 - memory intensive (keeps around every weight vector seen during training, so each one can vote)
- **Averaged Perceptron**
 - empirically similar performance to voted perceptron
 - can be implemented in a memory efficient way (running averages are efficient)
- **Kernel Perceptron**
 - Choose a kernel $K(x', x)$
 - Apply the **kernel trick** to Perceptron
 - Resulting algorithm is **still very simple**
- **Structured Perceptron**
 - Basic idea can also be applied when \mathbf{y} ranges over an exponentially large set
 - Mistake bound **does not** depend on the size of that set

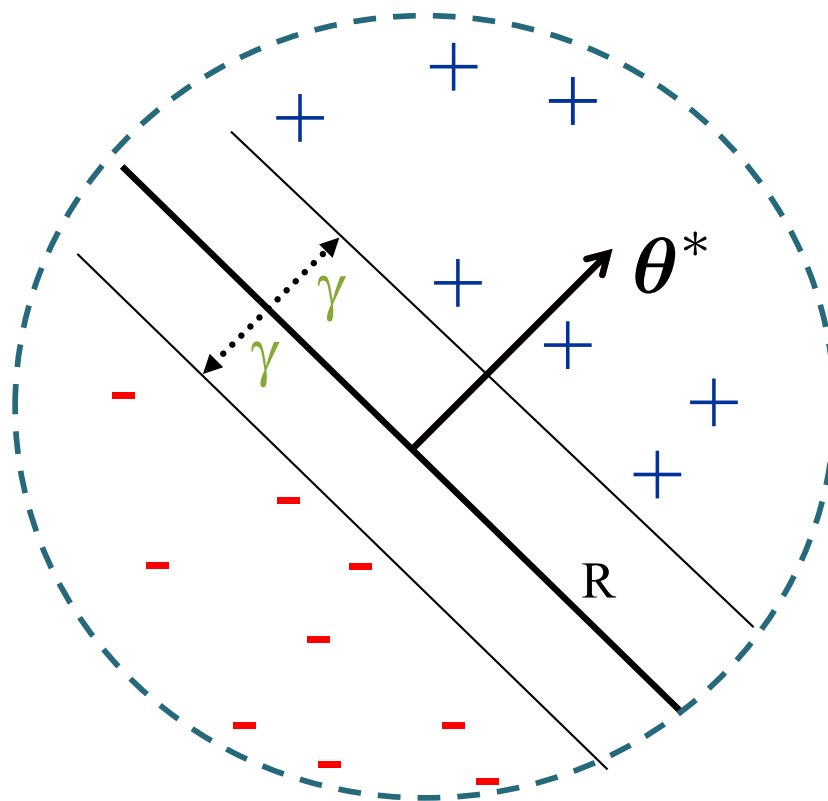
ANALYSIS OF PERCEPTRON

Analysis: Perceptron

Perceptron Mistake Bound

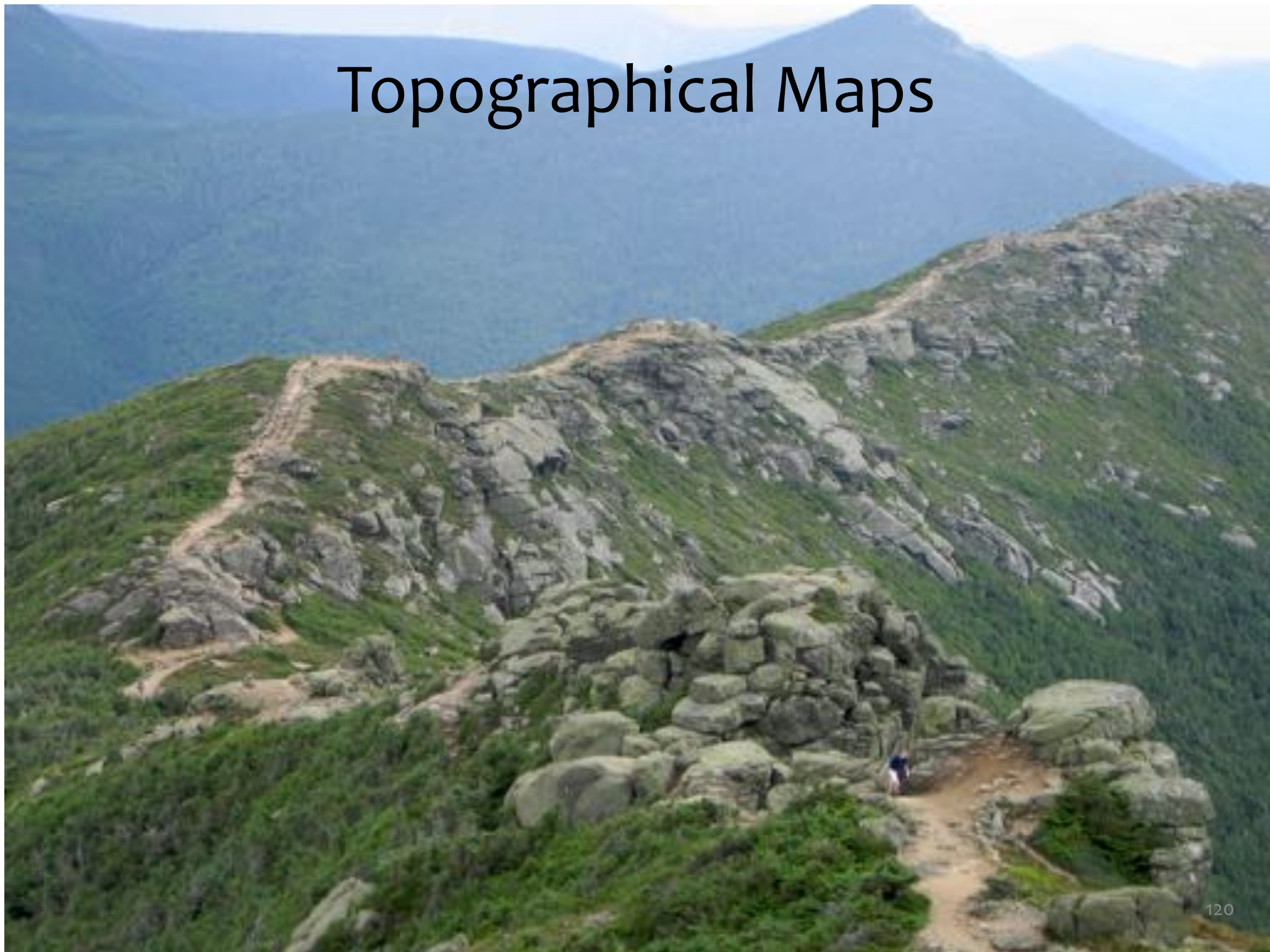
Guarantee: If data has margin γ and all points inside a ball of radius R , then Perceptron makes $\leq (R/\gamma)^2$ mistakes.

(Normalized margin: multiplying all points by 100, or dividing all points by 100, doesn't change the number of mistakes; algo is invariant to scaling.)

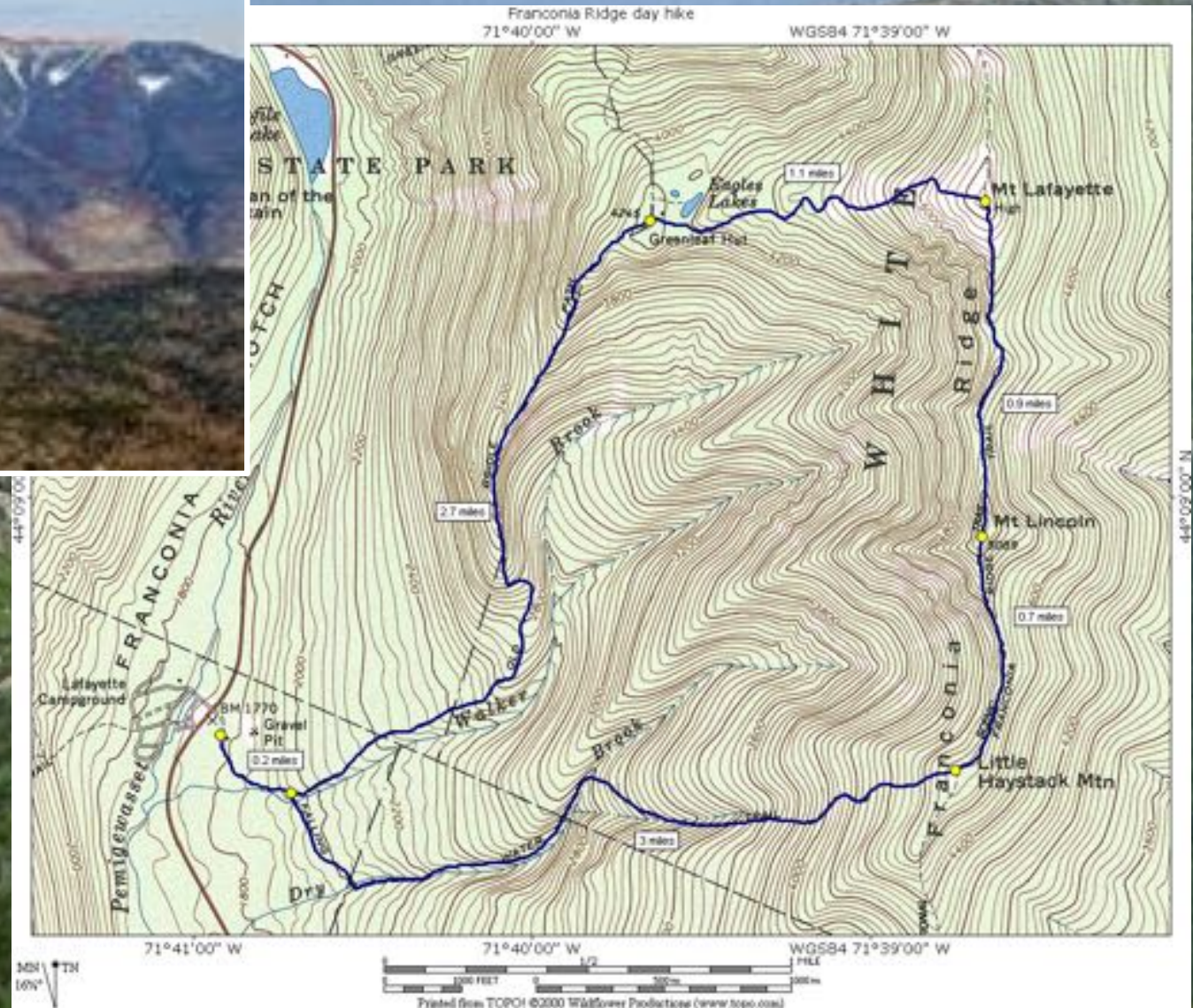


OPTIMIZATION FOR ML

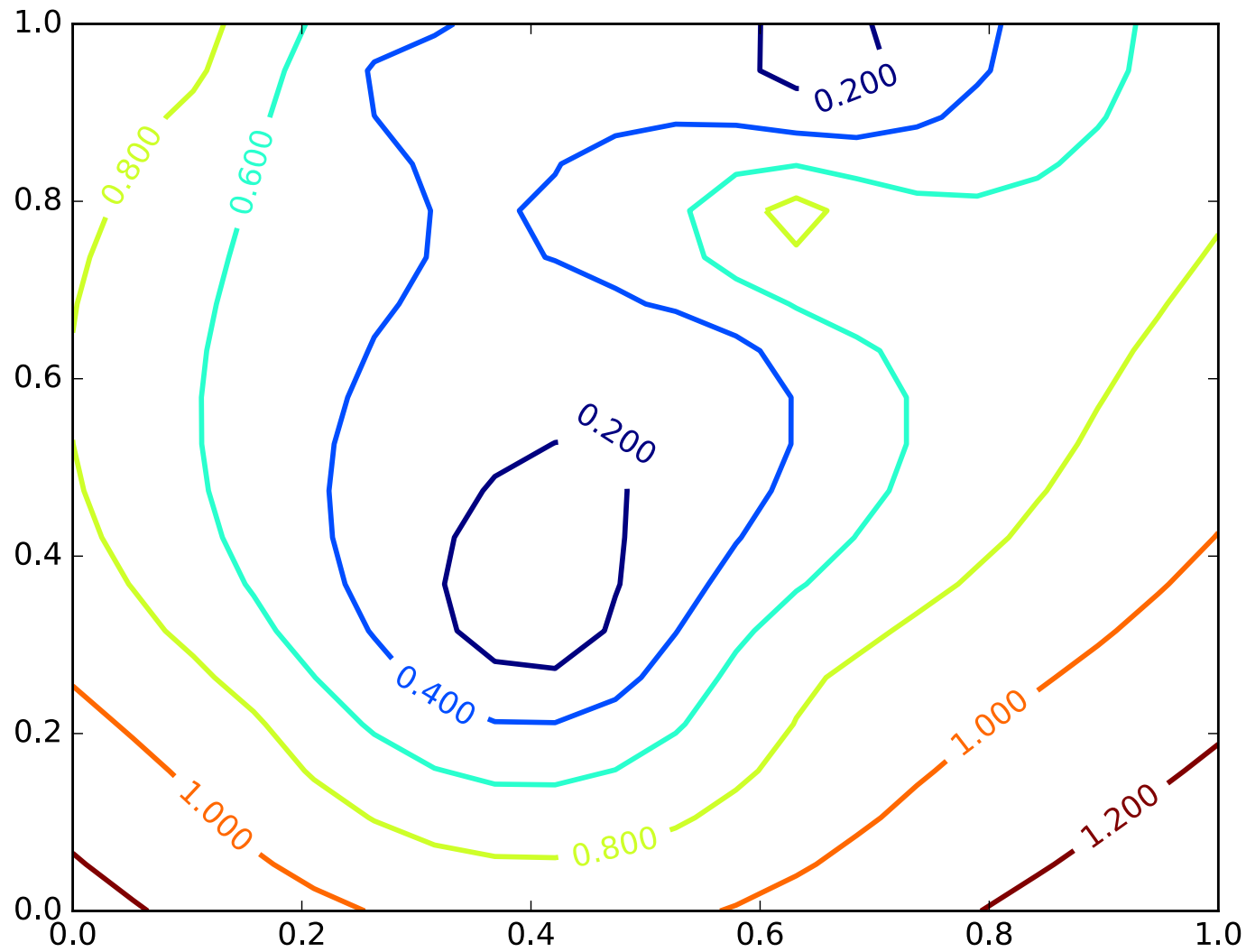
Topographical Maps



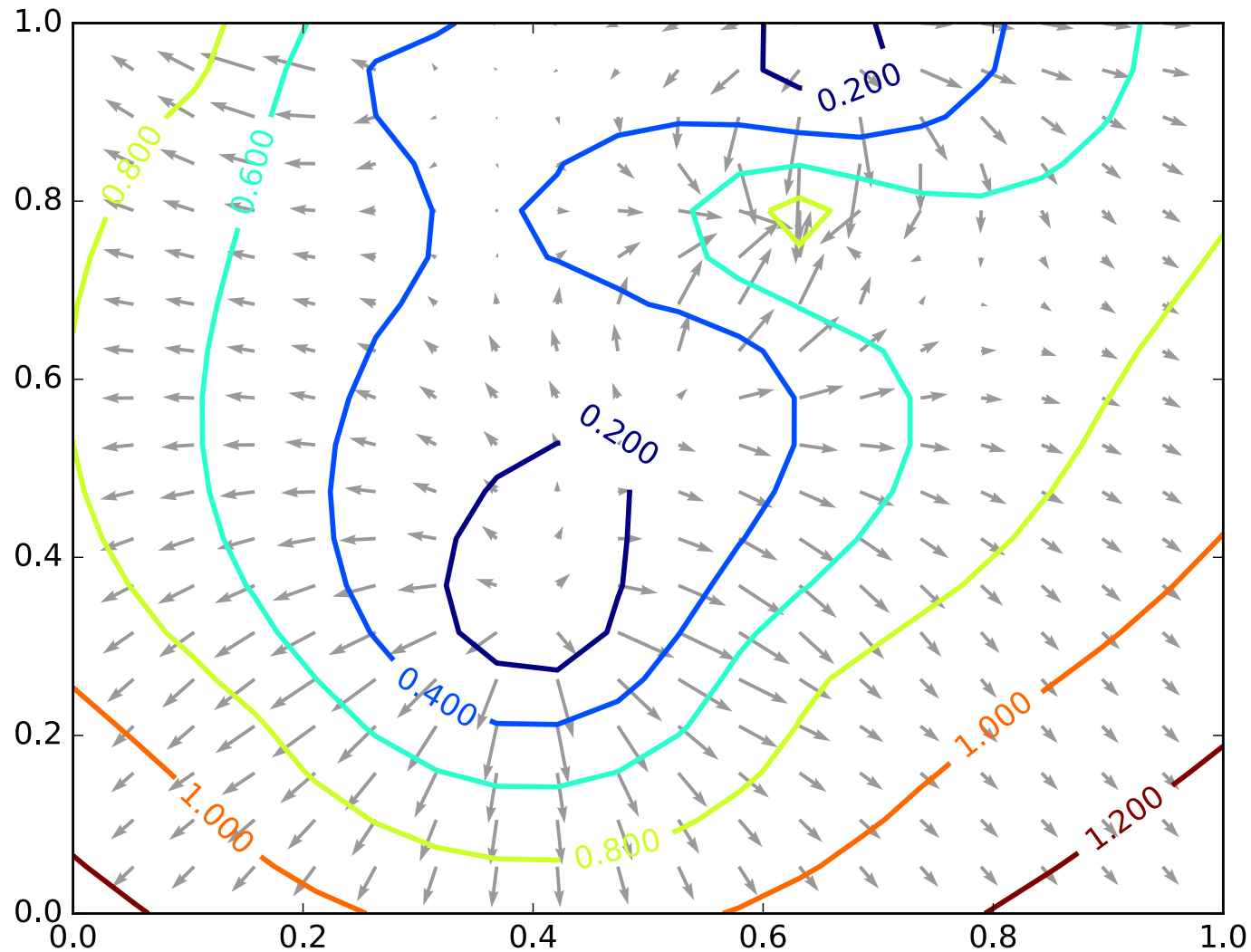
Topographical Maps



Gradients

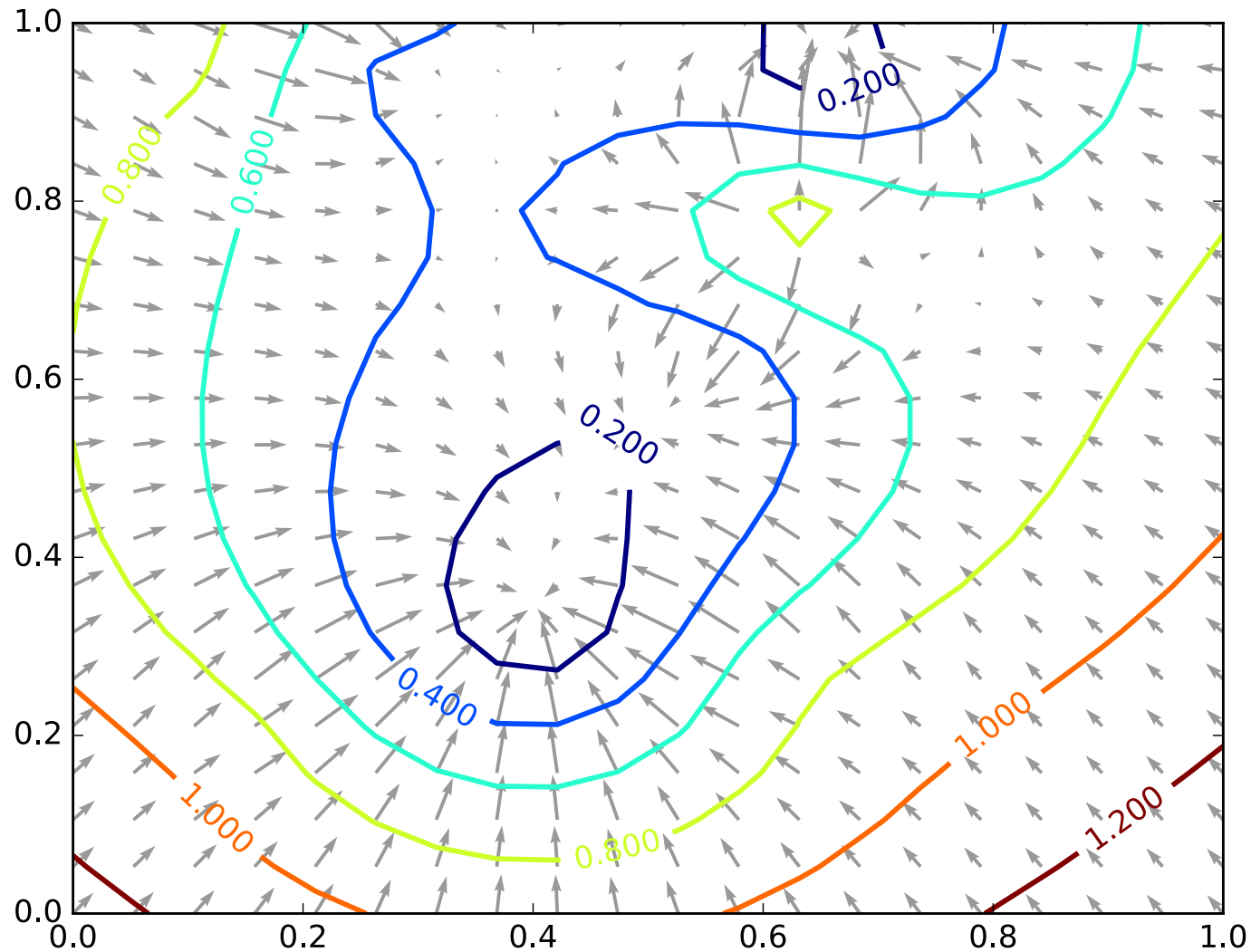


Gradients



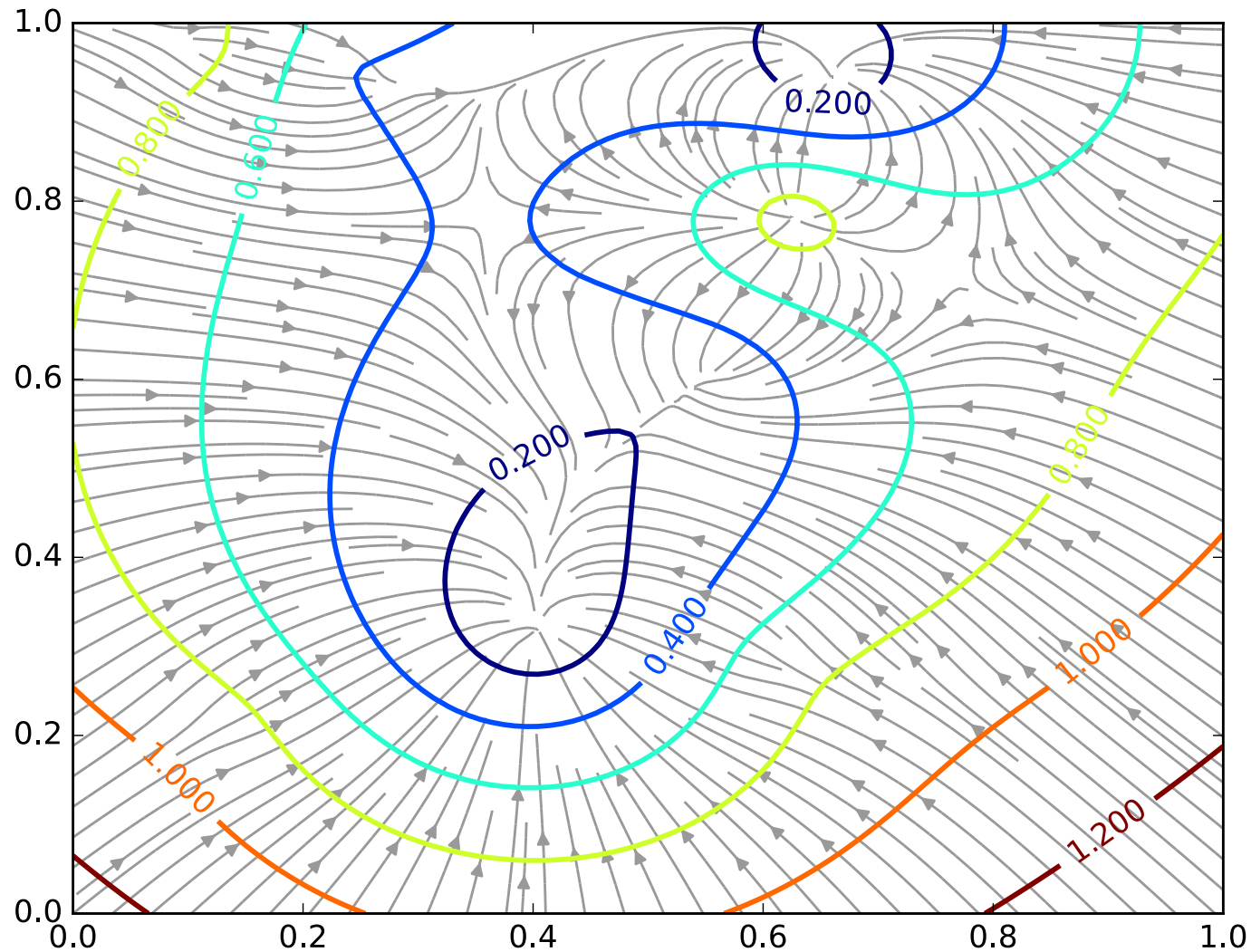
These are the **gradients** that
Gradient **Ascent** would follow.

(Negative) Gradients



These are the **negative** gradients that Gradient **Descent** would follow.

(Negative) Gradient *Paths*



Shown are the **paths** that Gradient Descent would follow if it were making **infinitesimally small steps**.

Gradient Descent

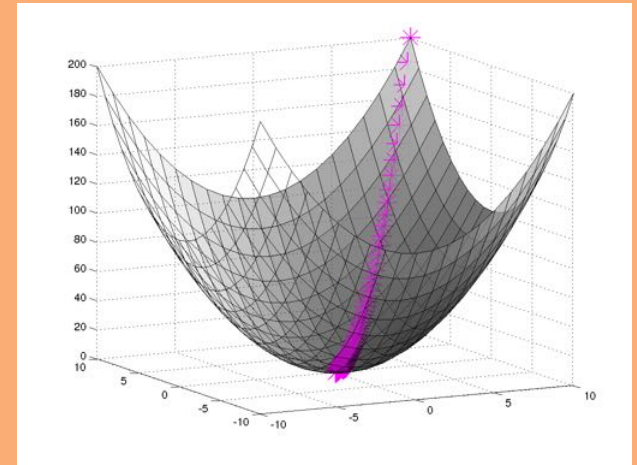
Chalkboard

- Example: 2D gradients
- Algorithm
- Details: starting point, stopping criterion, line search

Gradient Descent

Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \lambda \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



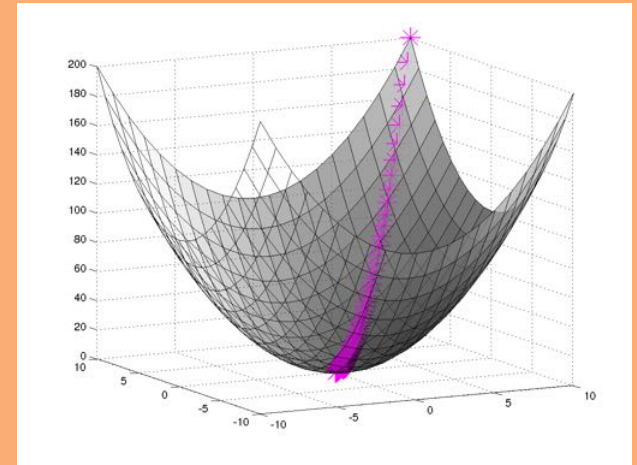
In order to apply GD to Linear Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J(\theta) \end{bmatrix}$$

Gradient Descent

Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )  
2:    $\theta \leftarrow \theta^{(0)}$   
3:   while not converged do  
4:      $\theta \leftarrow \theta - \lambda \nabla_{\theta} J(\theta)$   
5:   return  $\theta$ 
```



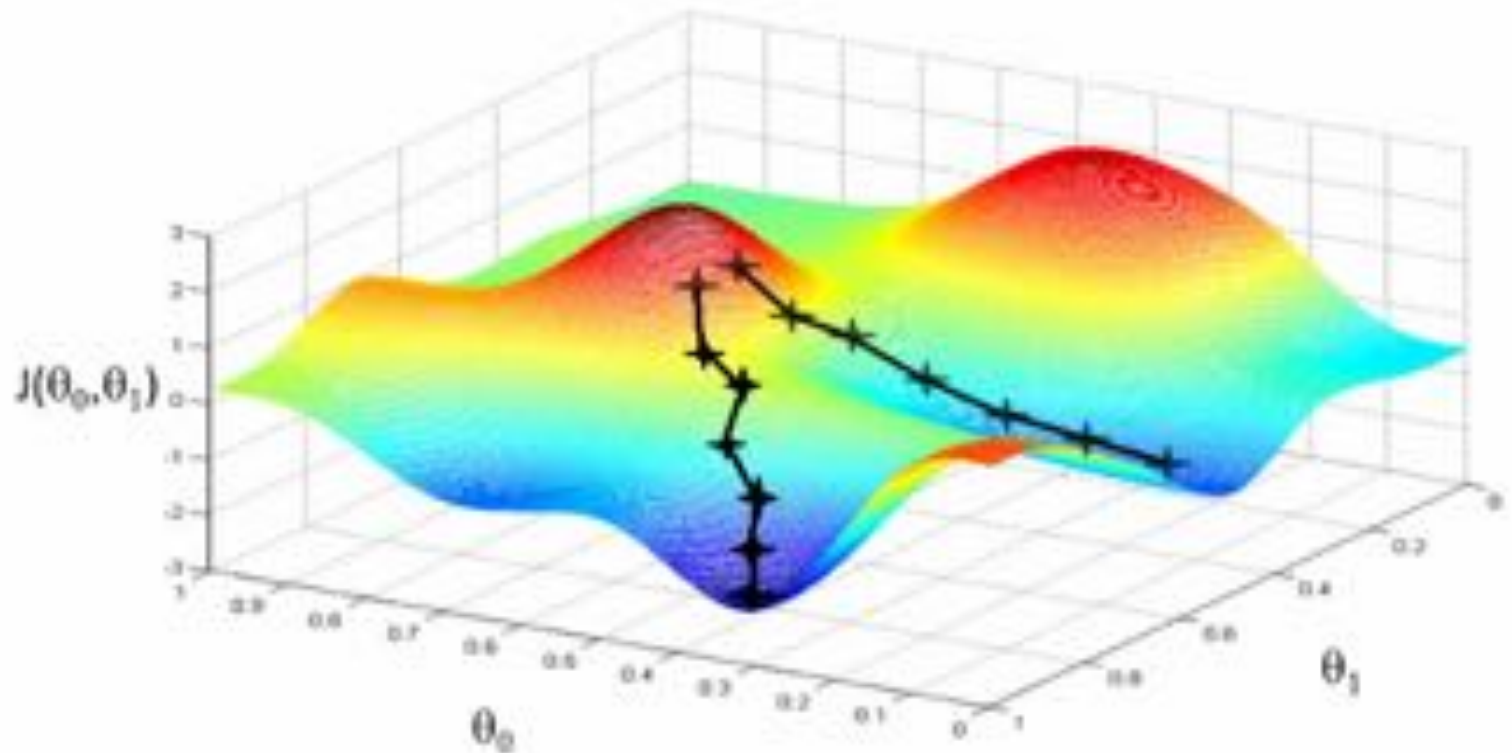
There are many possible ways to detect **convergence**. For example, we could check whether the L2 norm of the gradient is below some small tolerance.

$$\|\nabla_{\theta} J(\theta)\|_2 \leq \epsilon$$

Alternatively we could check that the reduction in the objective function from one iteration to the next is small.

Pros and cons of gradient descent

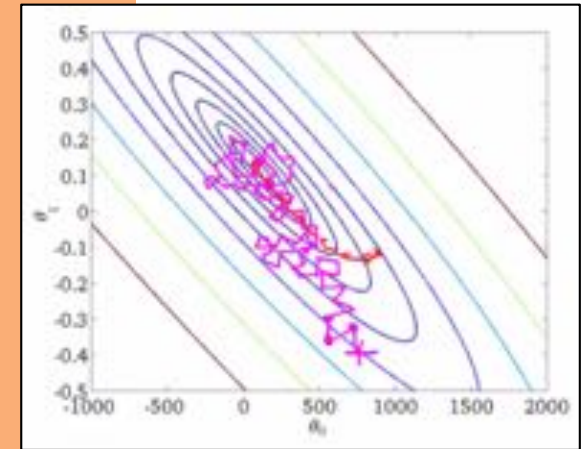
- Simple and often quite effective on ML tasks
- Often very scalable
- Only applies to smooth functions (differentiable)
- Might find a local minimum, rather than a global one



Stochastic Gradient Descent (SGD)

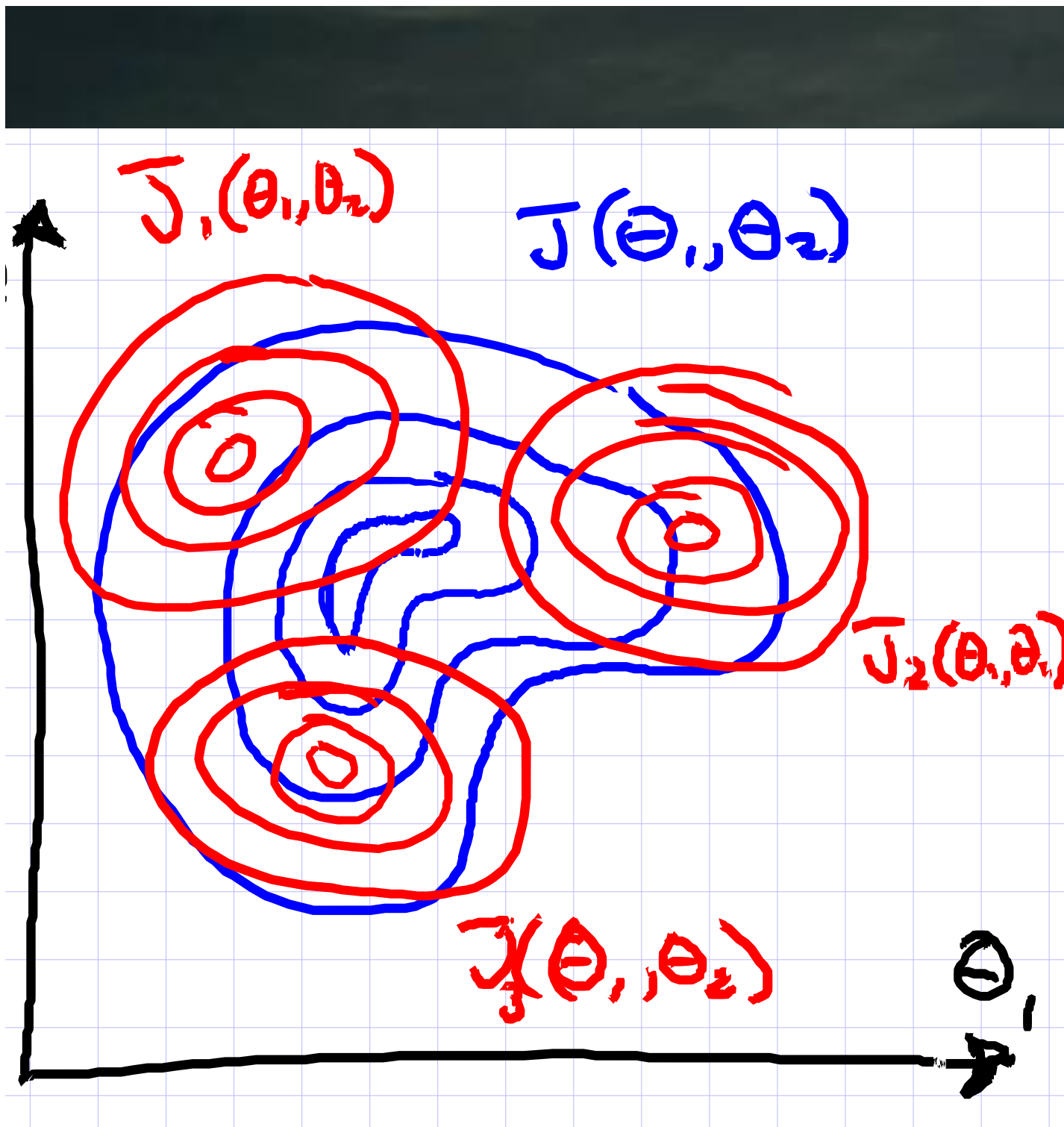
Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$



Expectations of Gradients

$$\frac{dJ(\vec{\theta})}{d\theta_j} = \frac{1}{N} \sum_{i=1}^N \frac{d}{d\theta_j} (J_i(\vec{\theta}))$$
$$\nabla J(\vec{\theta}) = \begin{bmatrix} \vdots \\ \text{jth} \\ \vdots \end{bmatrix} = \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta})$$

Recall: for any discrete r.v. X

$$E_X[f(x)] \triangleq \sum_x P(X=x) f(x)$$

Q: What is the expected value of a randomly chosen $\nabla J_i(\vec{\theta})$?

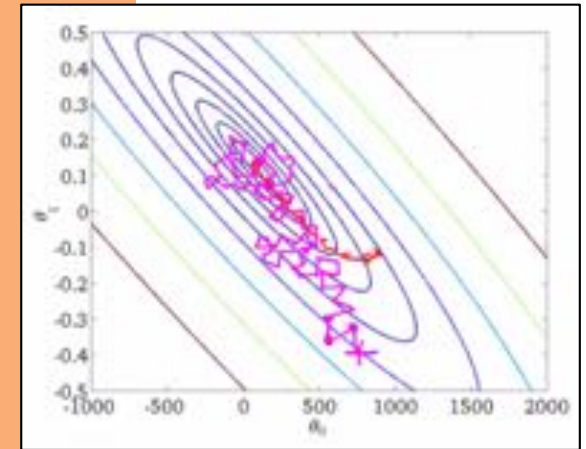
Let $I \sim \text{Uniform}(\{1, \dots, N\})$
 $\Rightarrow P(I=i) = \frac{1}{N}$ if $i \in \{1, \dots, N\}$

$$\begin{aligned} E_I[\nabla J_I(\vec{\theta})] &= \sum_{i=1}^N P(I=i) \nabla J_i(\vec{\theta}) \\ &= \frac{1}{N} \sum_{i=1}^N \nabla J_i(\vec{\theta}) \\ &= \nabla J(\vec{\theta}) \end{aligned}$$

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$ 
6:   return  $\theta$ 
```



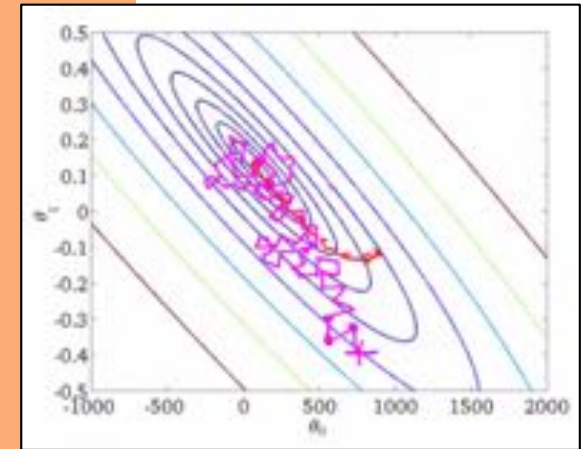
We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

Stochastic Gradient Descent (SGD)

Algorithm 2 Stochastic Gradient Descent (SGD)

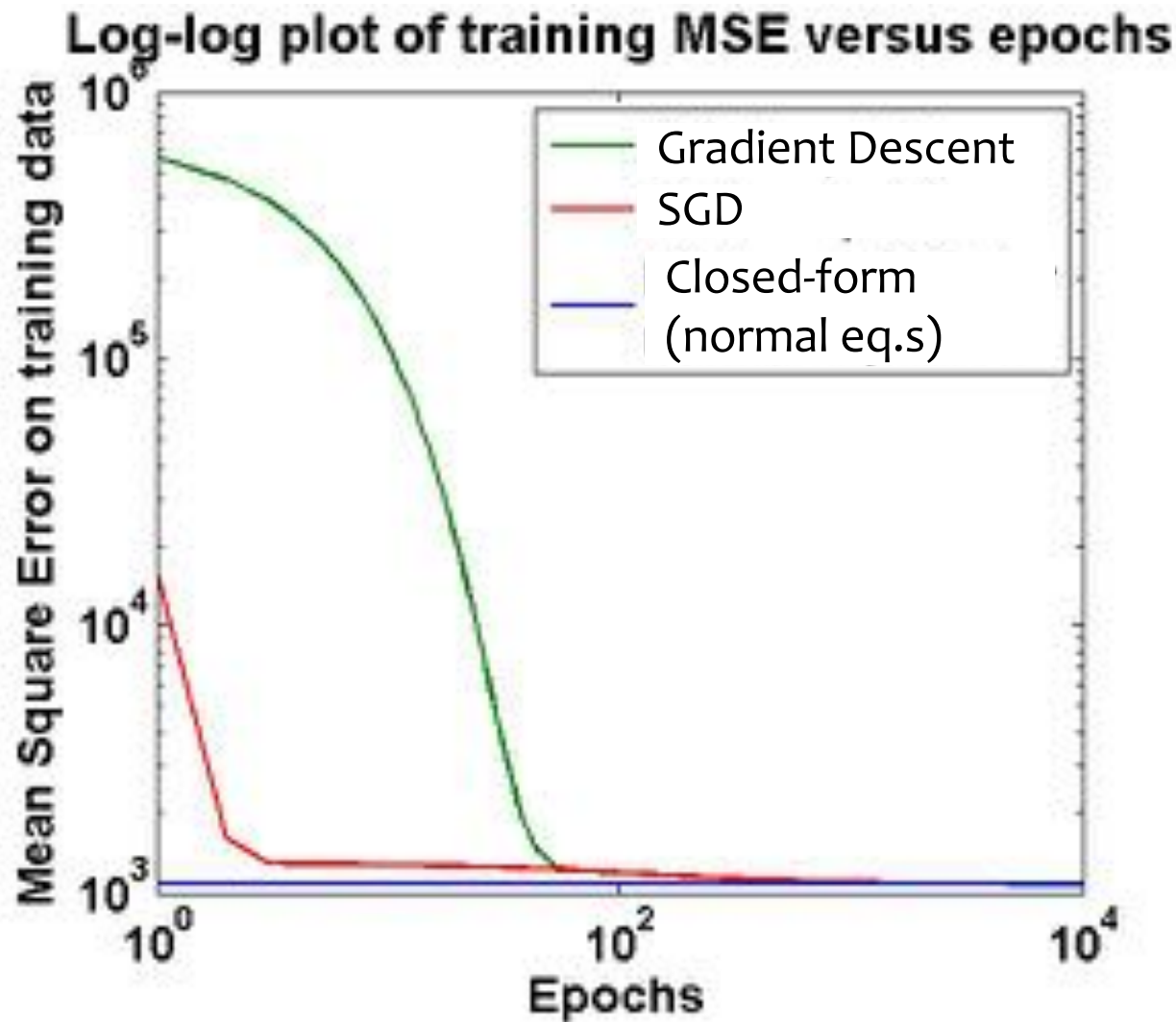
```
1: procedure SGD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:       for  $k \in \{1, 2, \dots, K\}$  do
6:          $\theta_k \leftarrow \theta_k - \lambda \frac{d}{d\theta_k} J^{(i)}(\theta)$ 
7:   return  $\theta$ 
```



We need a per-example objective:

$$\text{Let } J(\theta) = \sum_{i=1}^N J^{(i)}(\theta)$$

Convergence Curves



- Def: an **epoch** is a single pass through the training data
- 1. For GD, only **one update** per epoch
- 2. For SGD, **N updates** per epoch
 $N = (\# \text{ train examples})$

- SGD reduces MSE much more rapidly than GD
- For GD / SGD, training MSE is initially large due to uninformed initialization

Optimization Objectives

You should be able to...

- Apply gradient descent to optimize a function
- Apply stochastic gradient descent (SGD) to optimize a function
- Apply knowledge of zero derivatives to identify a closed-form solution (if one exists) to an optimization problem
- Distinguish between convex, concave, and nonconvex functions
- Obtain the gradient (and Hessian) of a (twice) differentiable function

PROBABILISTIC LEARNING

Probabilistic Learning

Function Approximation

Previously, we assumed that our output was generated using a **deterministic target function**:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} = c^*(\mathbf{x}^{(i)})$$

Our goal was to learn a hypothesis $h(\mathbf{x})$ that best approximates $c^*(\mathbf{x})$

Probabilistic Learning

Today, we assume that our output is **sampled** from a conditional **probability distribution**:

$$\mathbf{x}^{(i)} \sim p^*(\cdot)$$
$$y^{(i)} \sim p^*(\cdot | \mathbf{x}^{(i)})$$

Our goal is to learn a probability distribution $p(y|\mathbf{x})$ that best approximates $p^*(y|\mathbf{x})$

Robotic Farming

	Deterministic	Probabilistic
Classification (binary output)	Is this a picture of a wheat kernel?	Is this plant drought resistant?
Regression (continuous output)	How many wheat kernels are in this picture?	What will the yield of this plant be?



Maximum Likelihood Estimation

The principle of Maximum Likelihood Estimation (MLE):

Choose parameters that make the data "most likely".

Assumptions: Data generated iid from distribution $p^*(x | \vec{\theta}^*)$
and comes from a family of distributions parameterized
 $\theta \in \Theta$ \swarrow set of possible parameters

Formally:

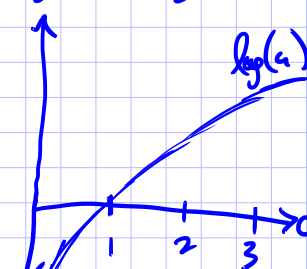
$$\begin{aligned}\theta_{MLE} &= \underset{\theta \in \Theta}{\operatorname{argmax}} p(D|\theta) \\ &= \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(D|\theta) \\ &= \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)\end{aligned}$$

usually
a continuous
optimization

where $\ell(\theta) \triangleq \log p(D|\theta)$
 \swarrow
'log-likelihood'

\swarrow treat as function of θ
where D is constant

since log is monotonic



$$\begin{aligned}\log(a_1) &< \log(a_2) \\ \text{iff } a_1 &< a_2 \\ \Rightarrow \log(f(a_1)) &< \log(f(a_2)) \\ \text{iff } f(a_1) &< f(a_2)\end{aligned}$$

MOTIVATION: LOGISTIC REGRESSION

Example: Image Classification

- ImageNet LSVRC-2010 contest:
 - **Dataset:** 1.2 million labeled images, 1000 classes
 - **Task:** Given a new image, label it with the correct class
 - **Multiclass** classification problem
- Examples from <http://image-net.org/>

Bird

Warm-blooded egg-laying vertebrates characterized by feathers and forelimbs modified as wings

2126
pictures92.85%
Popularity
Percentile

- marine animal, marine creature, sea animal, sea creature (1)
- scavenger (1)
- biped (0)
- predator, predatory animal (1)
- larva (49)
- acrodont (0)
- feeder (0)
- stunt (0)
- chordate (3087)
 - tunicate, urochordate, urochord (6)
 - cephalochordate (1)
 - vertebrate, craniate (3077)
 - mammal, mammalian (1169)
 - bird (871)
 - dickeybird, dickey-bird, dickybird, dicky-bird (0)
 - cock (1)
 - hen (0)
 - nester (0)
 - night bird (1)
 - bird of passage (0)
 - protoavis (0)
 - archaeopteryx, archeopteryx, Archaeopteryx lithographi
 - Sinornis (0)
 - libero-mesornis (0)
 - archaeornis (0)
 - ratite, ratite bird, flightless bird (10)
 - carinate, carinate bird, flying bird (0)
 - passerine, passeriform bird (279)
 - nonpasserine bird (0)
 - bird of prey, raptor, raptorial bird (80)
 - gallinaceous bird, gallinacean (114)

Treemap Visualization

Images of the Synset

Downloads



German iris, *Iris kochii*Iris of northern Italy having deep blue-purple flowers; similar to but smaller than *Iris germanica*469
pictures49.6%
Popularity
Percentile

- halophyte (0)
- succulent (39)
- cultivar (0)
- cultivated plant (0)
- weed (54)
- evergreen, evergreen plant (0)
- deciduous plant (0)
- vine (272)
- creeper (0)
- woody plant, ligneous plant (1868)
- geophyte (0)
- desert plant, xerophyte, xerophytic plant, xerophile, xerophilic
- mesophyte, mesophytic plant (0)
- aquatic plant, water plant, hydrophyte, hydrophytic plant (11)
- tuberous plant (0)
- bulbous plant (179)
 - Iridaceous plant (27)
 - iris, flag, fleur-de-lis, sword lily (19)
 - bearded iris (4)
 - Florentine iris, orris, *Iris germanica florentina*, *Iris*
 - German iris, *Iris germanica* (0)
 - German iris, *Iris kochii* (0)
 - Dalmatian iris, *Iris pallida* (0)
 - beardless iris (4)
 - bulbous iris (0)
 - dwarf iris, *Iris cristata* (0)
 - stinking iris, gladdon, gladdon iris, stinking gladdwyn,
 - Persian iris, *Iris persica* (0)
 - yellow iris, yellow flag, yellow water flag, *Iris pseudo*
 - dwarf iris, vernal iris, *Iris verna* (0)
 - blue flag, *Iris versicolor* (0)

Treemap Visualization

Images of the Synset

Downloads



Court, courtyard

An area wholly or partly surrounded by walls or buildings; "the house was built around an inner court"

165
pictures

92.61%
Popularity
Percentile



Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32326)
 - plant, flora, plant life (4486)
 - geological formation, formation (175)
 - natural object (1112)
 - sport, athletics (176)
 - artifact, artefact (10504)
 - instrumentality, instrumentation (5494)
 - structure, construction (1405)
 - airdock, hangar, repair shed (0)
 - altar (1)
 - arcade, colonnade (1)
 - arch (31)
 - area (344)
 - aisle (0)
 - auditorium (1)
 - baggage claim (0)
 - box (1)
 - breakfast area, breakfast nook (0)
 - bullpen (0)
 - chancel, sanctuary, bema (0)
 - choir (0)
 - corner, nook (2)
 - court, courtyard (6)
 - atrium (0)
 - bailey (0)
 - cloister (0)
 - food court (0)
 - forecourt (0)
 - narvik (0)

Treemap Visualization

Images of the Synset

Downloads



Example: Image Classification

CNN for Image Classification

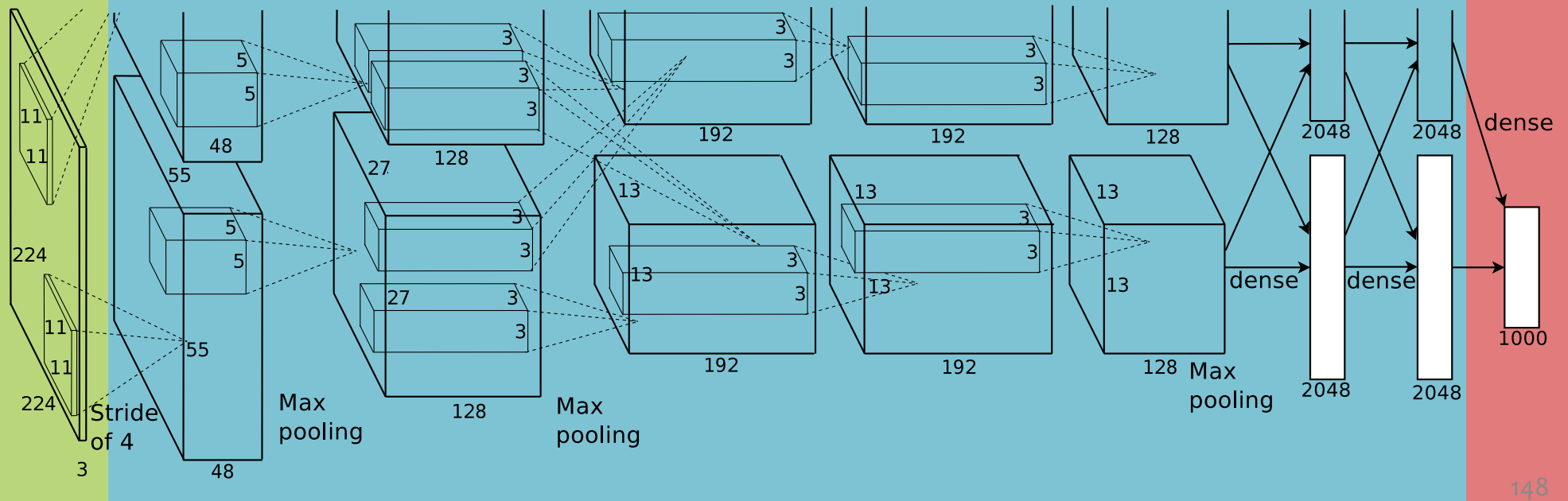
(Krizhevsky, Sutskever & Hinton, 2011)

17.5% error on ImageNet LSVRC-2010 contest

Input
image
(pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way
softmax



Example: Image Classification

CNN for Image Classification

(Krizhevsky, Sutskever & Hinton, 2011)

17.5% error on ImageNet LSVRC-2010 contest

Input
image
(pixels)

- Five convolutional layers (w/max-pooling)
- Three fully connected layers

1000-way softmax

The rest is *just*
some fancy
feature extraction
(discussed later in
the course)

This “softmax” layer is Logistic Regression!

dense

dense

1000

pooling

2048


2048

LOGISTIC REGRESSION

Logistic Regression

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \text{ where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{0, 1\}$$



We are back to
classification.

Despite the name
logistic **regression**.

Recall...

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:

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

Recall...

Background: Hyperplanes

Notation Trick: fold the bias b and the weights \mathbf{w} into a single vector $\boldsymbol{\theta}$ by prepending a constant to \mathbf{x} and increasing dimensionality by one!

Hyperplane (Definition 1):

$$\mathcal{H} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} = b\}$$

Hyperplane (Definition 2):

$$\mathcal{H} = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} = 0$$

$$\text{and } x_0 = 1\}$$

$$\boldsymbol{\theta} = [b, w_1, \dots, w_M]^T$$

Half-spaces:

$$\mathcal{H}^+ = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} > 0 \text{ and } x_0 = 1\}$$

$$\mathcal{H}^- = \{\mathbf{x} : \boldsymbol{\theta}^T \mathbf{x} < 0 \text{ and } x_0 = 1\}$$

Using gradient ascent for linear classifiers

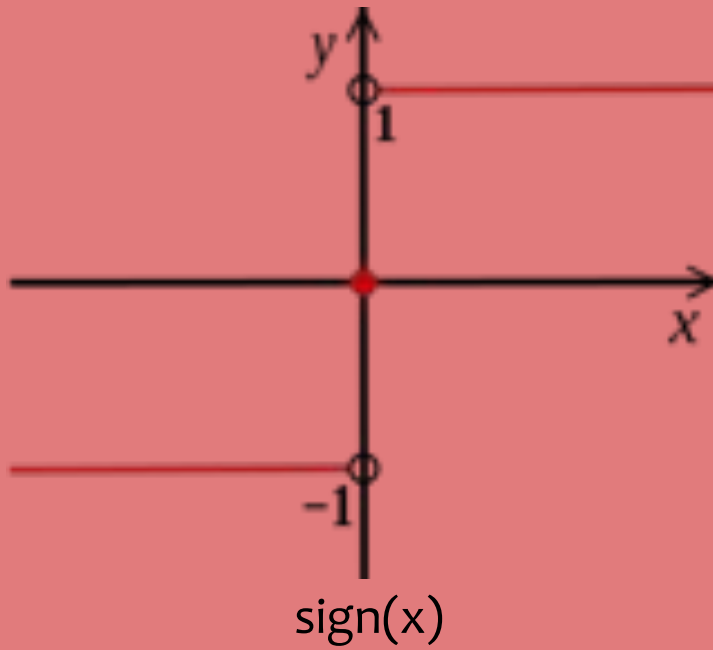
Key idea behind today's lecture:

1. Define a linear classifier (logistic regression)
2. Define an objective function (likelihood)
3. Optimize it with gradient descent to learn parameters
4. Predict the class with highest probability under the model

Using gradient ascent for linear classifiers

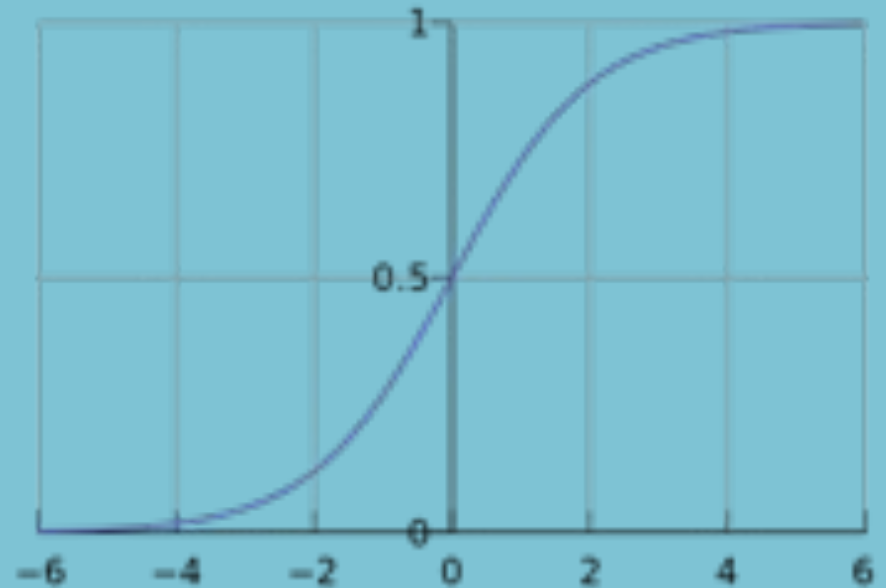
This decision function isn't differentiable:

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



Use a differentiable function instead:

$$p_{\boldsymbol{\theta}}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

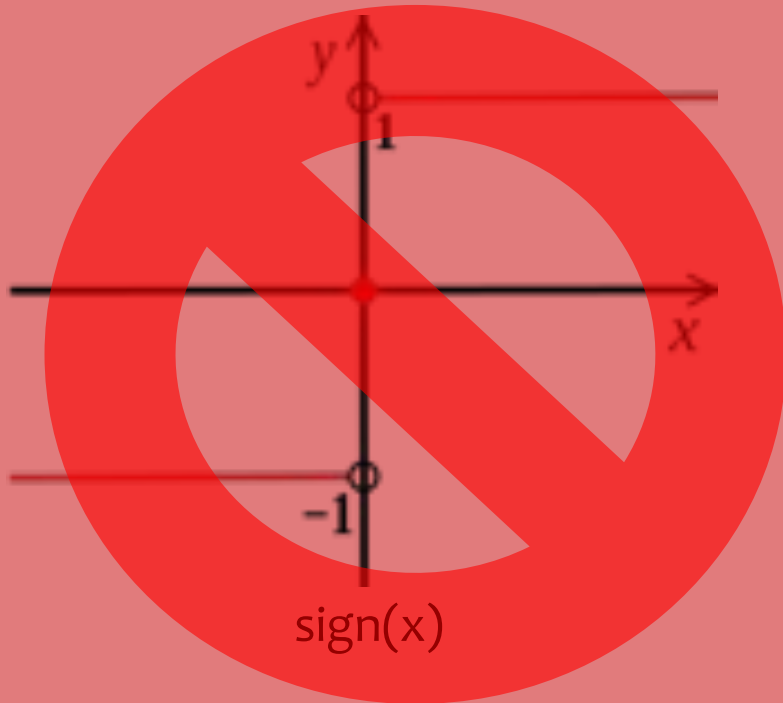


$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

Using gradient ascent for linear classifiers

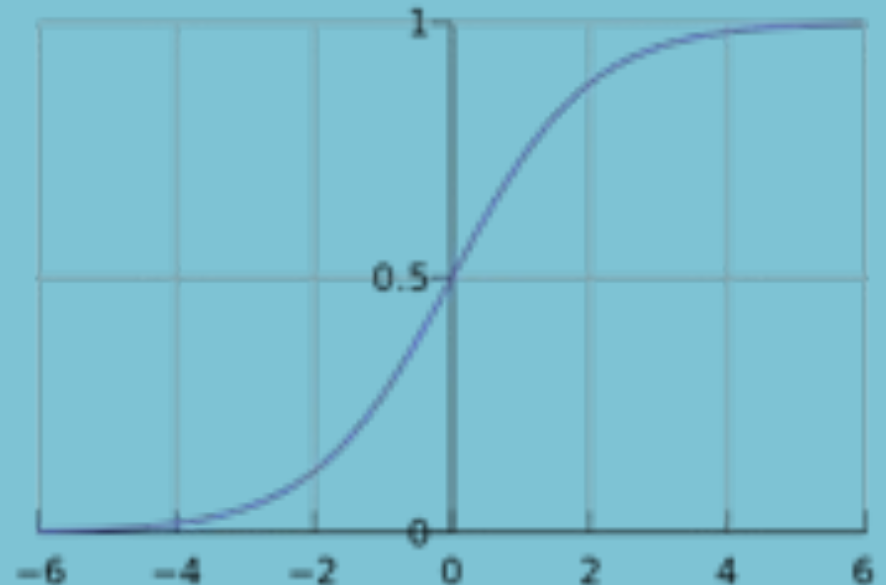
This decision function isn't differentiable:

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



Use a differentiable function instead:

$$p_{\boldsymbol{\theta}}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$



$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

Logistic Regression

Data: Inputs are continuous vectors of length M . Outputs are discrete.

$$\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^N \text{ where } \mathbf{x} \in \mathbb{R}^M \text{ and } y \in \{0, 1\}$$

Model: Logistic function applied to dot product of parameters with input vector.

$$p_{\boldsymbol{\theta}}(y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})}$$

Learning: finds the parameters that minimize some objective function. $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} J(\boldsymbol{\theta})$

Prediction: Output is the most probable class.

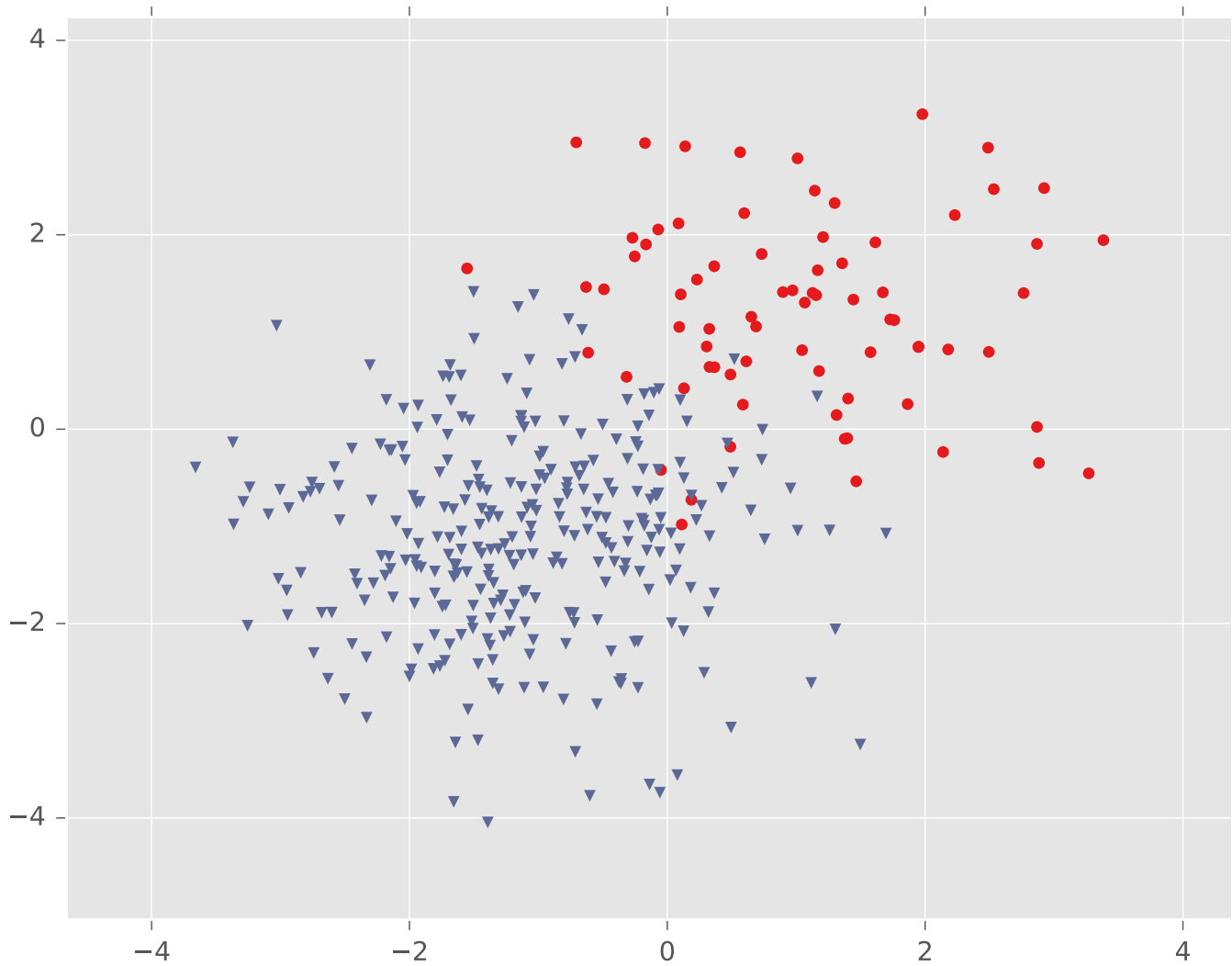
$$\hat{y} = \underset{y \in \{0,1\}}{\operatorname{argmax}} p_{\boldsymbol{\theta}}(y|\mathbf{x})$$

Logistic Regression

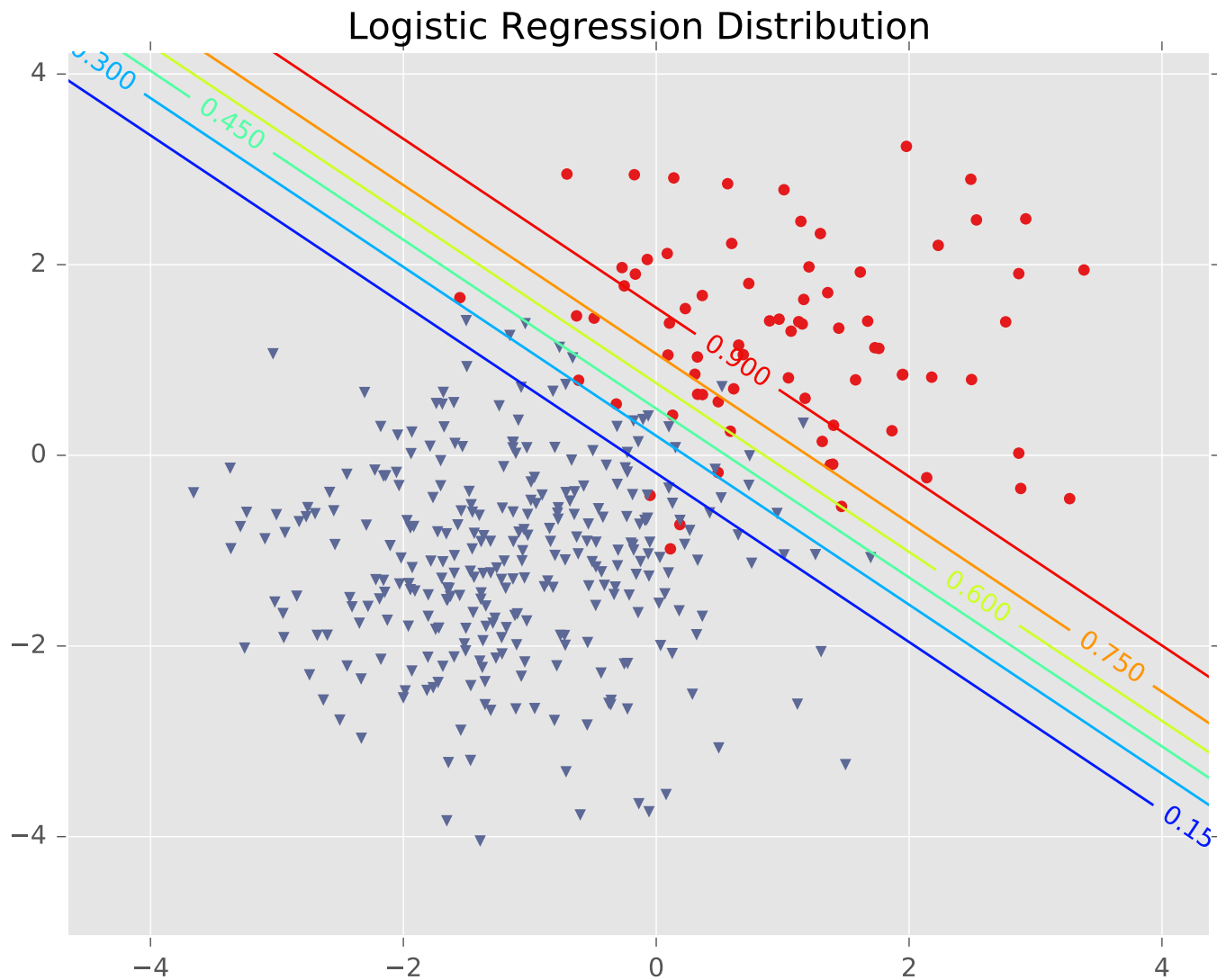
Whiteboard

- Bernoulli interpretation
- Logistic Regression Model
- Decision boundary

Logistic Regression

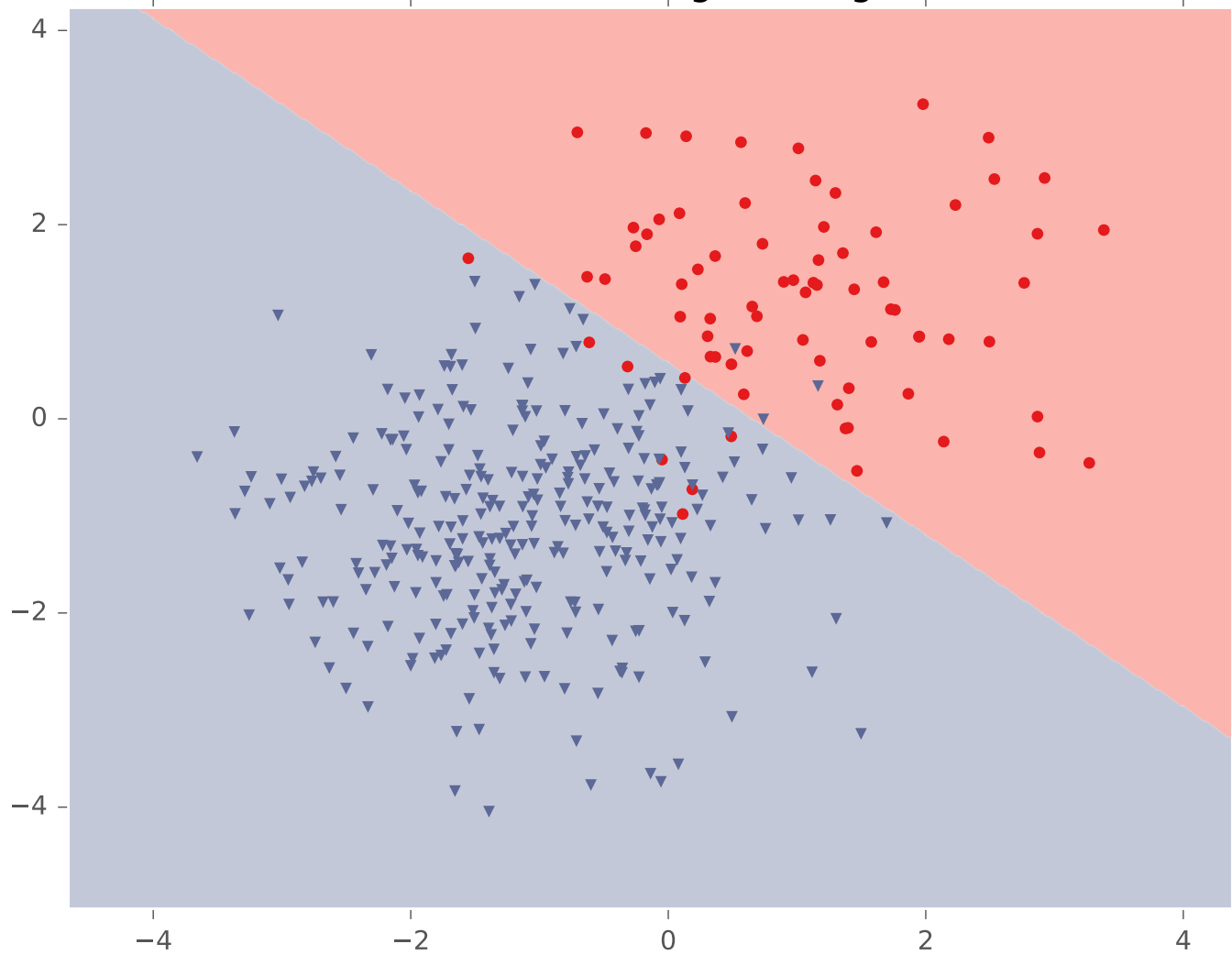


Logistic Regression



Logistic Regression

Classification with Logistic Regression



LEARNING LOGISTIC REGRESSION

Maximum Conditional Likelihood Estimation

Learning: finds the parameters that minimize some objective function.

$$\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

We minimize the *negative* log conditional likelihood:

$$J(\theta) = -\log \prod_{i=1}^N p_{\theta}(y^{(i)} | \mathbf{x}^{(i)})$$

Why?

1. We can't maximize likelihood (as in Naïve Bayes) because we don't have a joint model $p(\mathbf{x}, y)$
2. It worked well for Linear Regression (least squares is MCLE)

Maximum Conditional Likelihood Estimation

Learning: Four approaches to solving $\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$

Approach 1: Gradient Descent

(take larger – more certain – steps opposite the gradient)

Approach 2: Stochastic Gradient Descent (SGD)

(take many small steps opposite the gradient)

Approach 3: Newton's Method

(use second derivatives to better follow curvature)

Approach 4: Closed Form???

(set derivatives equal to zero and solve for parameters)

Maximum Conditional Likelihood Estimation

Learning: Four approaches to solving $\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$

Approach 1: Gradient Descent

(take larger – more certain – steps opposite the gradient)

Approach 2: Stochastic Gradient Descent (SGD)

(take many small steps opposite the gradient)

Approach 3: Newton's Method

(use second derivatives to better follow curvature)

~~**Approach 4:** Closed Form???~~

~~(set derivatives equal to zero and solve for parameters)~~

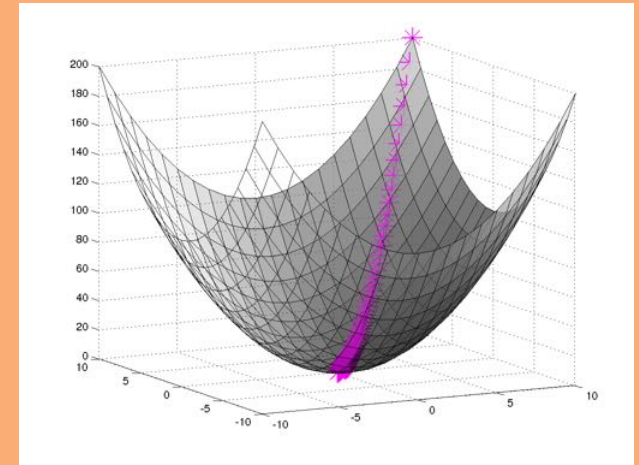
Logistic Regression does not have a closed form solution for MLE parameters.

Recall...

Gradient Descent

Algorithm 1 Gradient Descent

```
1: procedure GD( $\mathcal{D}$ ,  $\theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:      $\theta \leftarrow \theta - \lambda \nabla_{\theta} J(\theta)$ 
5:   return  $\theta$ 
```



In order to apply GD to Logistic Regression all we need is the **gradient** of the objective function (i.e. vector of partial derivatives).

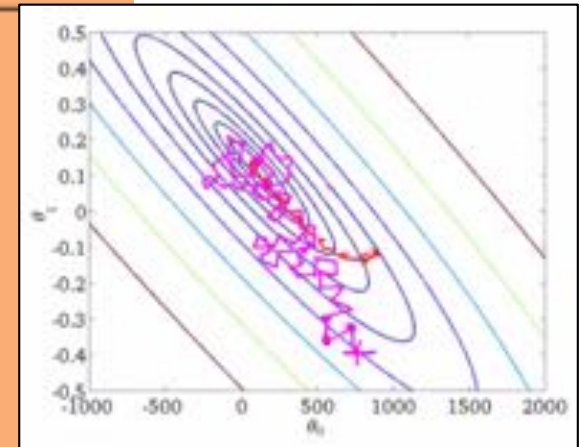
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{d}{d\theta_1} J(\theta) \\ \frac{d}{d\theta_2} J(\theta) \\ \vdots \\ \frac{d}{d\theta_M} J(\theta) \end{bmatrix}$$

Stochastic Gradient Descent (SGD)

Recall...

Algorithm 1 Stochastic Gradient Descent (SGD)

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )  
2:    $\theta \leftarrow \theta^{(0)}$   
3:   while not converged do  
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do  
5:        $\theta \leftarrow \theta - \lambda \nabla_{\theta} J^{(i)}(\theta)$   
6:   return  $\theta$ 
```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

$$\text{Let } J(\boldsymbol{\theta}) = \sum_{i=1}^N J^{(i)}(\boldsymbol{\theta})$$
$$\text{where } J^{(i)}(\boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(y^i | \mathbf{x}^i).$$

GRADIENT FOR LOGISTIC REGRESSION

Learning for Logistic Regression

Whiteboard

- Partial derivative for Logistic Regression
- Gradient for Logistic Regression

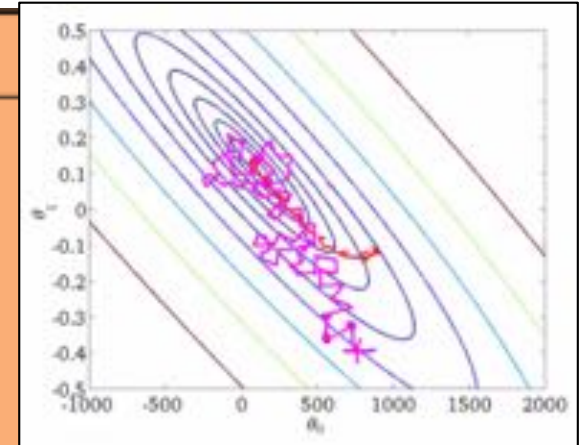
Details: Picking learning rate

- Use grid-search in log-space over small values on a tuning set:
 - e.g., 0.01, 0.001, ...
- Sometimes, decrease after each pass:
 - e.g factor of $1/(1 + dt)$, $t=\text{epoch}$
 - sometimes $1/t^2$
- Fancier techniques I won't talk about:
 - Adaptive gradient: scale gradient differently for each dimension (Adagrad, ADAM,)

SGD for Logistic Regression

Algorithm 1 SGD for Logistic Regression

```
1: procedure SGD( $\mathcal{D}, \theta^{(0)}$ )
2:    $\theta \leftarrow \theta^{(0)}$ 
3:   while not converged do
4:     for  $i \in \text{shuffle}(\{1, 2, \dots, N\})$  do
5:        $\theta \leftarrow \theta - \lambda(y^{(i)} - \rho^{(i)})\mathbf{x}^{(i)}$ 
6:       where  $\rho^{(i)} := 1/(1 + \exp(-\theta^T \mathbf{x}))$ 
7:   return  $\theta$ 
```



We can also apply SGD to solve the MCLE problem for Logistic Regression.

We need a per-example objective:

$$\text{Let } J(\boldsymbol{\theta}) = \sum_{i=1}^N J^{(i)}(\boldsymbol{\theta})$$
$$\text{where } J^{(i)}(\boldsymbol{\theta}) = -\log p_{\boldsymbol{\theta}}(y^i | \mathbf{x}^i).$$

Summary

1. Discriminative classifiers directly model the **conditional**, $p(y|x)$
2. Logistic regression is a **simple linear classifier**, that retains a **probabilistic semantics**
3. Parameters in LR are learned by **iterative optimization** (e.g. SGD)

Logistic Regression Objectives

You should be able to...

- Apply the principle of maximum likelihood estimation (MLE) to learn the parameters of a probabilistic model
- Given a discriminative probabilistic model, derive the conditional log-likelihood, its gradient, and the corresponding Bayes Classifier
- Explain the practical reasons why we work with the **log** of the likelihood
- Implement logistic regression for binary or multiclass classification
- Prove that the decision boundary of binary logistic regression is linear
- For linear regression, show that the parameters which minimize squared error are equivalent to those that maximize conditional likelihood

MULTINOMIAL LOGISTIC REGRESSION

Multinomial Logistic Regression

Chalkboard

- Background: Multinomial distribution
- Definition: Multi-class classification
- Geometric intuitions
- Multinomial logistic regression model
- Generative story
- Reduction to binary logistic regression
- Partial derivatives and gradients
- Applying Gradient Descent and SGD
- Implementation w/ sparse features

Debug that Program!

In-Class Exercise: *Think-Pair-Share*

Debug the following program which is (incorrectly) attempting to run SGD for multinomial logistic regression

Buggy Program:

```
while not converged:
    for i in shuffle([1,...,N]):
        for k in [1,...,K]:
            theta[k] = theta[k] - lambda * grad(x[i], y[i],
theta, k)
```

Assume: `grad(x[i], y[i], theta, k)` returns the gradient of the negative log-likelihood of the training example $(x[i], y[i])$ with respect to vector `theta[k]`. `lambda` is the learning rate. `N` = # of examples. `K` = # of output classes. `M` = # of features. `theta` is a `K` by `M` matrix.

Debug that Program!

In-Class Exercise: *Think-Pair-Share*

Debug the following program which is (incorrectly) attempting to run SGD for multinomial logistic regression

Buggy Program:

```
while not converged:
    for i in shuffle([1,...,N]):
        for k in [1,...,K]:
            for m in [1,..., M]:
                theta[k,m] = theta[k,m] + lambda * grad(x[i],
y[i], theta, k,m)
```

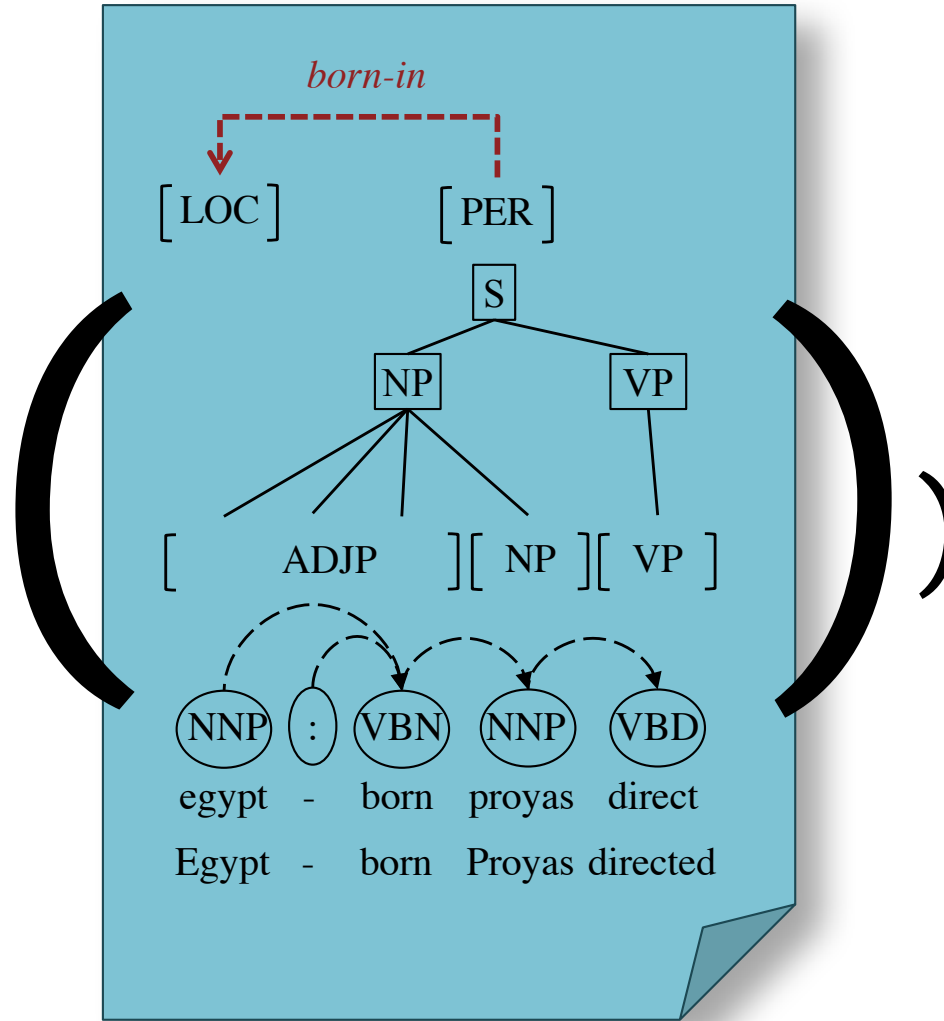
Assume: `grad(x[i], y[i], theta, k, m)` returns the partial derivative of the negative log-likelihood of the training example $(x[i], y[i])$ with respect to $\theta[k, m]$. `lambda` is the learning rate. `N` = # of examples. `K` = # of output classes. `M` = # of features. `theta` is a `K` by `M` matrix.

FEATURE ENGINEERING

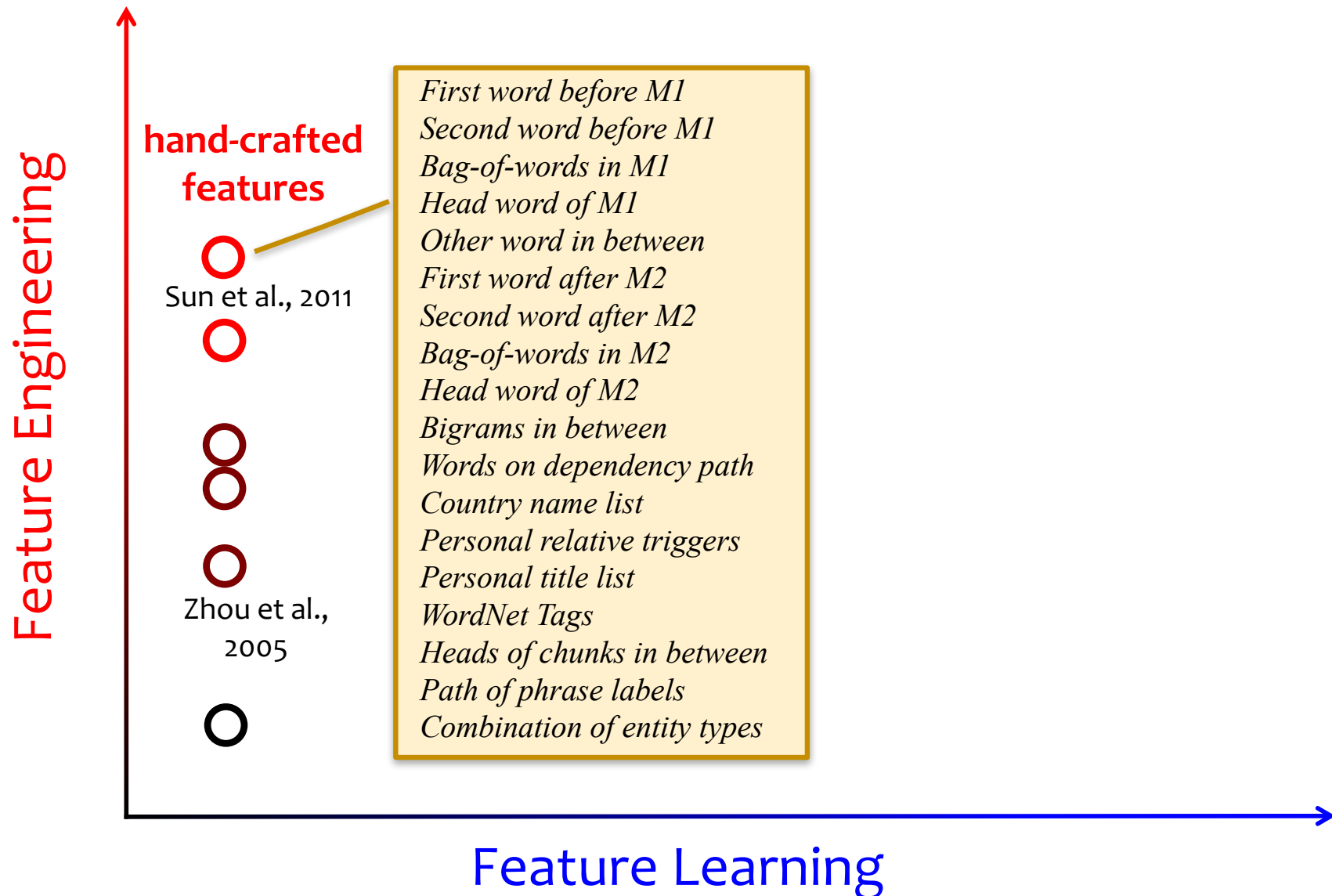
Handcrafted Features

$$p(y|x) \propto$$

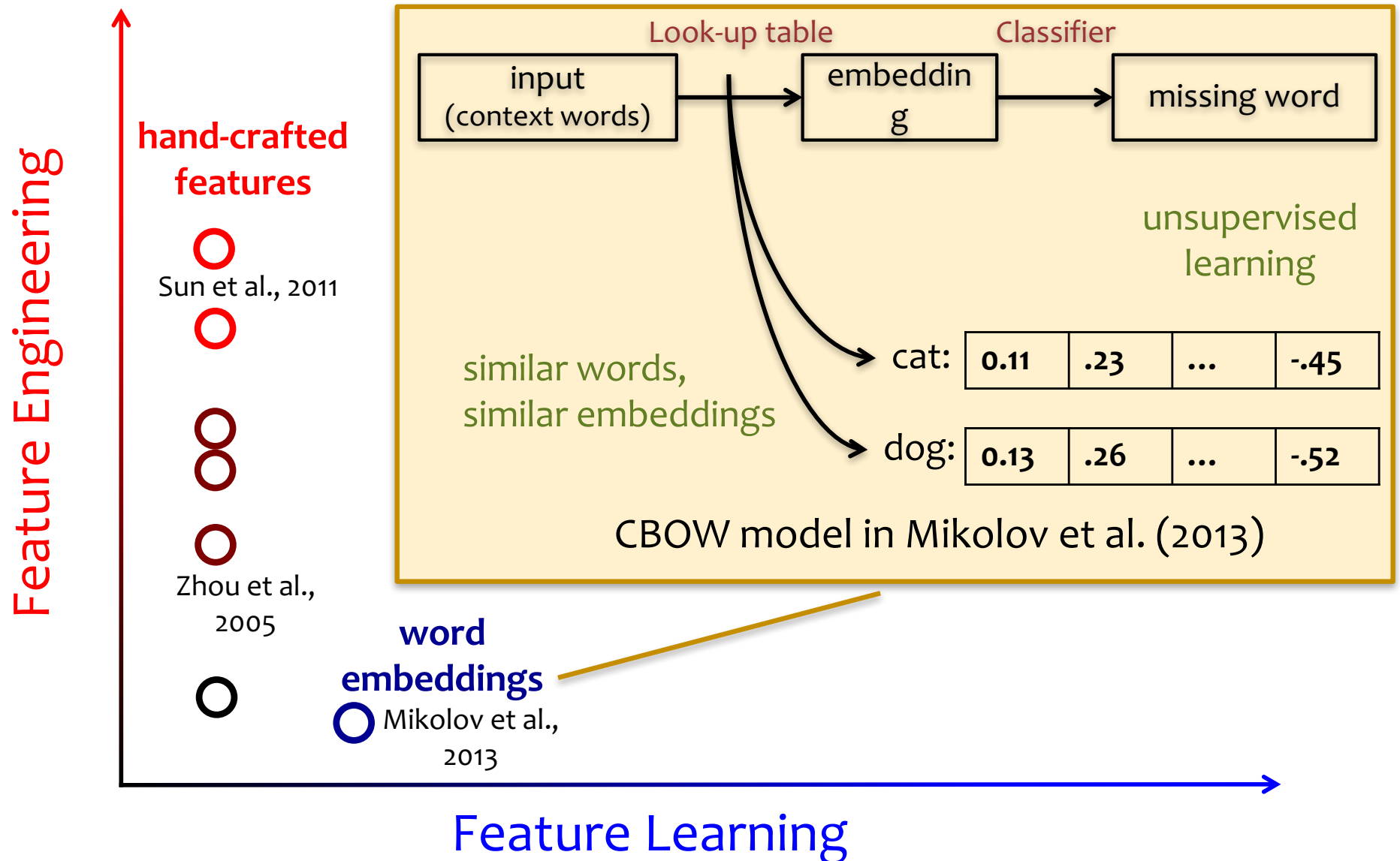
$$\exp(\Theta_y \cdot f$$



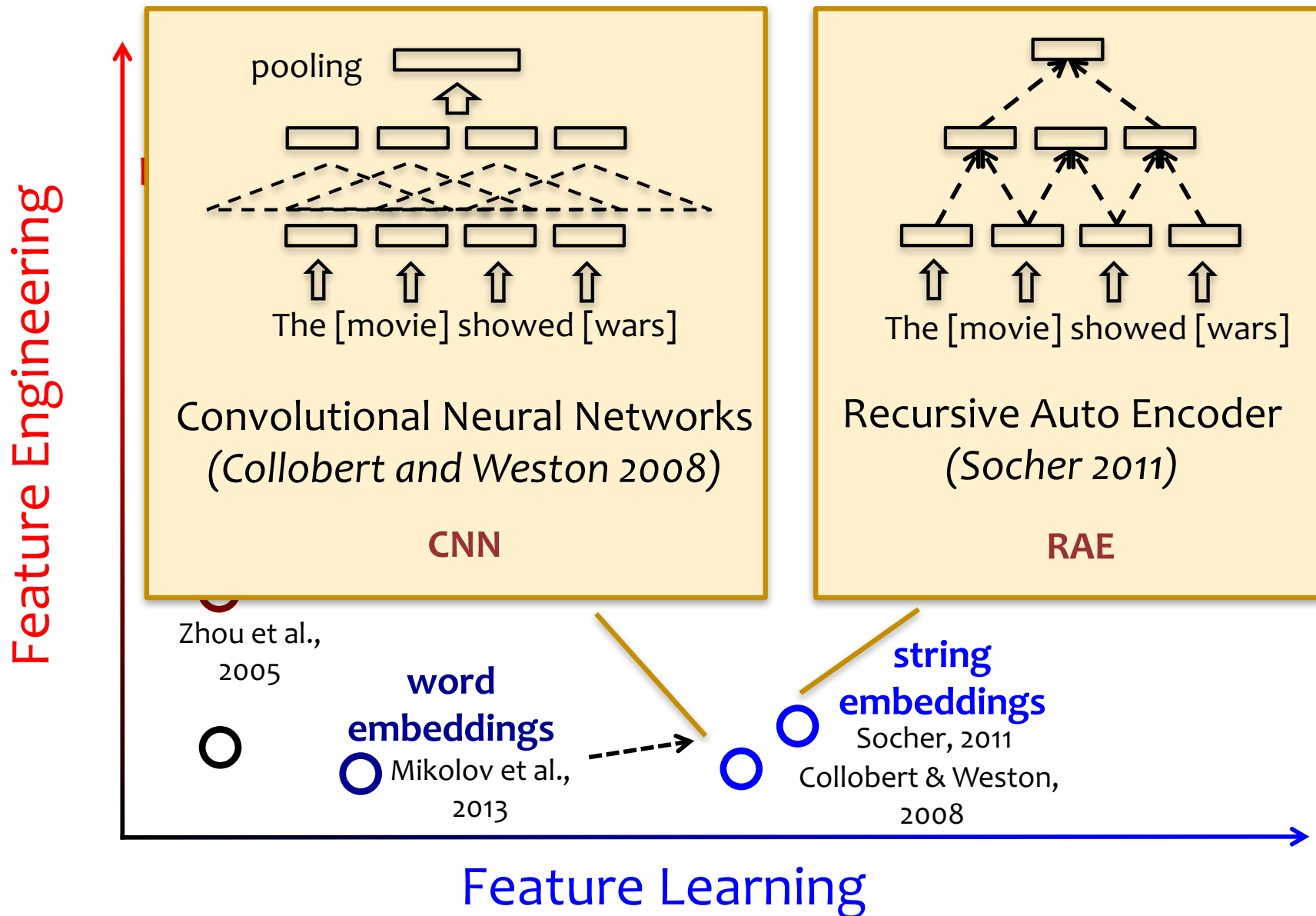
Where do features come from?



Where do features come from?



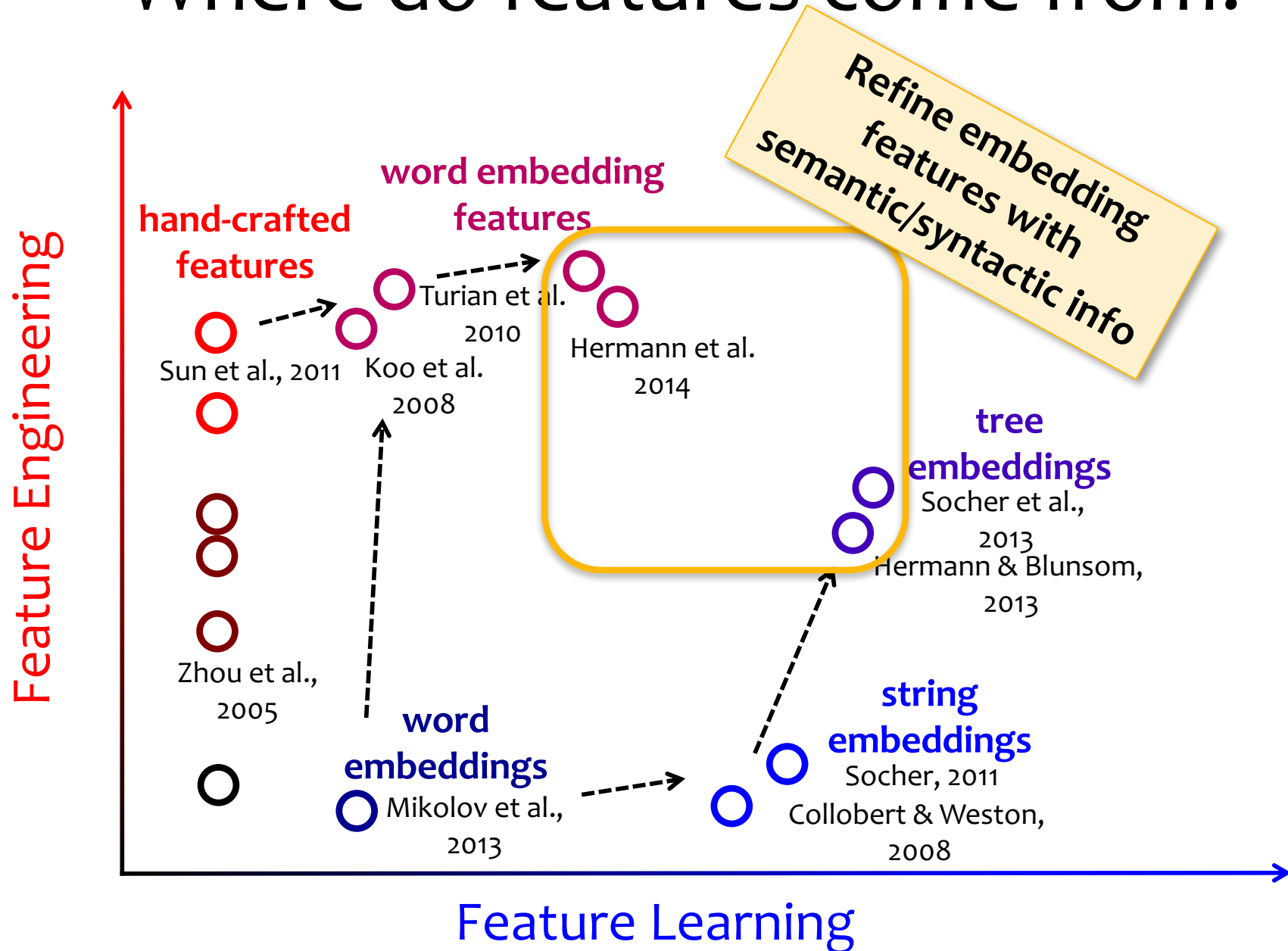
Where do features come from?



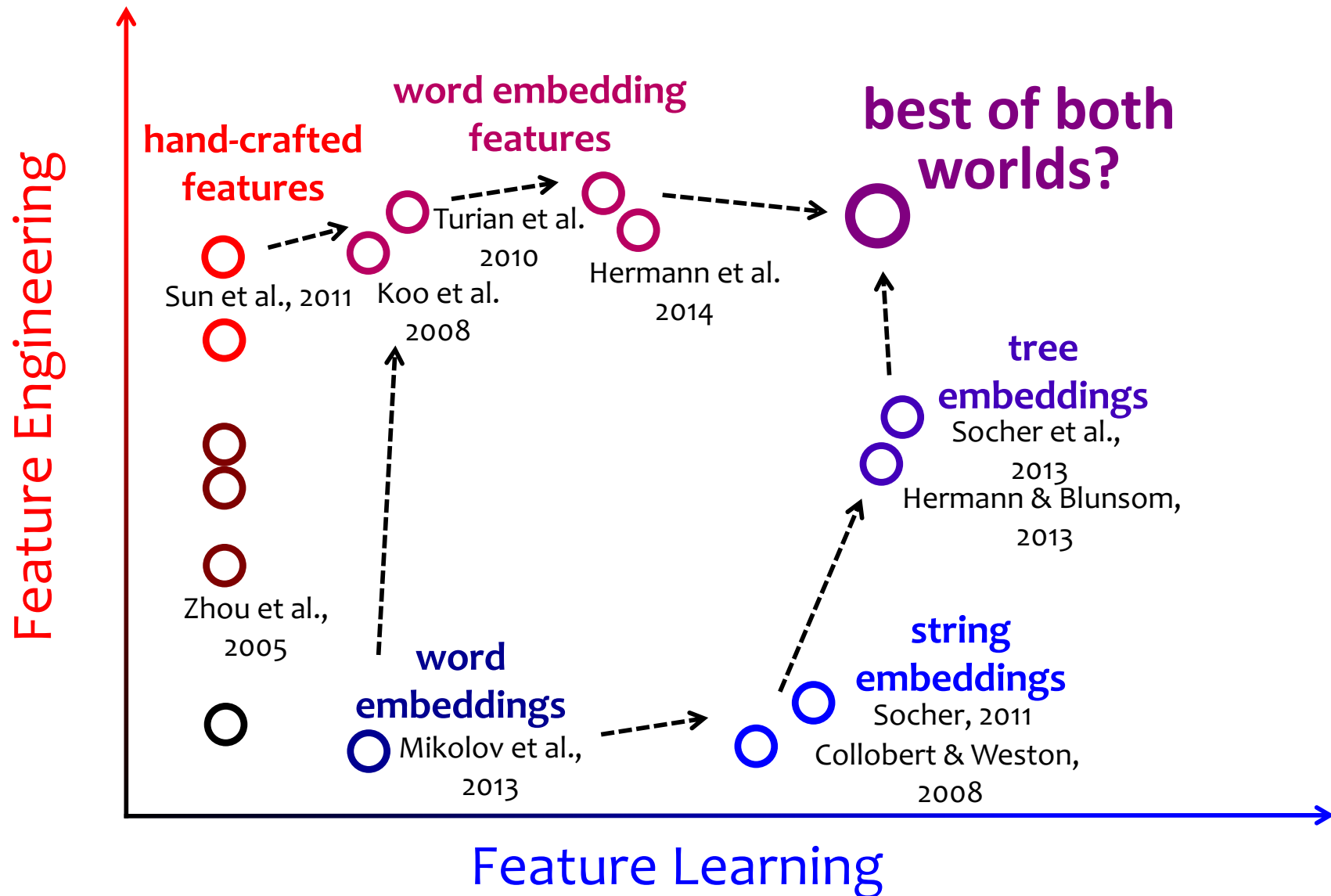
Feature Engineering



Where do features come from?



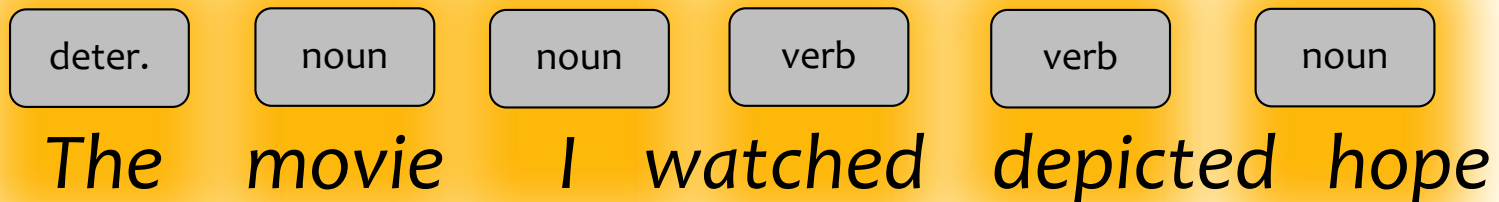
Where do features come from?



Feature Engineering for NLP

Suppose you build a logistic regression model to predict a part-of-speech (POS) tag for each word in a sentence.

What features should you use?



Feature Engineering for NLP

Per-word Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
<code>is-capital(w_i)</code>	1	0	1	0	0	0
<code>endswith(w_i, "e")</code>	1	1	0	0	0	1
<code>endswith(w_i, "d")</code>	0	0	0	1	1	0
<code>endswith(w_i, "ed")</code>	0	0	0	1	1	0
<code>$w_i == \text{"aardvark"}$</code>	0	0	0	0	0	0
<code>$w_i == \text{"hope"}$</code>	0	0	0	0	0	1
...

deter. noun noun verb verb noun

The movie I watched depicted hope

Feature Engineering for NLP

Context Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
...
$w_i == \text{"watched"}$	0	0	0	1	0	0
$w_{i+1} == \text{"watched"}$	0	0	1	0	0	0
$w_{i-1} == \text{"watched"}$	0	0	0	0	1	0
$w_{i+2} == \text{"watched"}$	0	1	0	0	0	0
$w_{i-2} == \text{"watched"}$	0	0	0	0	0	1
...

deter.	noun	noun	verb	verb	noun
<i>The</i>	<i>movie</i>	<i>I</i>	<i>watched</i>	<i>depicted</i>	<i>hope</i>

Feature Engineering for NLP

Context Features:

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$	$x^{(6)}$
...
$w_i == \text{"I"}$	0	0	1	0	0	0
$w_{i+1} == \text{"I"}$	0	1	0	0	0	0
$w_{i-1} == \text{"I"}$	0	0	0	1	0	0
$w_{i+2} == \text{"I"}$	1	0	0	0	0	0
$w_{i-2} == \text{"I"}$	0	0	0	0	1	0
...

deter.	noun	noun	verb	verb	noun
<i>The</i>	<i>movie</i>	<i>I</i>	<i>watched</i>	<i>depicted</i>	<i>hope</i>

Feature Engineering for NLP

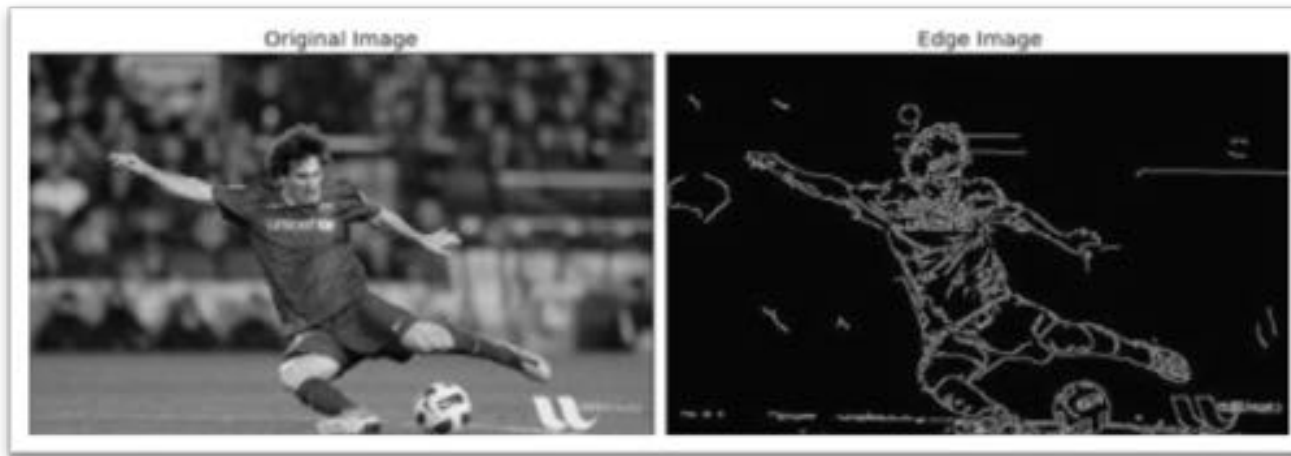
Table 3. Tagging accuracies with different feature templates and other changes on the *WSJ* 19-21 development set.

Model	Feature Templates	# Feats	Sent. Acc.	Token Acc.	Unk. Acc.
3GRAMMEMM	See text	248,798	52.07%	96.92%	88.99%
NAACL 2003	See text and [1]	460,552	55.31%	97.15%	88.61%
Replication	See text and [1]	460,551	55.62%	97.18%	88.92%
Replication'	+rareFeatureThresh = 5	482,364	55.67%	97.19%	88.96%
5W	+ $\langle t_0, w_{-2} \rangle, \langle t_0, w_2 \rangle$	730,178	56.23%	97.20%	89.03%
5WSHAPES	+ $\langle t_0, s_{-1} \rangle, \langle t_0, s_0 \rangle, \langle t_0, s_{+1} \rangle$	731,661	56.52%	97.25%	89.81%
5WSHAPESDS	+ distributional similarity	737,955	56.79%	97.28%	90.46%

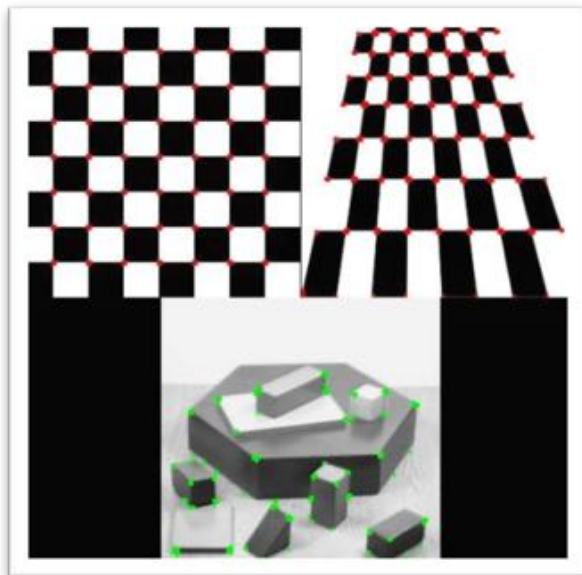
deter. noun noun verb verb noun
The movie I watched depicted hope

Feature Engineering for CV

Edge detection (Canny)

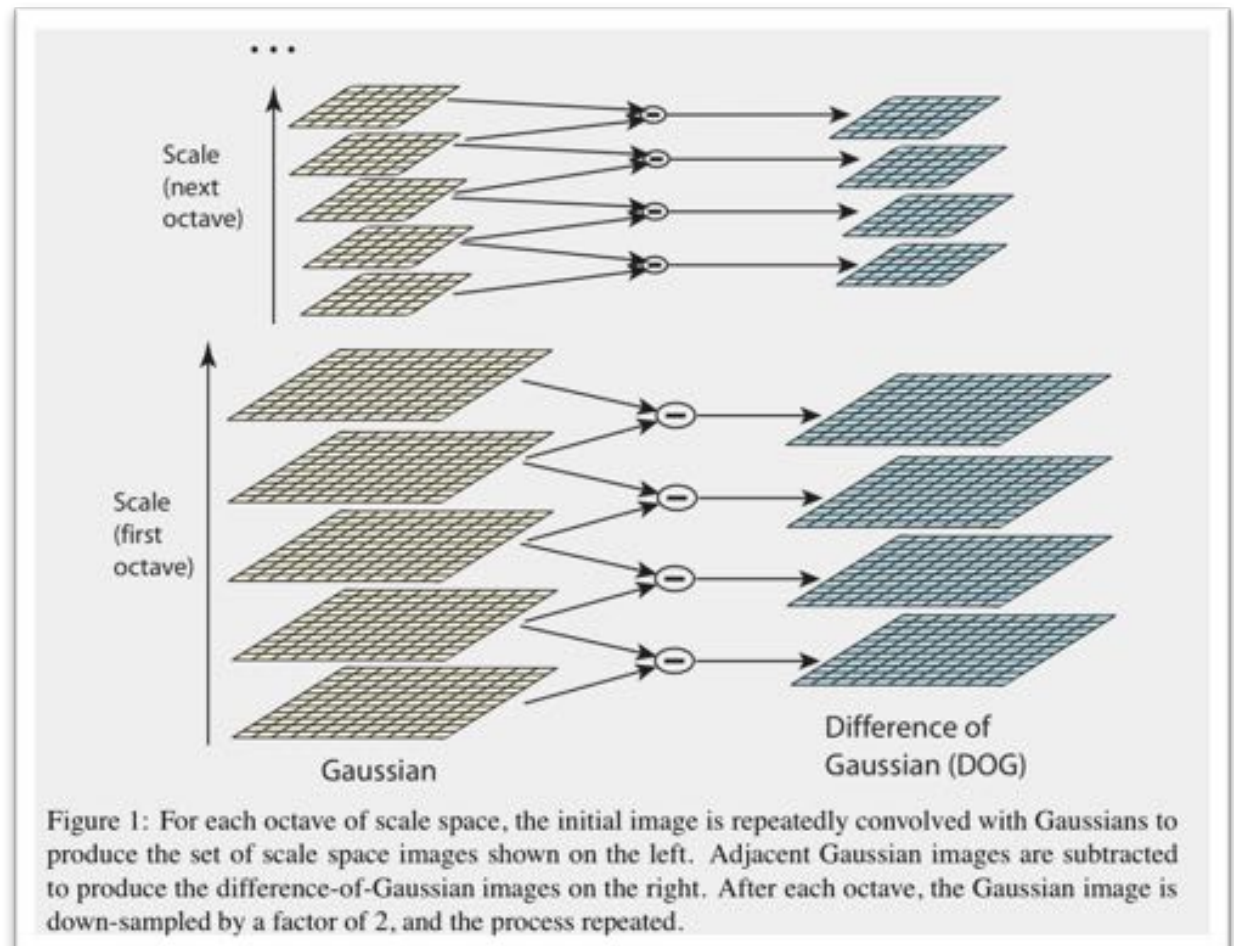


Corner Detection (Harris)



Feature Engineering for CV

Scale Invariant Feature Transform (SIFT)



NON-LINEAR FEATURES

Nonlinear Features

- aka. “nonlinear basis functions”
- So far, input was always $\mathbf{x} = [x_1, \dots, x_M]$
- **Key Idea:** let input be some function of \mathbf{x}
 - original input: $\mathbf{x} \in \mathbb{R}^M$
 - new input: $\mathbf{x}' \in \mathbb{R}^{M'}$ where $M' > M$ (usually)
 - define $\mathbf{x}' = b(\mathbf{x}) = [b_1(\mathbf{x}), b_2(\mathbf{x}), \dots, b_{M'}(\mathbf{x})]$
where $b_i : \mathbb{R}^M \rightarrow \mathbb{R}$ is any function

- **Examples:** ($M = 1$)

polynomial

$$b_j(x) = x^j \quad \forall j \in \{1, \dots, J\}$$

radial basis function

$$b_j(x) = \exp\left(\frac{-(x - \mu_j)^2}{2\sigma_j^2}\right)$$

sigmoid

$$b_j(x) = \frac{1}{1 + \exp(-\omega_j x)}$$

log

$$b_j(x) = \log(x)$$

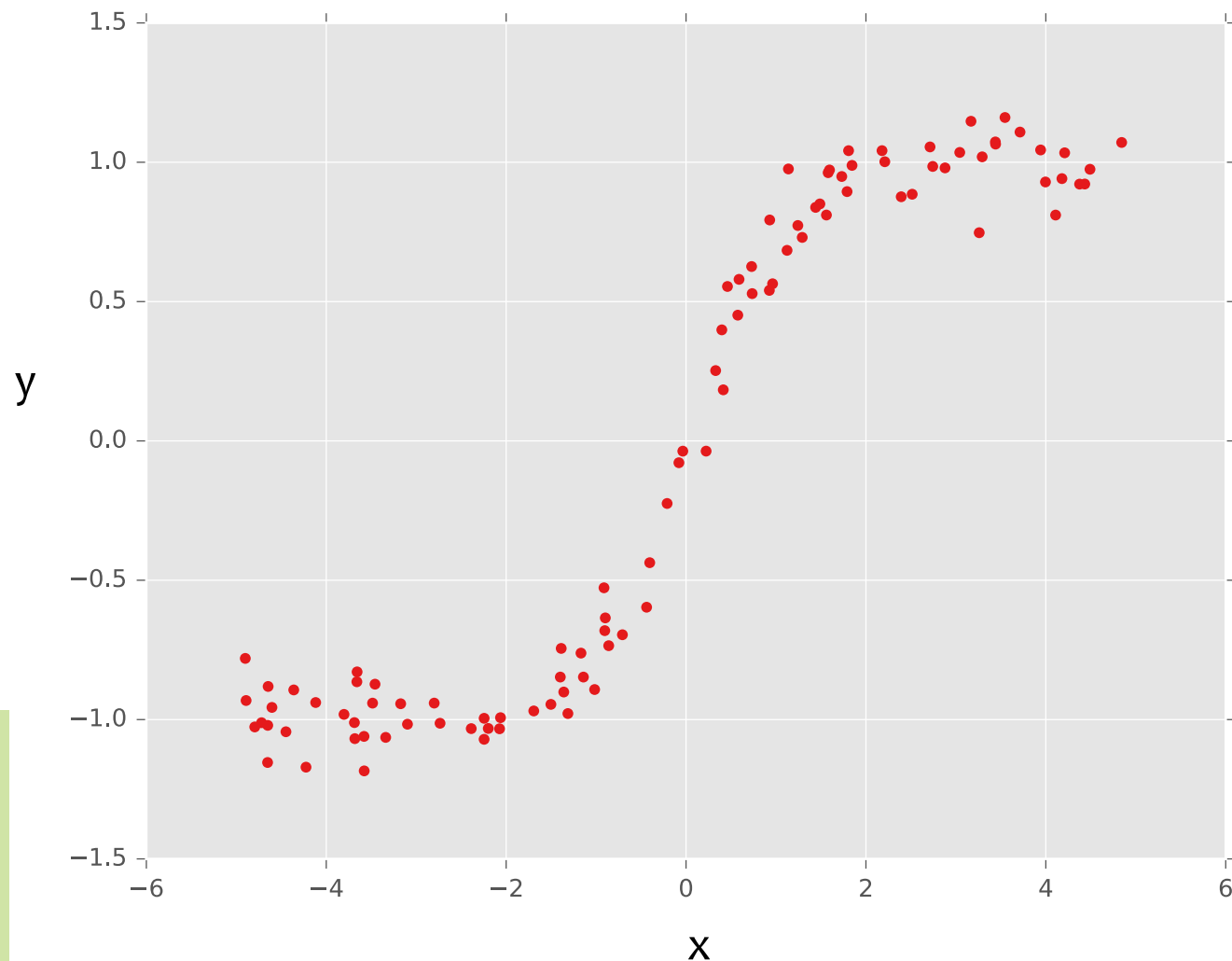
For a linear model:
still a linear function
of $b(\mathbf{x})$ even though a
nonlinear function of
 \mathbf{x}

Examples:

- Perceptron
- Linear regression
- Logistic regression

Example: Linear Regression

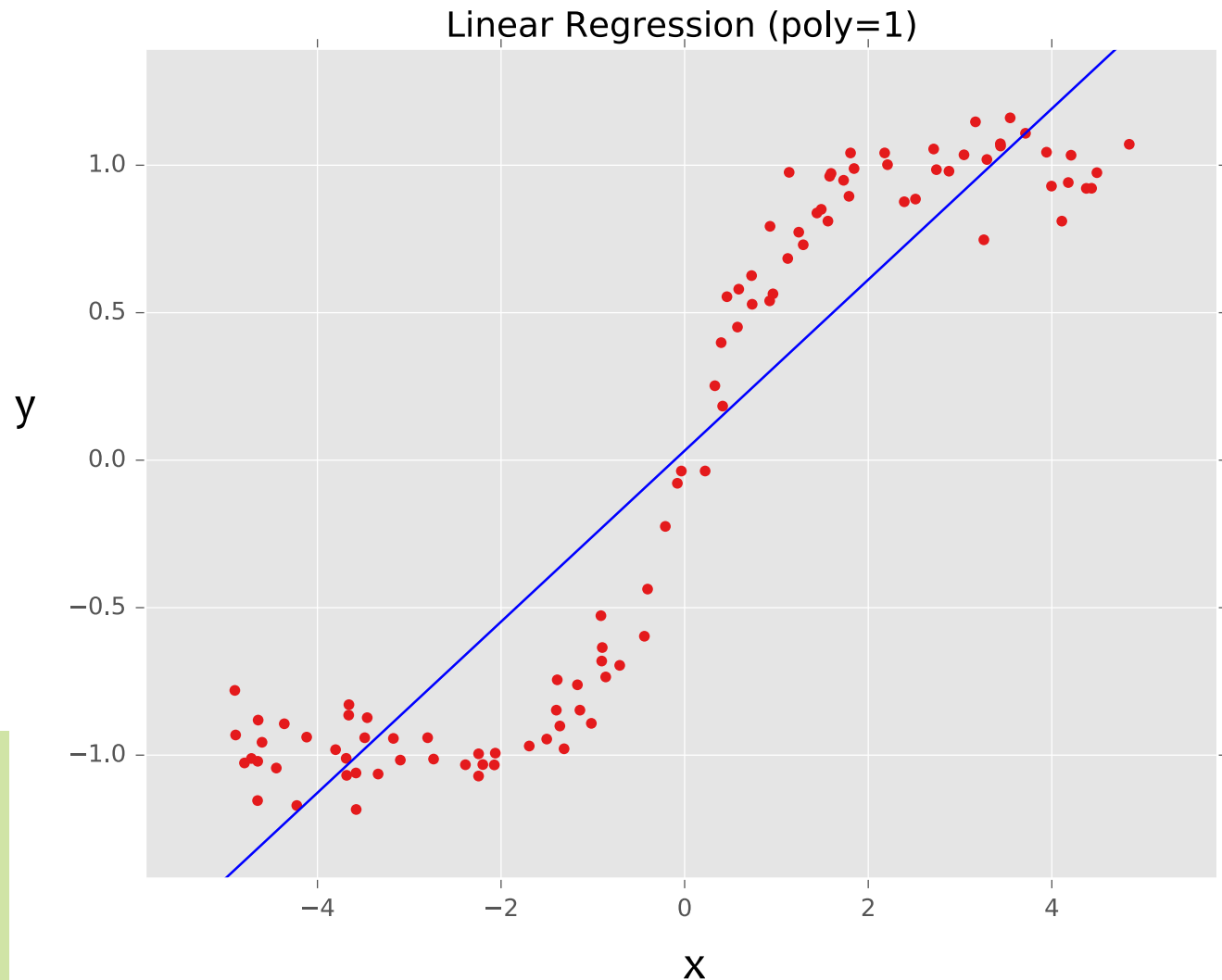
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

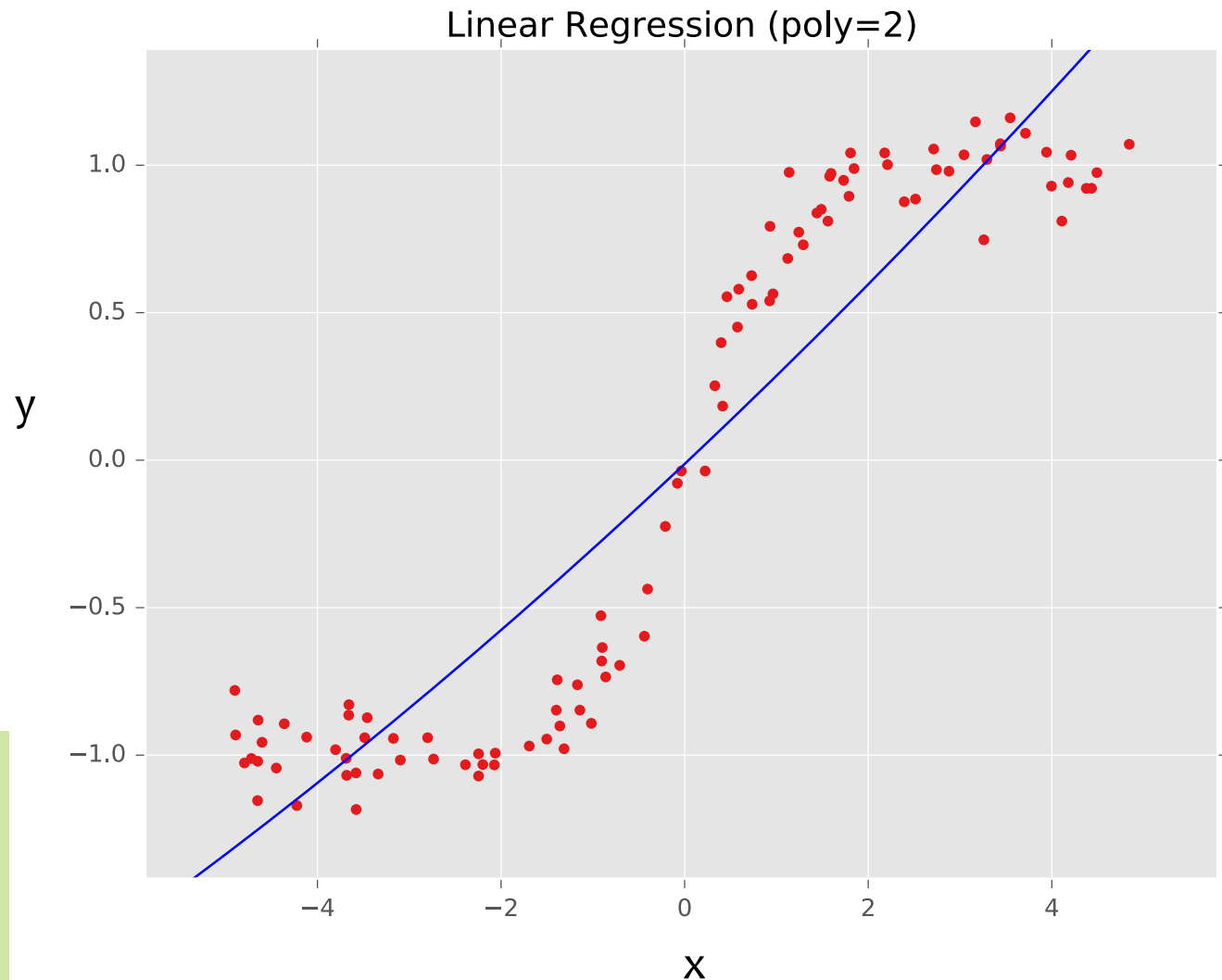
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

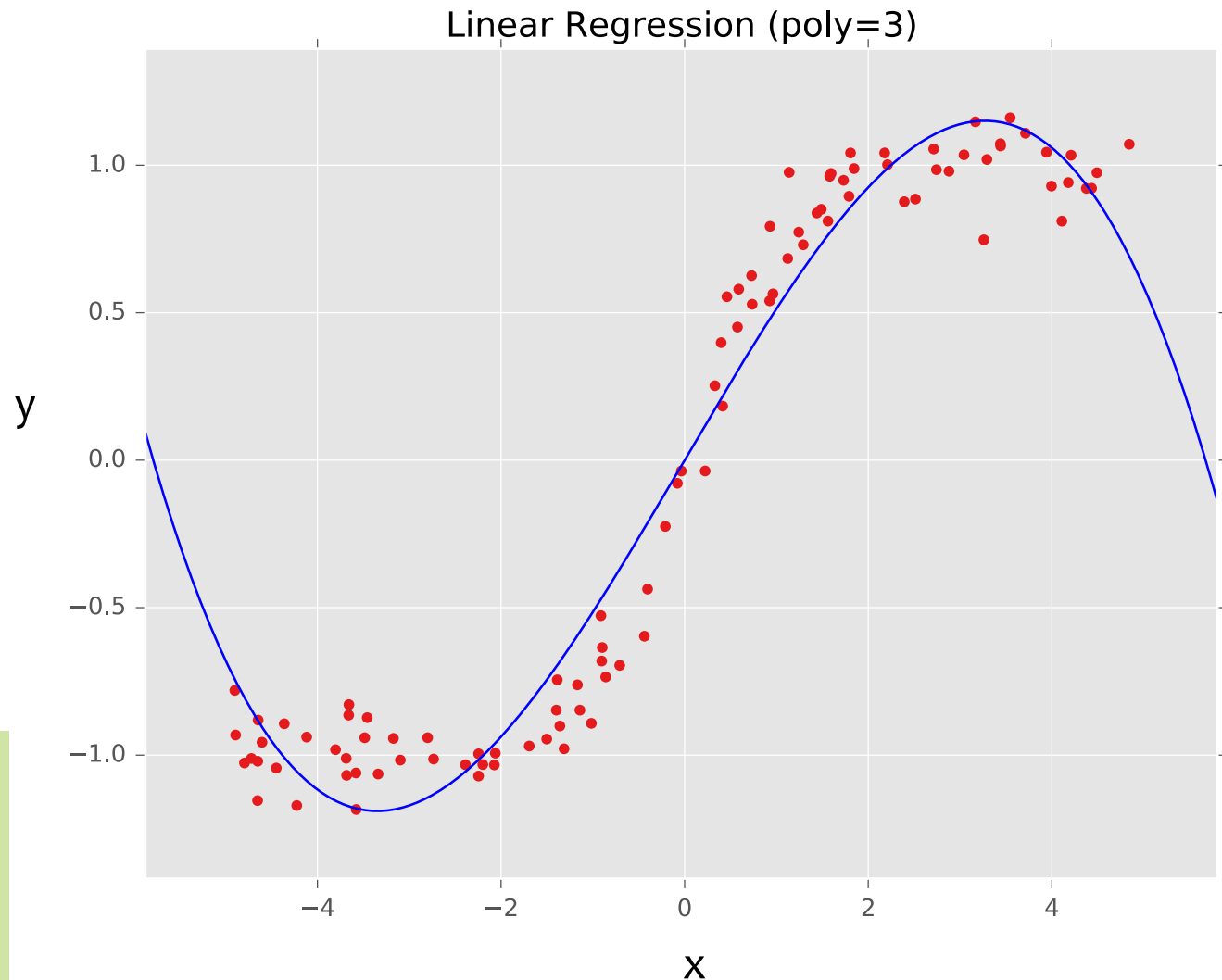
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

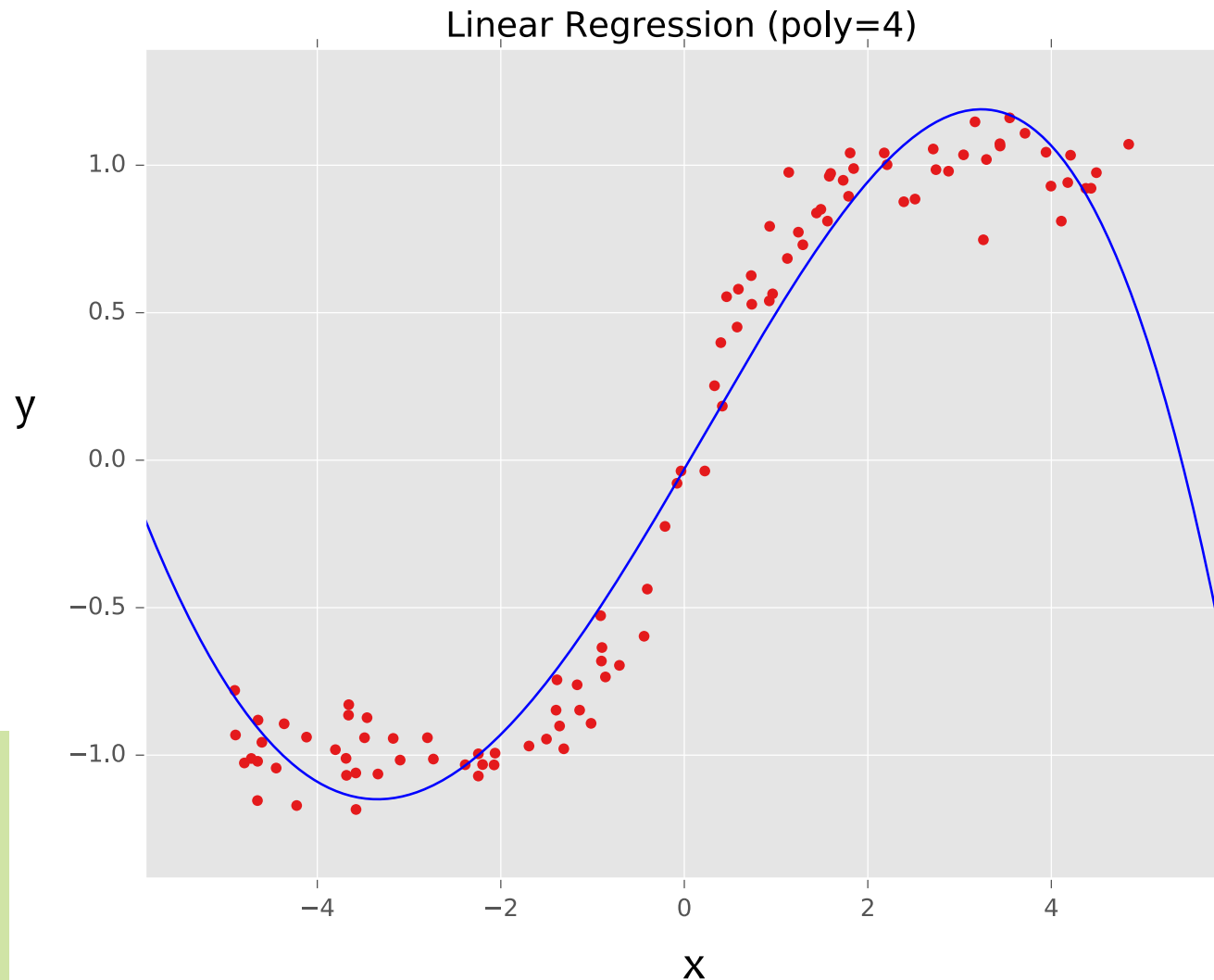
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

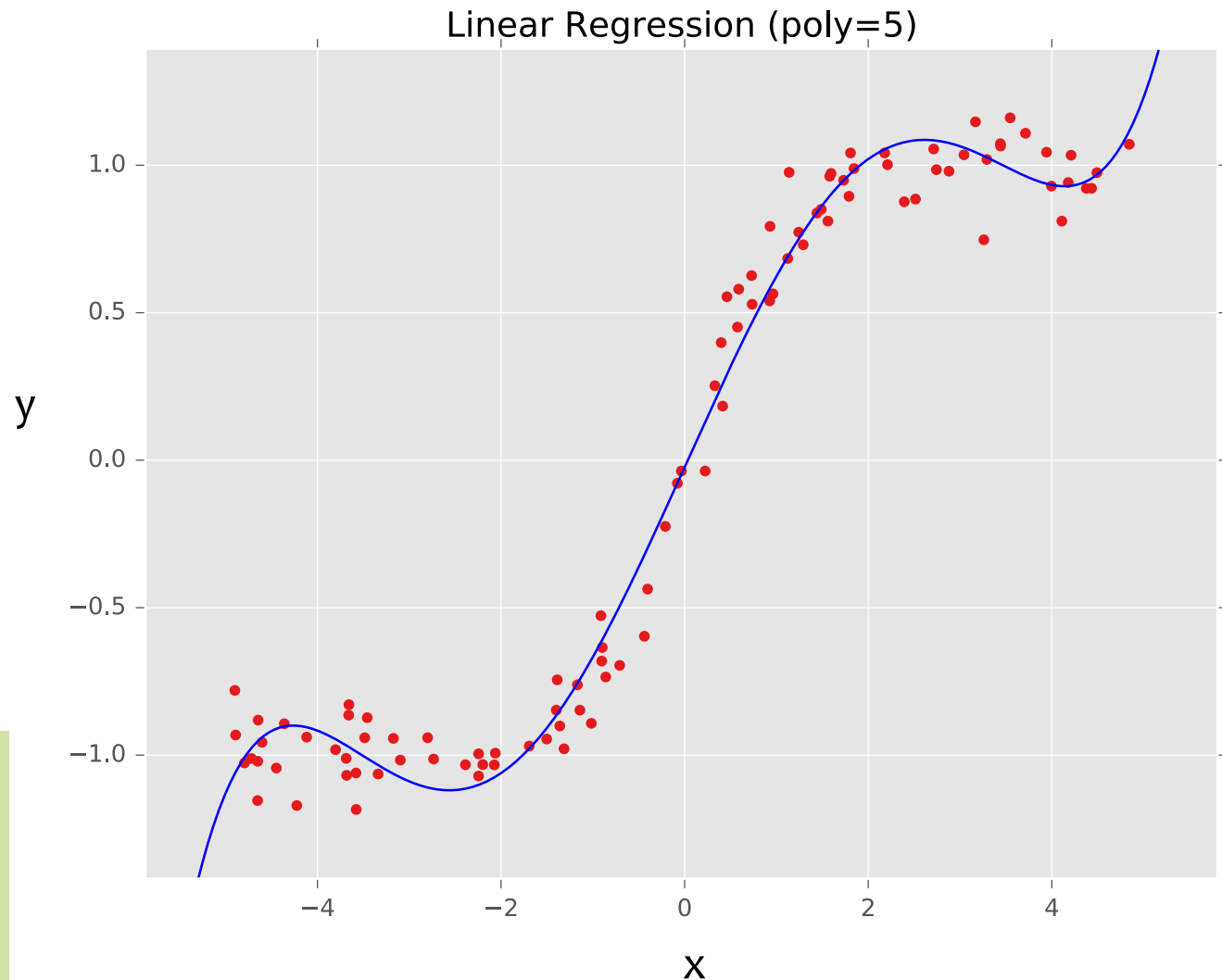
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

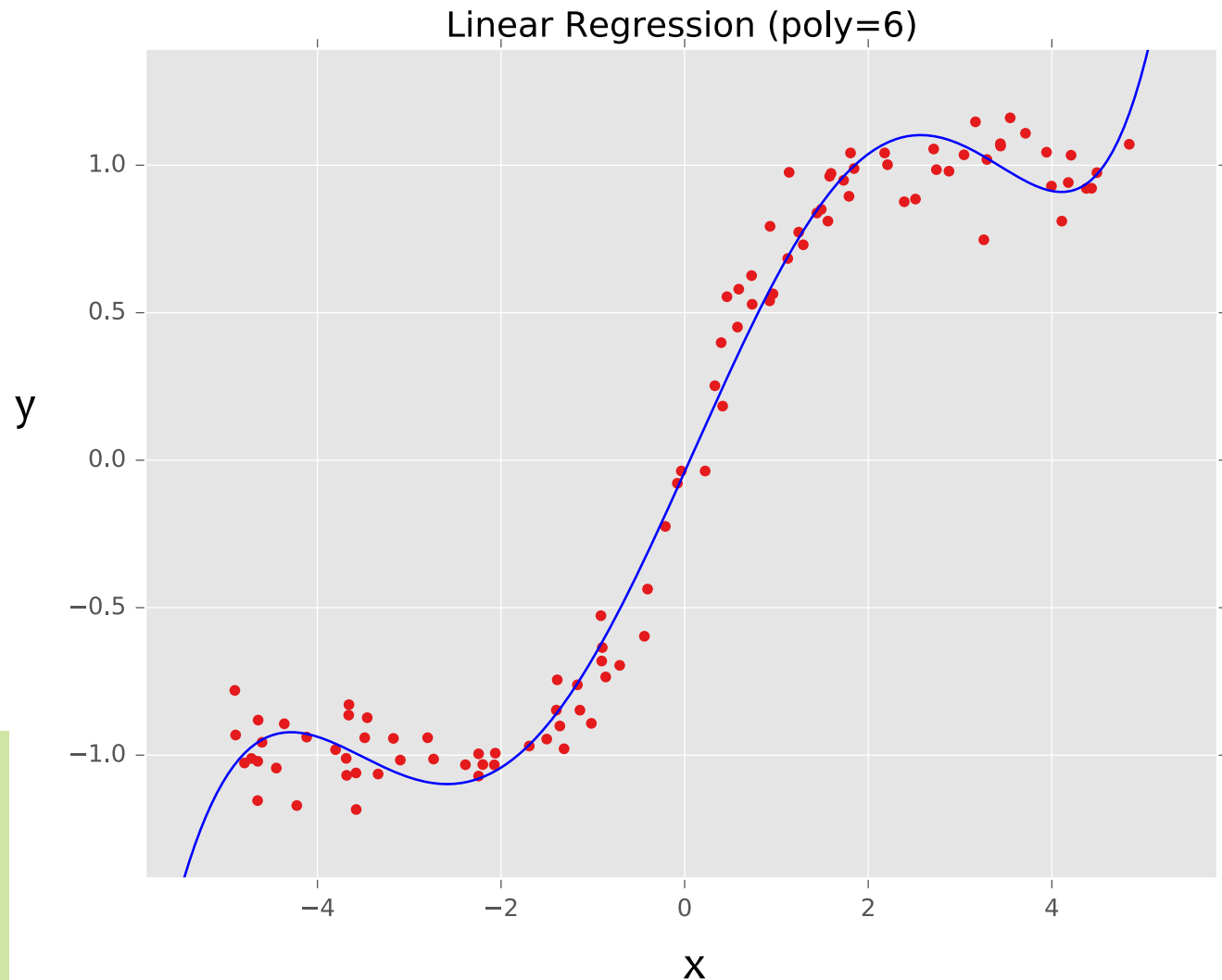
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

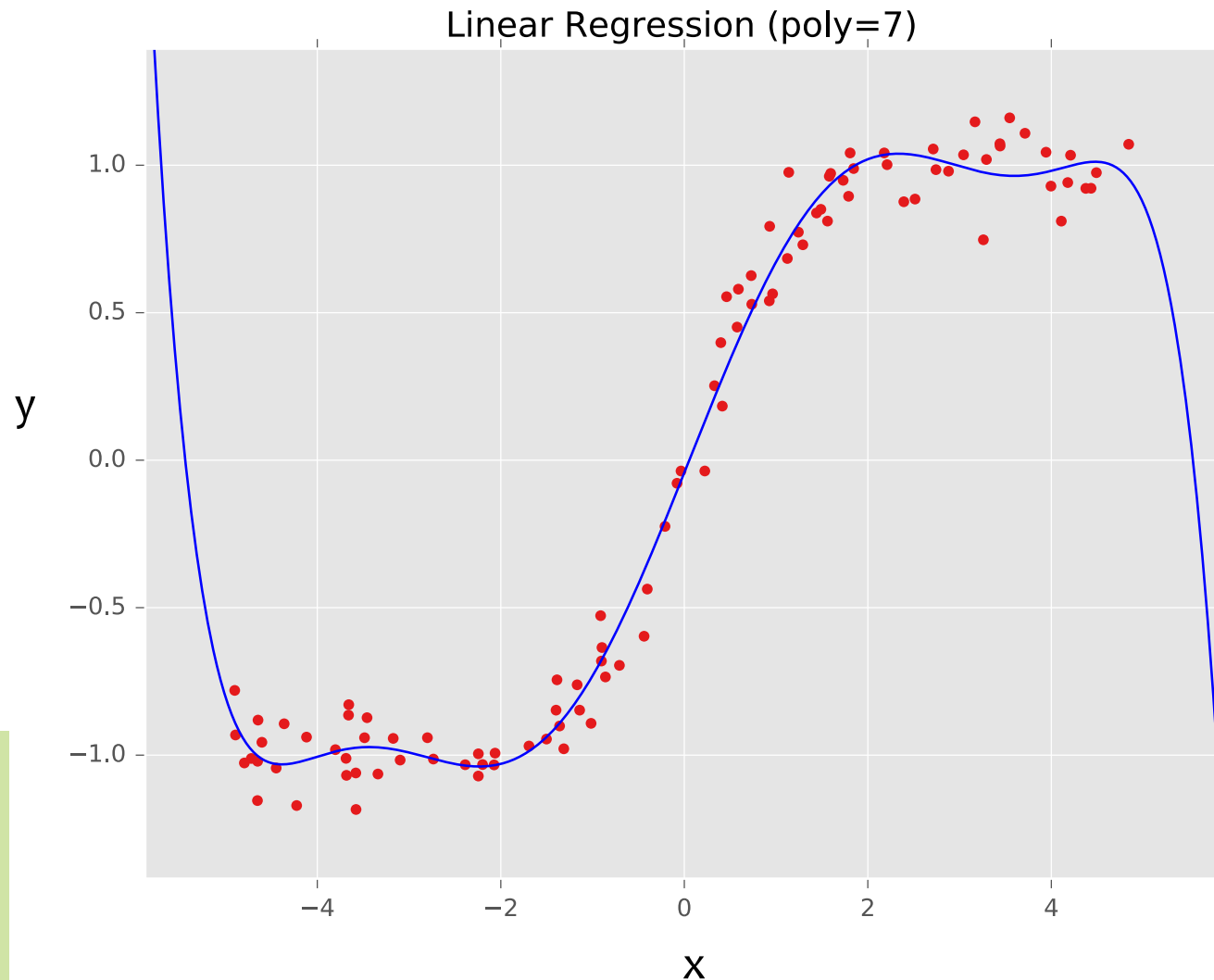
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

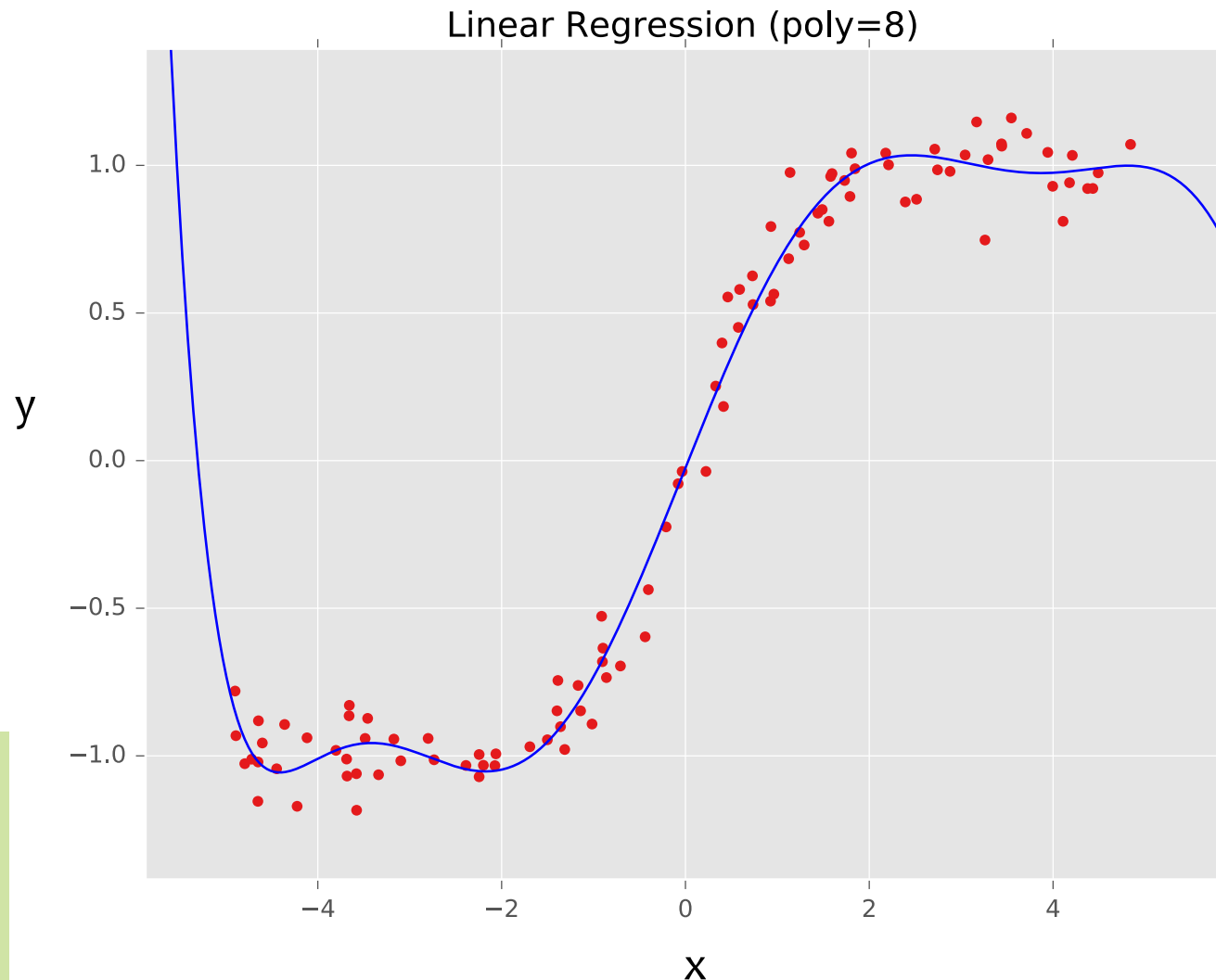
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

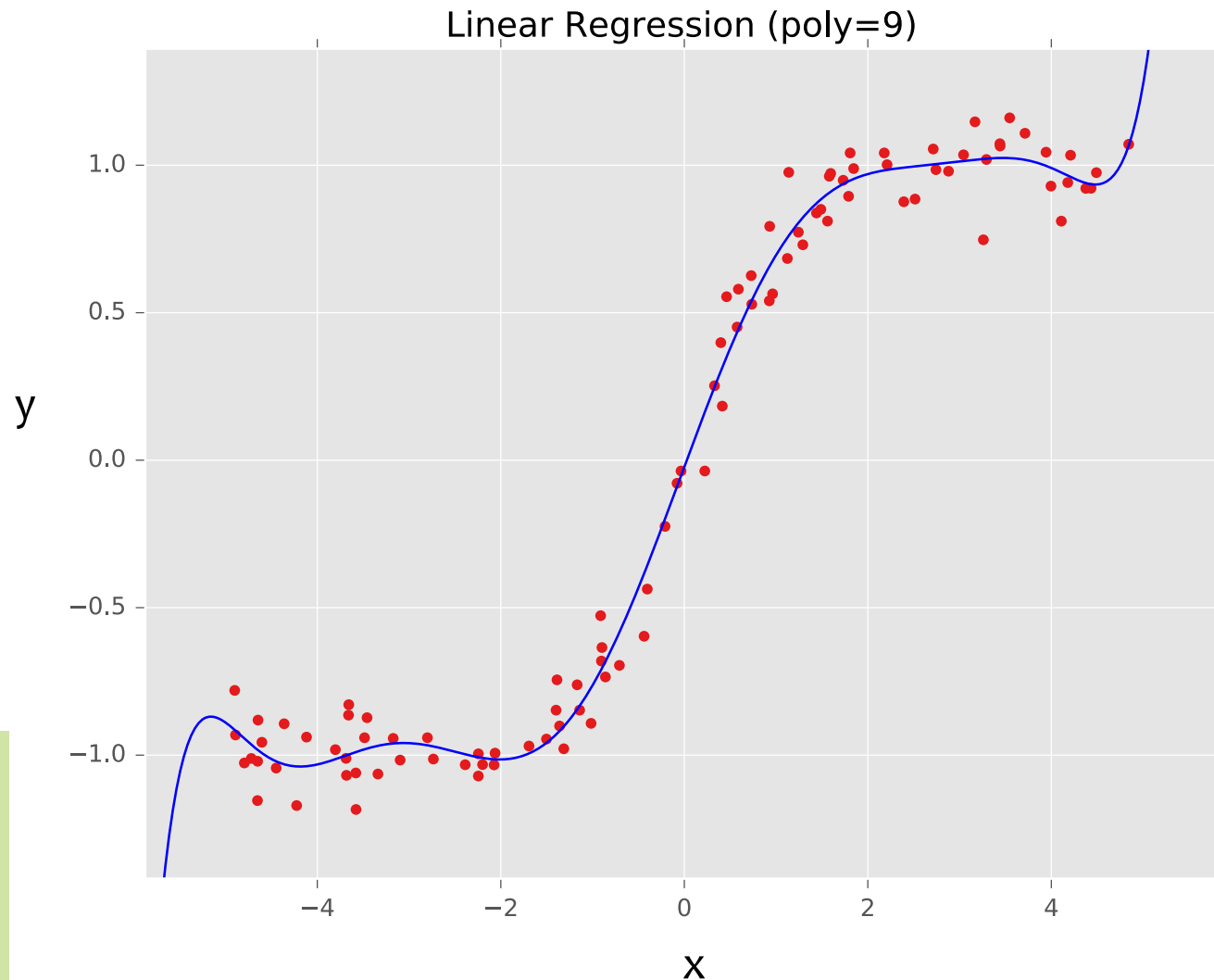
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

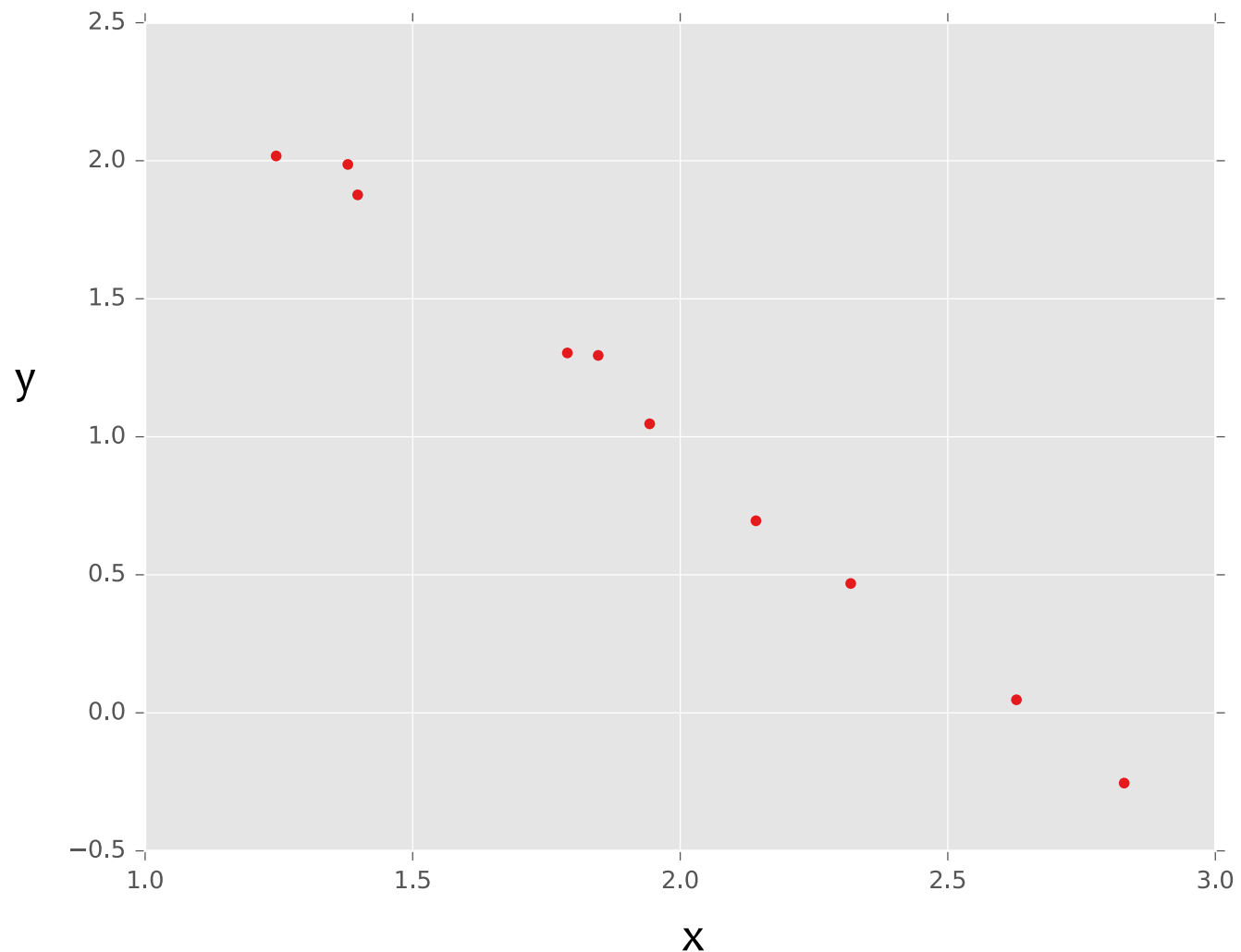
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
 $y = \tanh(x) + \text{noise}$

Example: Linear Regression

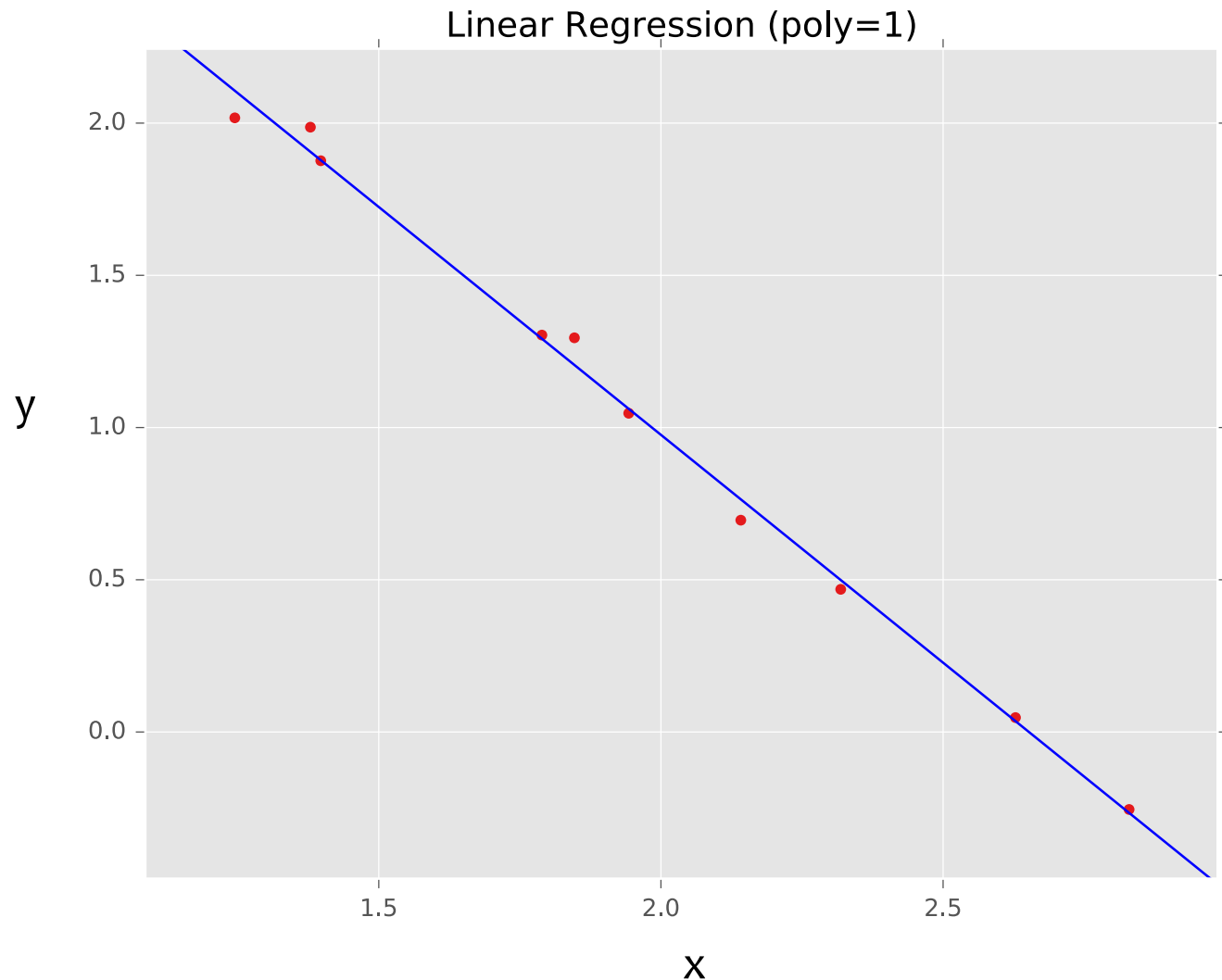
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

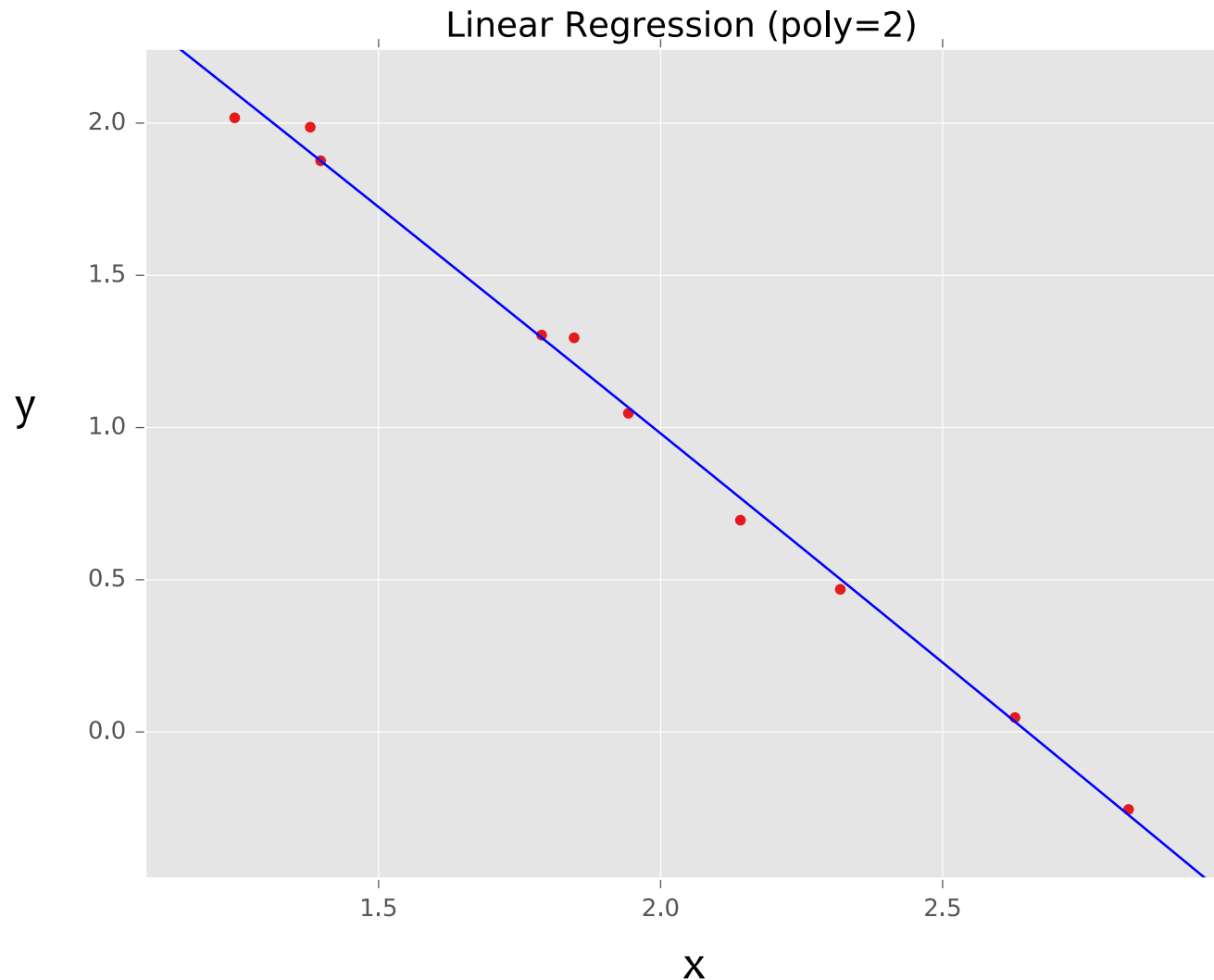
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

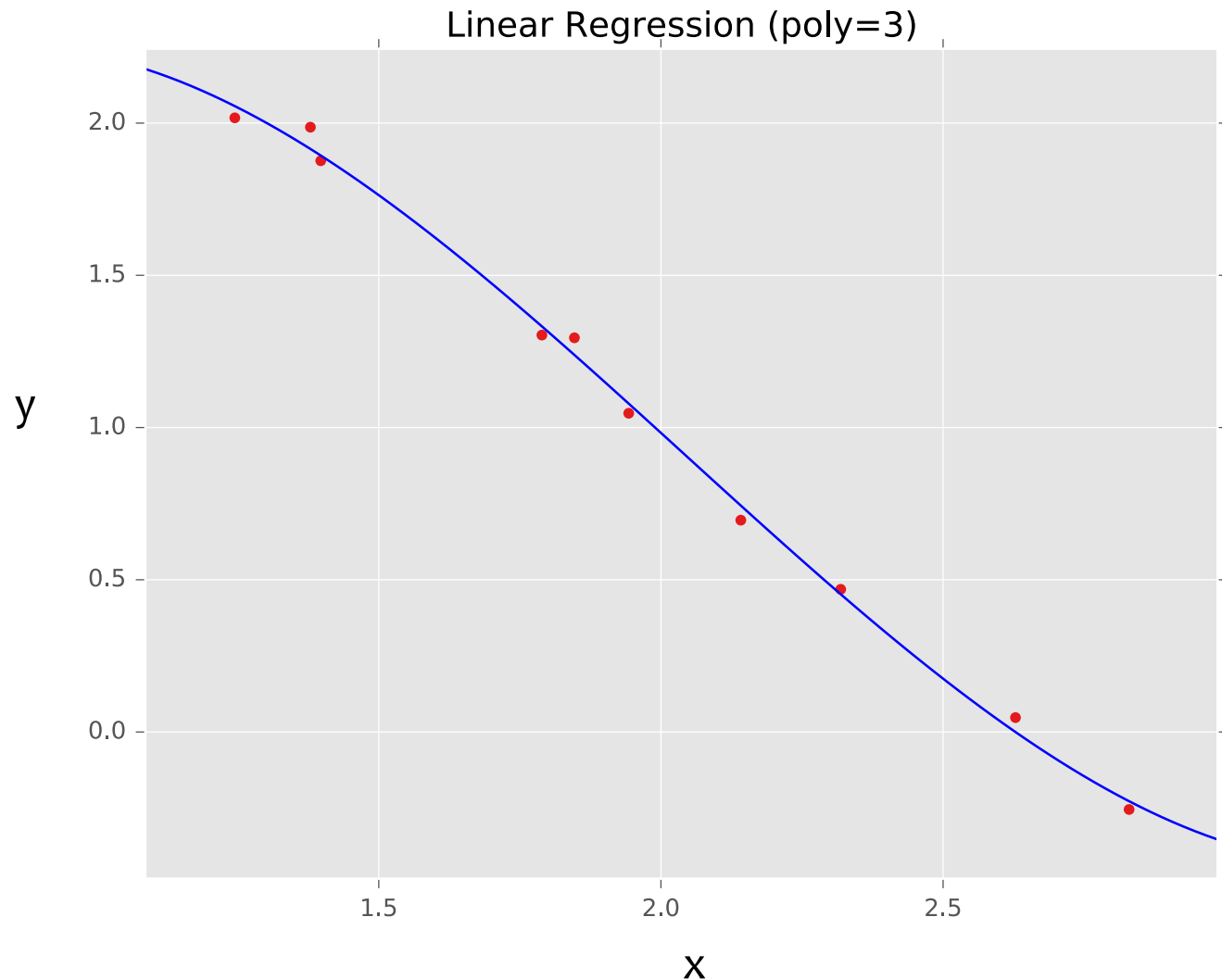
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

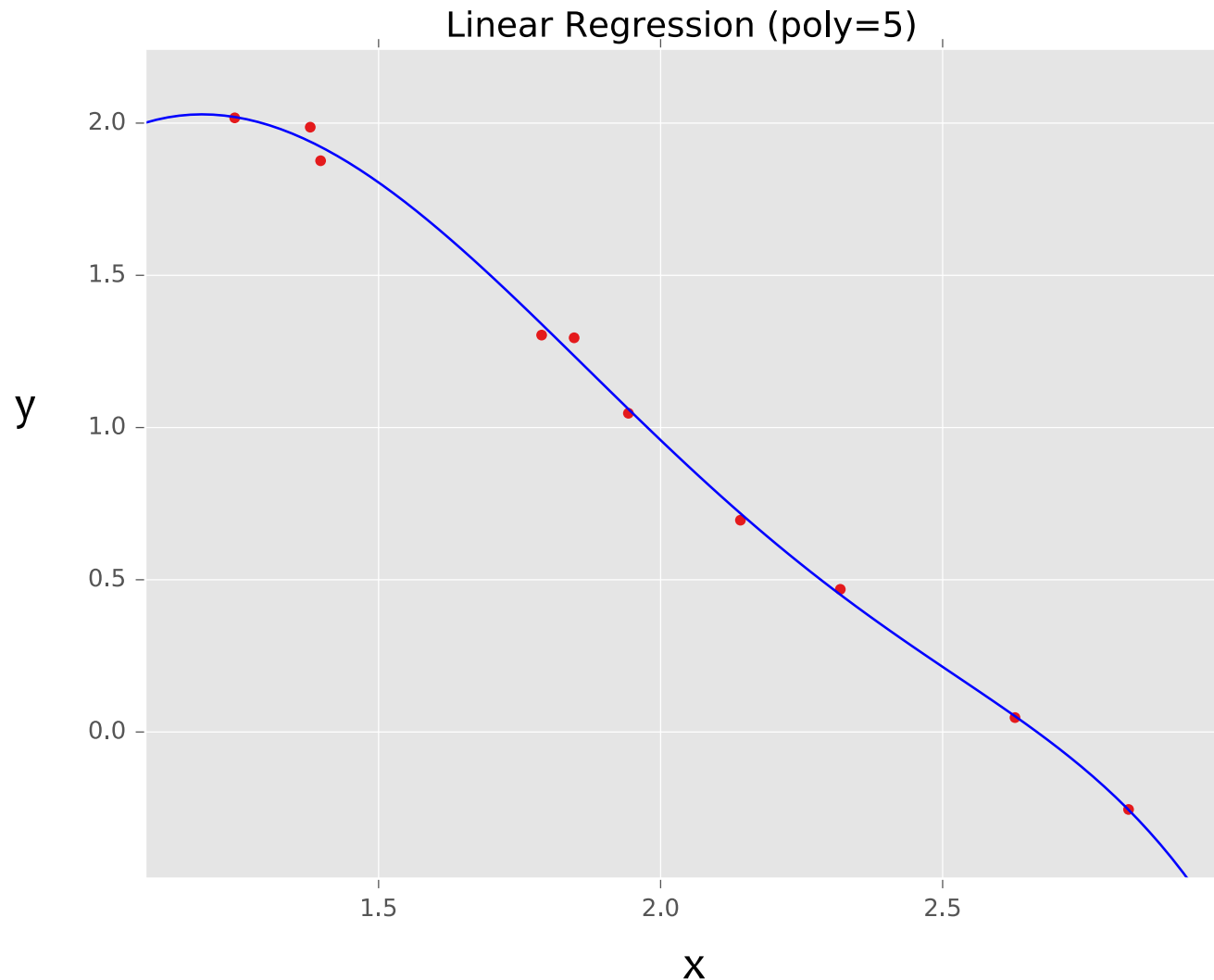
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

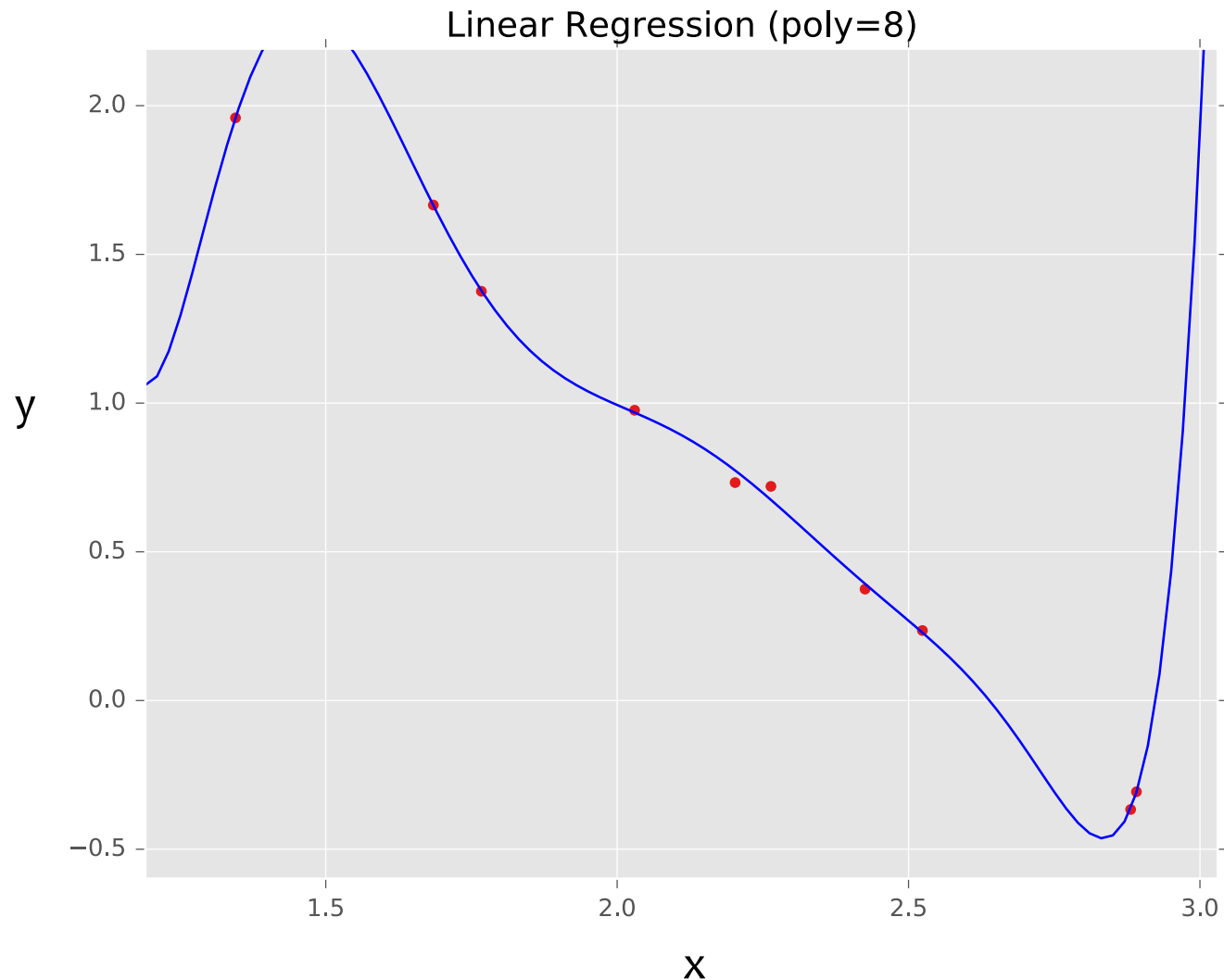
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

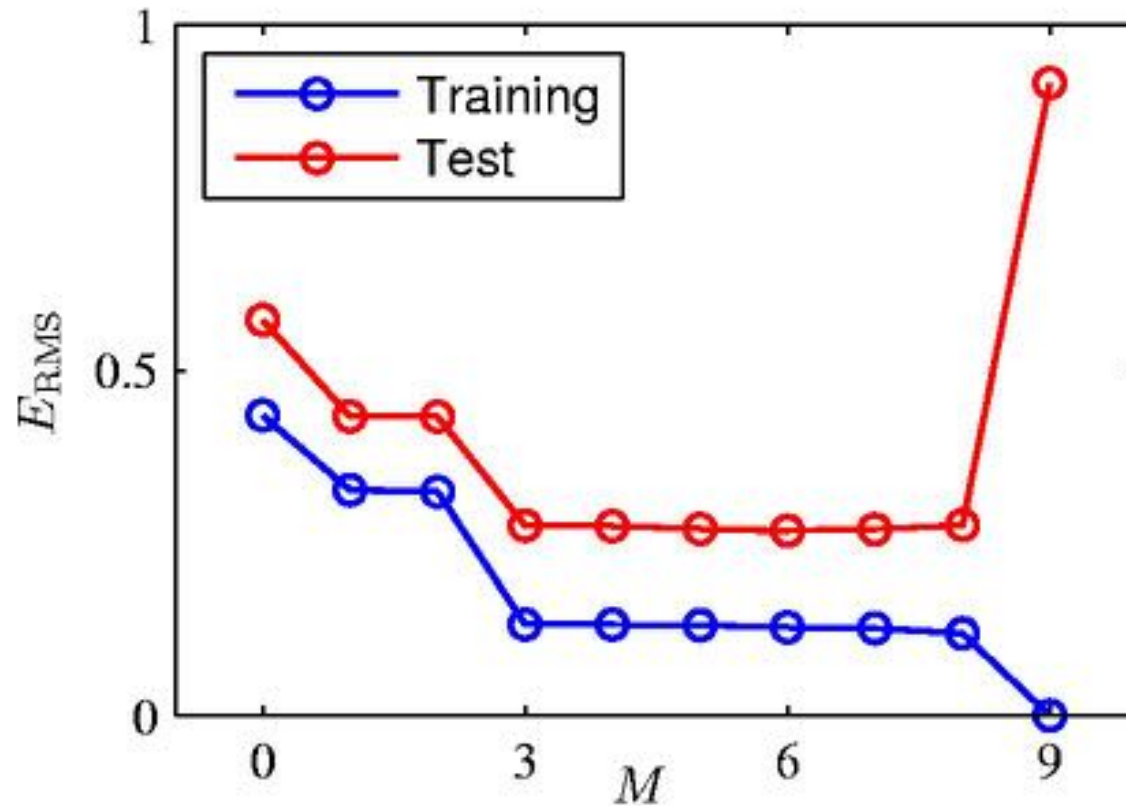
Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Over-fitting



Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
θ_0	0.19	0.82	0.31	0.35
θ_1		-1.27	7.99	232.37
θ_2			-25.43	-5321.83
θ_3			17.37	48568.31
θ_4				-231639.30
θ_5				640042.26
θ_6				-1061800.52
θ_7				1042400.18
θ_8				-557682.99
θ_9				125201.43

Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function

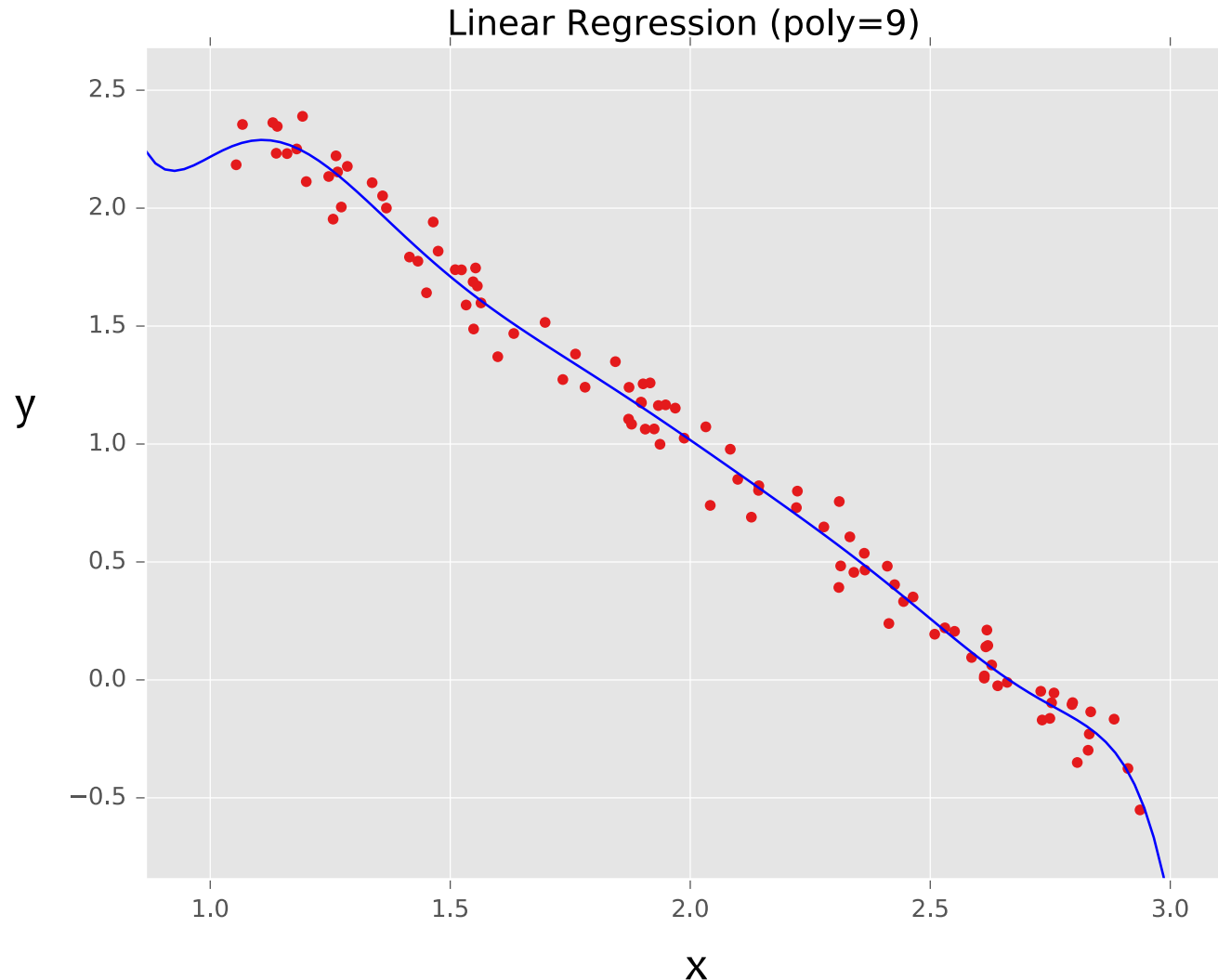


true “unknown”
target function is
linear with
negative slope
and gaussian
noise

Example: Linear Regression

Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function

Same as before, but now
with $N = 100$ points



true “unknown”
target function is
linear with
negative slope
and gaussian
noise

REGULARIZATION

Overfitting

Definition: The problem of **overfitting** is when the model captures the noise in the training data instead of the underlying structure

Overfitting can occur in all the models we've seen so far:

- KNN (e.g. when k is small)
- Naïve Bayes (e.g. without a prior)
- Linear Regression (e.g. with basis function)
- Logistic Regression (e.g. with many rare features)

Motivation: Regularization

Example: Stock Prices

- Suppose we wish to predict Google's stock price at time $t+1$
- **What features should we use?**
(putting all computational concerns aside)
 - Stock prices of all other stocks at times $t, t-1, t-2, \dots, t-k$
 - Mentions of Google with positive / negative sentiment words in all newspapers and social media outlets
- Do we believe that **all** of these features are going to be useful?



Motivation: Regularization

- **Occam's Razor:** prefer the simplest hypothesis
- What does it mean for a hypothesis (or model) to be **simple**?
 1. small number of features (**model selection**)
 2. small number of “important” features (**shrinkage**)

Regularization

Chalkboard

- L2, L1, L0 Regularization
- Example: Linear Regression

Regularization

Don't Regularize the Bias (Intercept) Parameter!

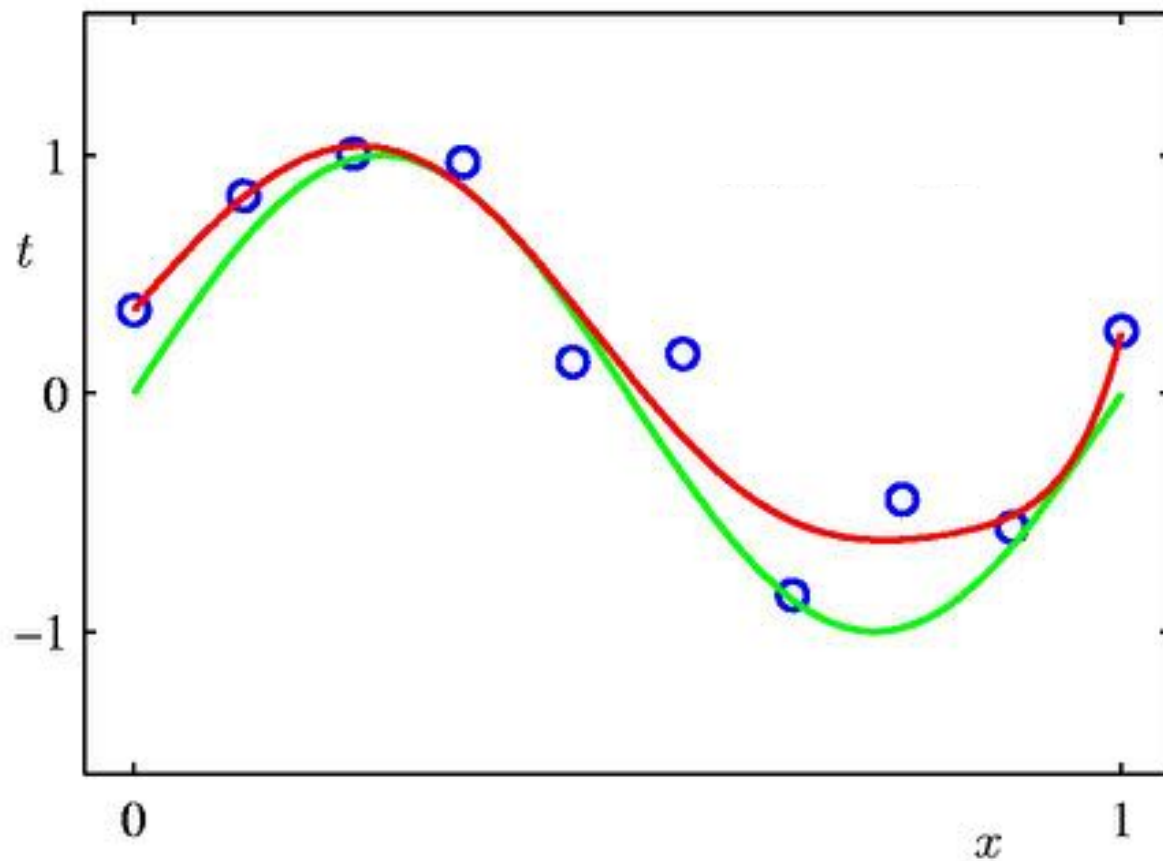
- In our models so far, the bias / intercept parameter is usually denoted by θ_0 -- that is, the parameter for which we fixed $x_0 = 1$
- Regularizers always avoid penalizing this bias / intercept parameter
- Why? Because otherwise the learning algorithms wouldn't be invariant to a shift in the y-values

Whitening Data

- It's common to *whiten* each feature by subtracting its mean and dividing by its variance
- For regularization, this helps all the features be penalized in the same units
(e.g. convert both centimeters and kilometers to z-scores)

Regularization:

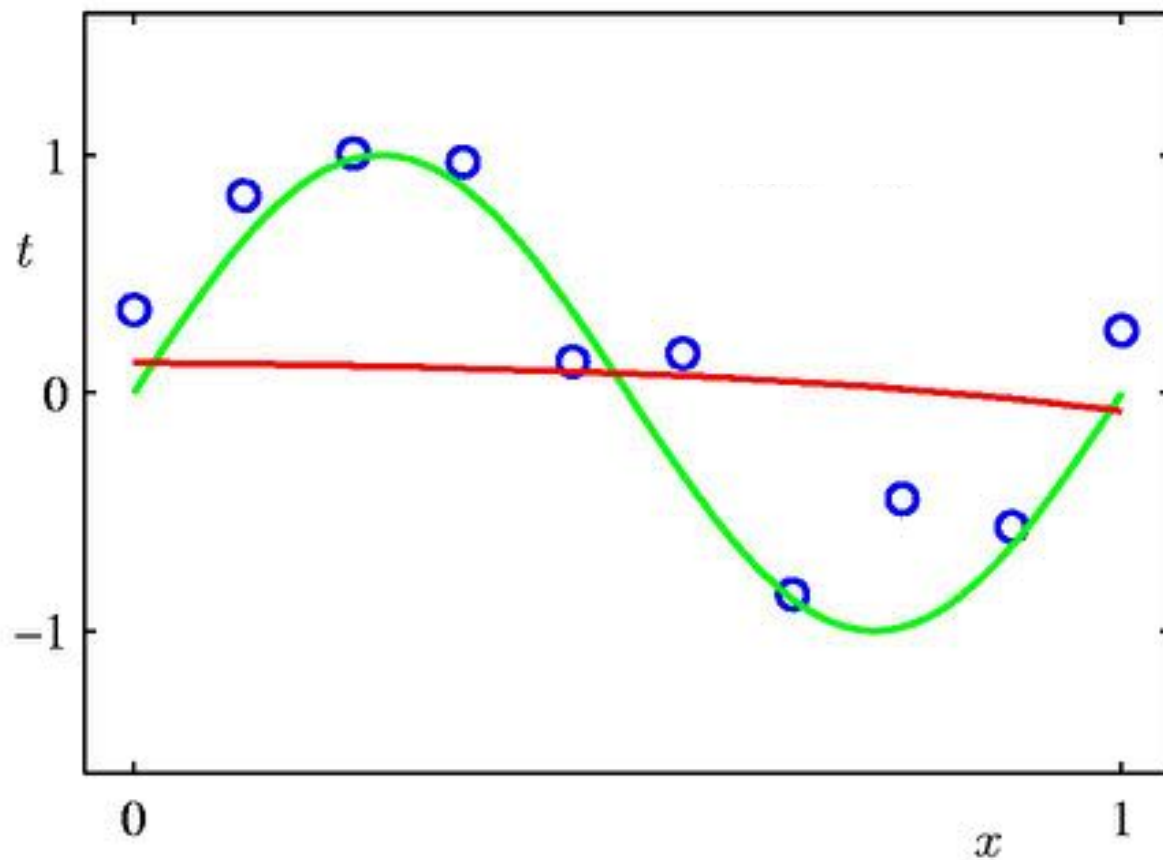
$$\ln \lambda = +.18$$



Polynomial Coefficients

	none	exp(18)	huge
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53	0.00
w_9^*	125201.43	72.68	0.01

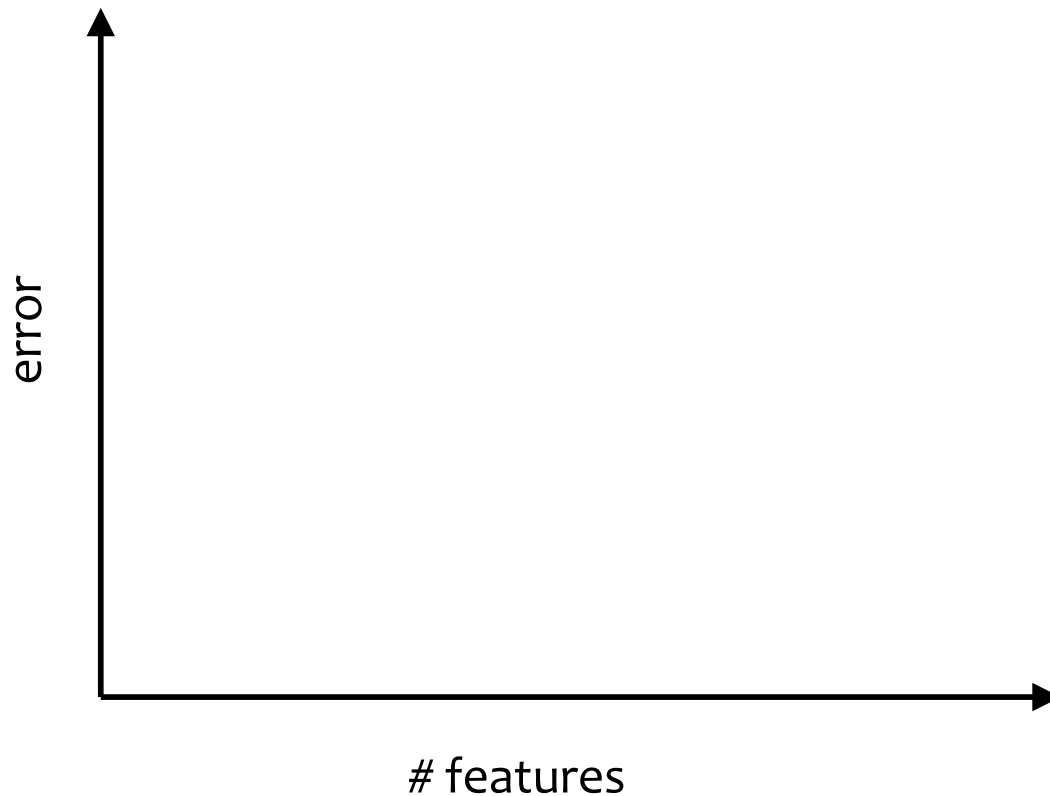
Over Regularization:



Regularization Exercise

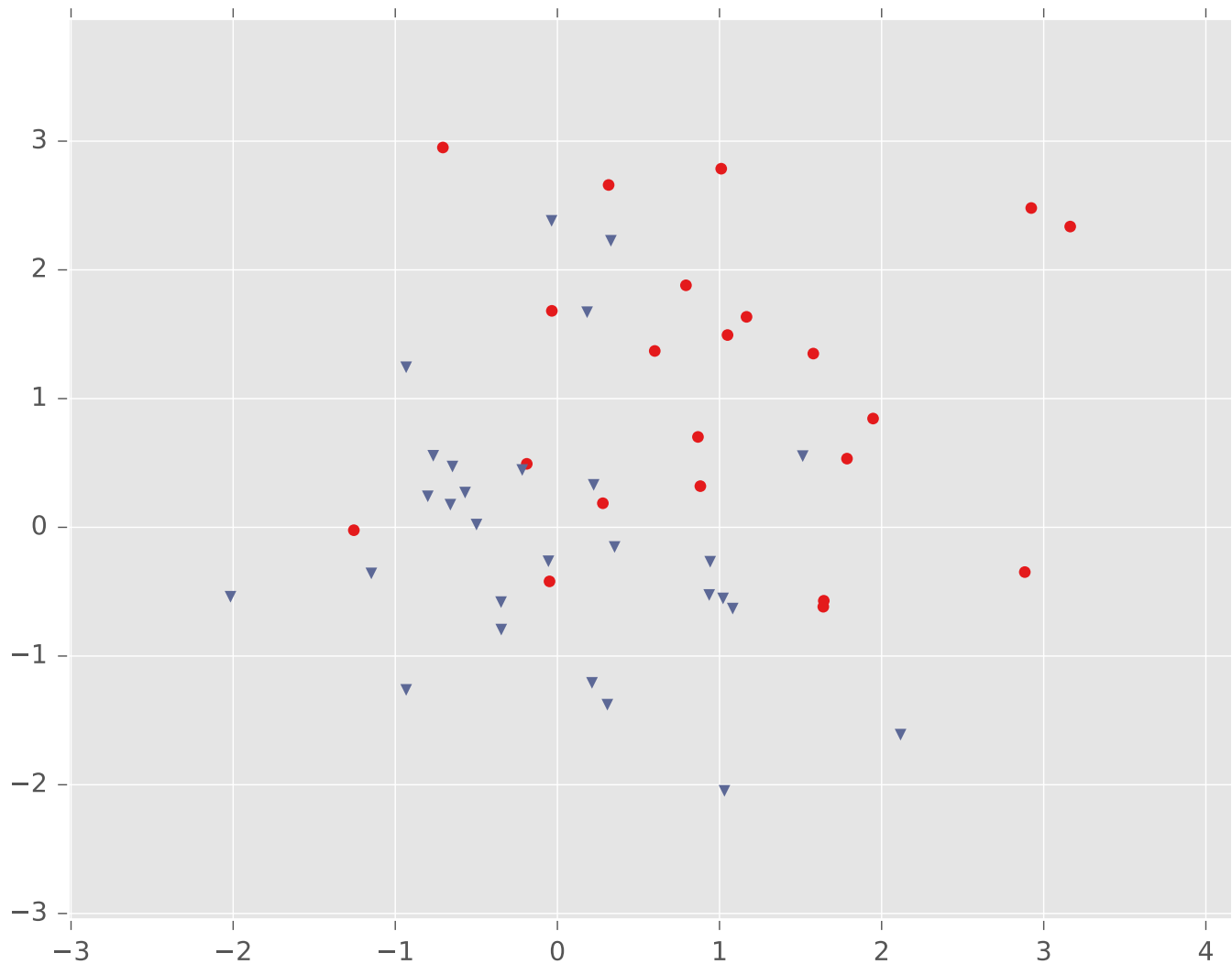
In-class Exercise

1. Plot train error vs. # features (cartoon)
2. Plot test error vs. # features (cartoon)



Example: Logistic Regression

Training
Data

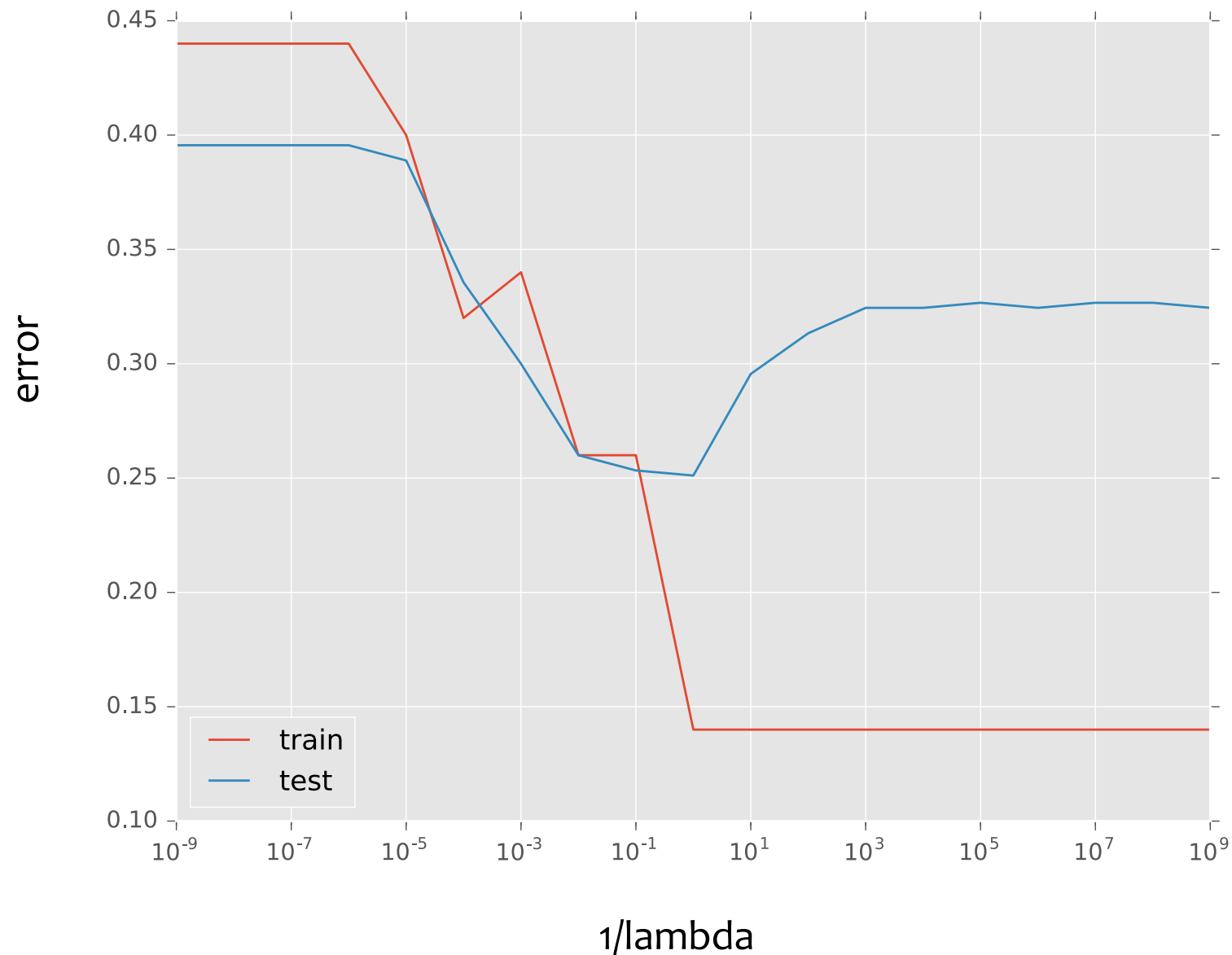


Example: Logistic Regression

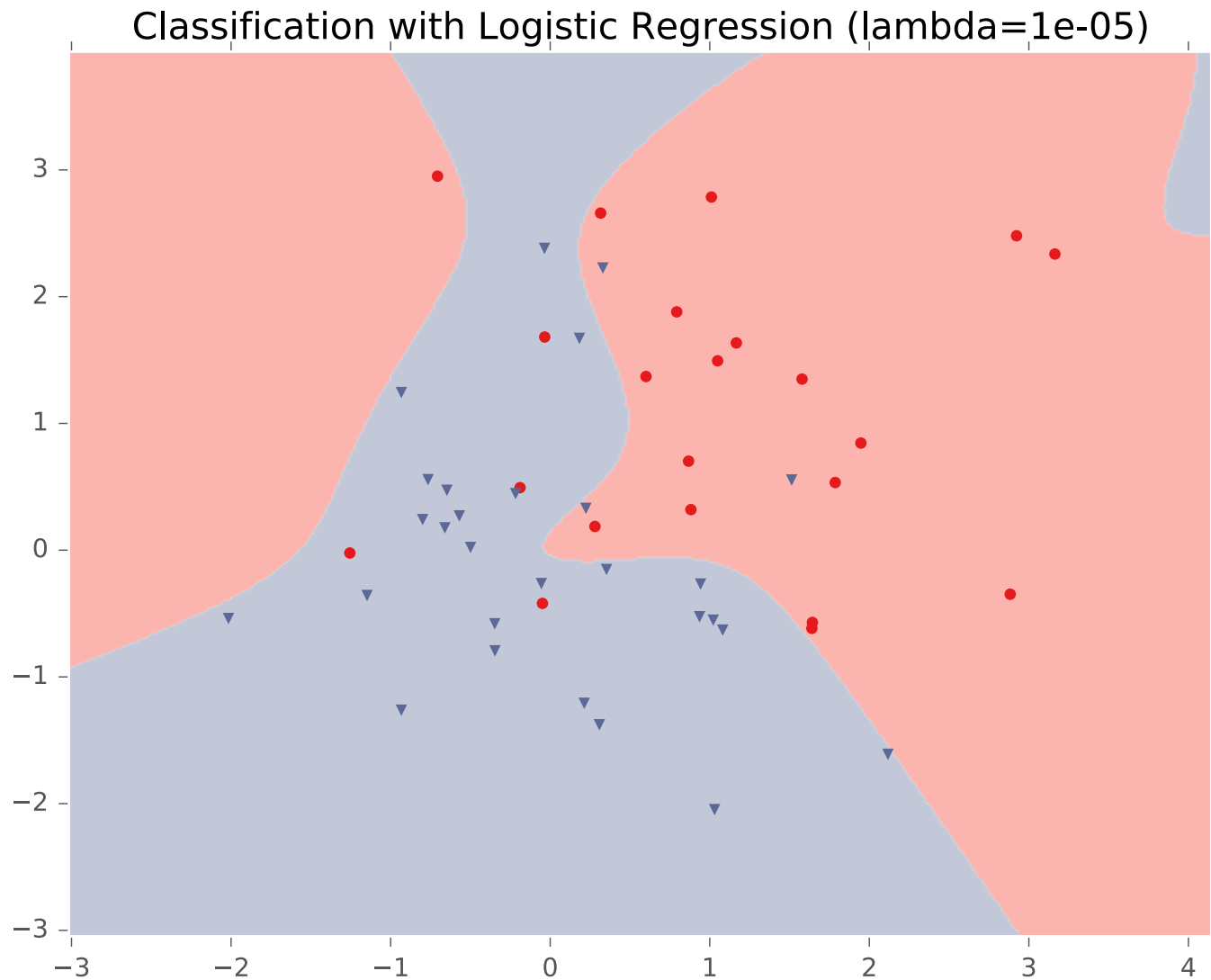
Test
Data



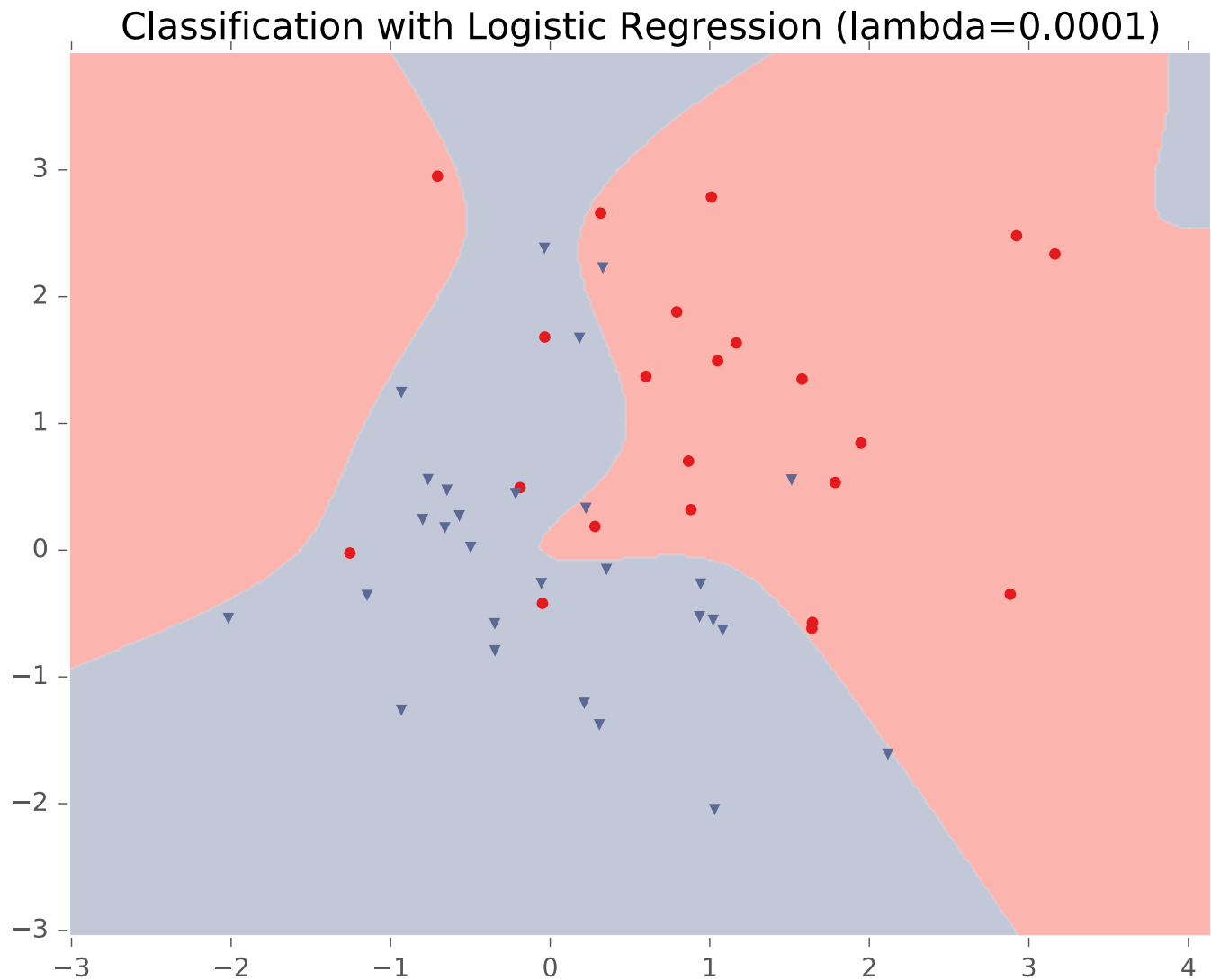
Example: Logistic Regression



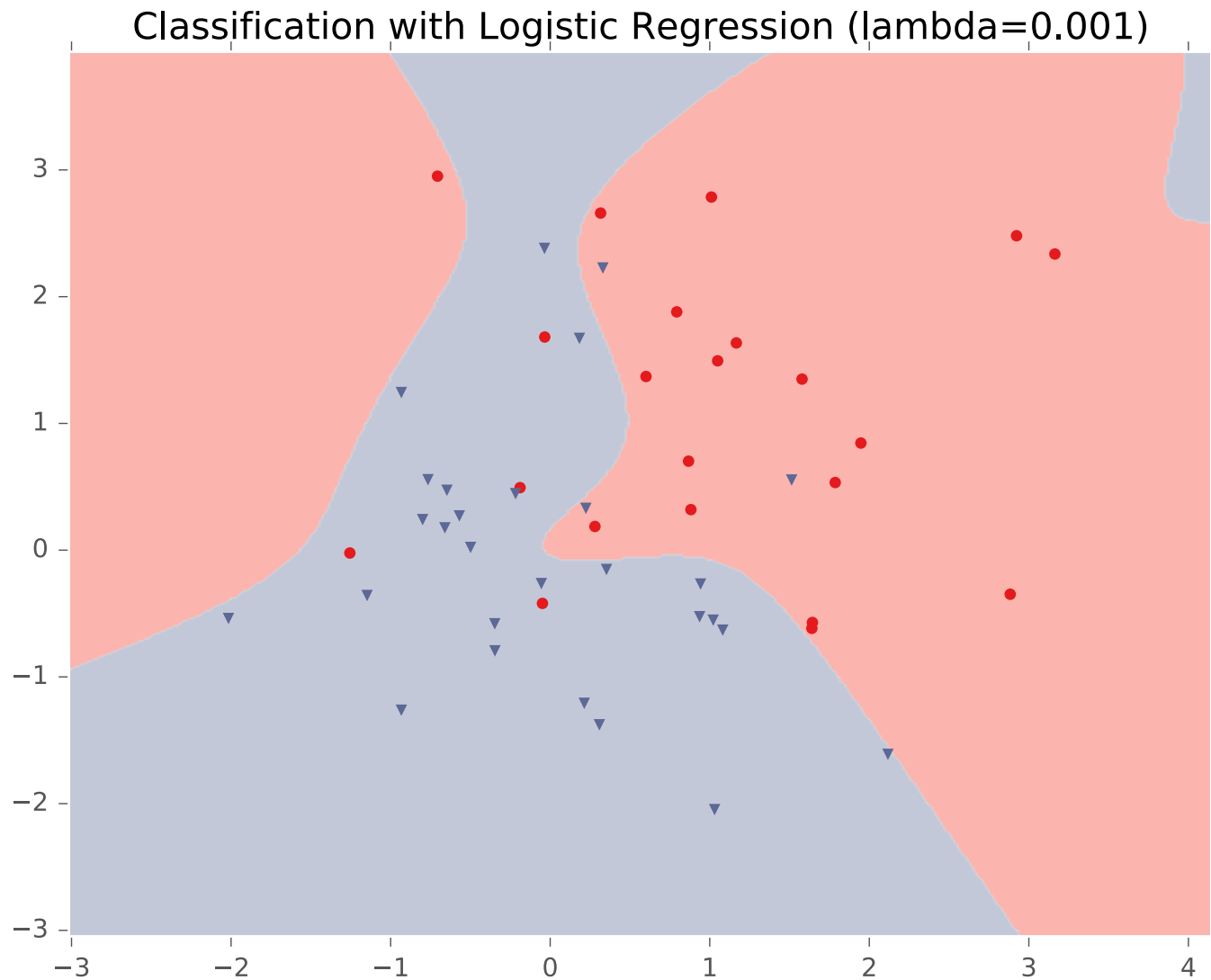
Example: Logistic Regression



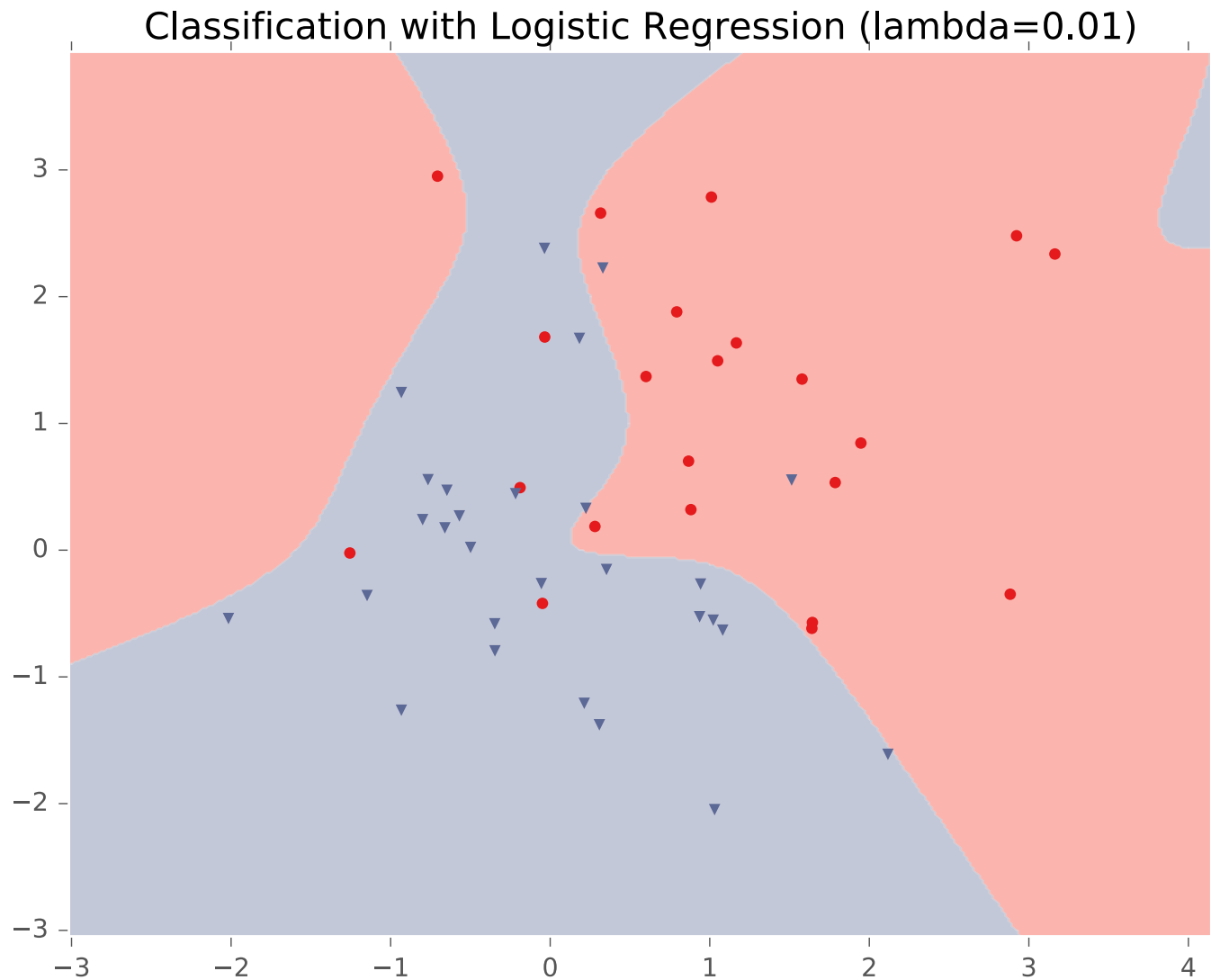
Example: Logistic Regression



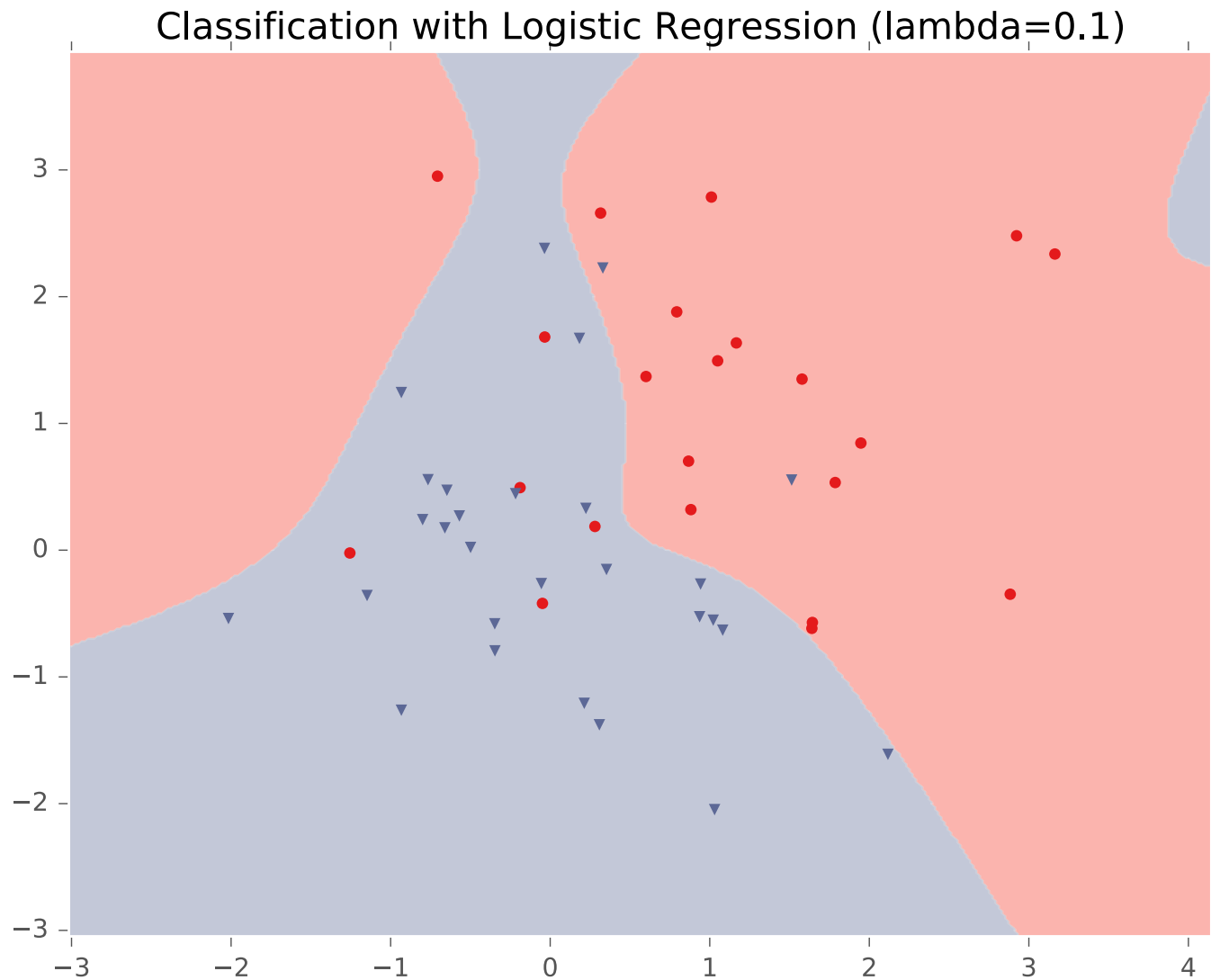
Example: Logistic Regression



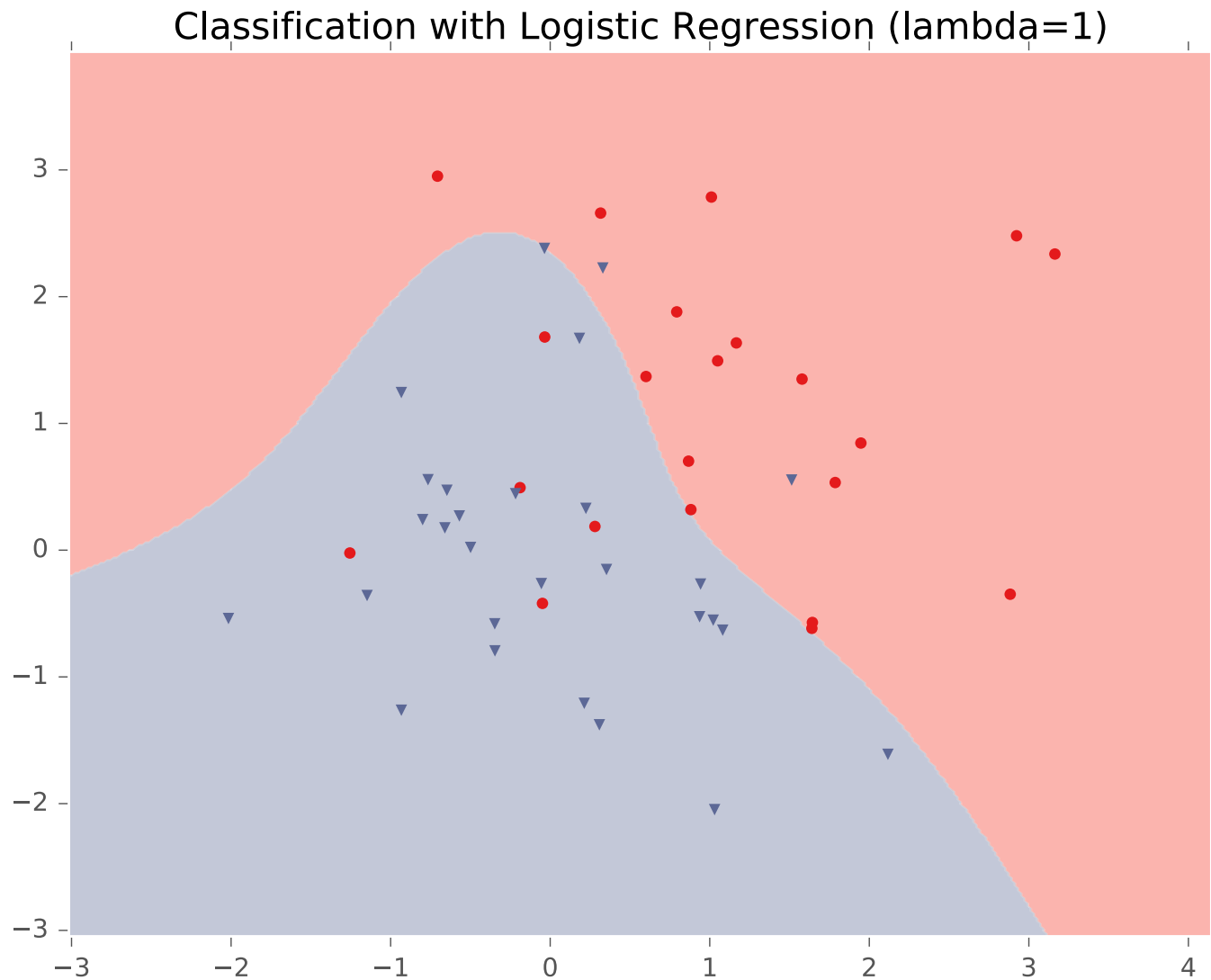
Example: Logistic Regression



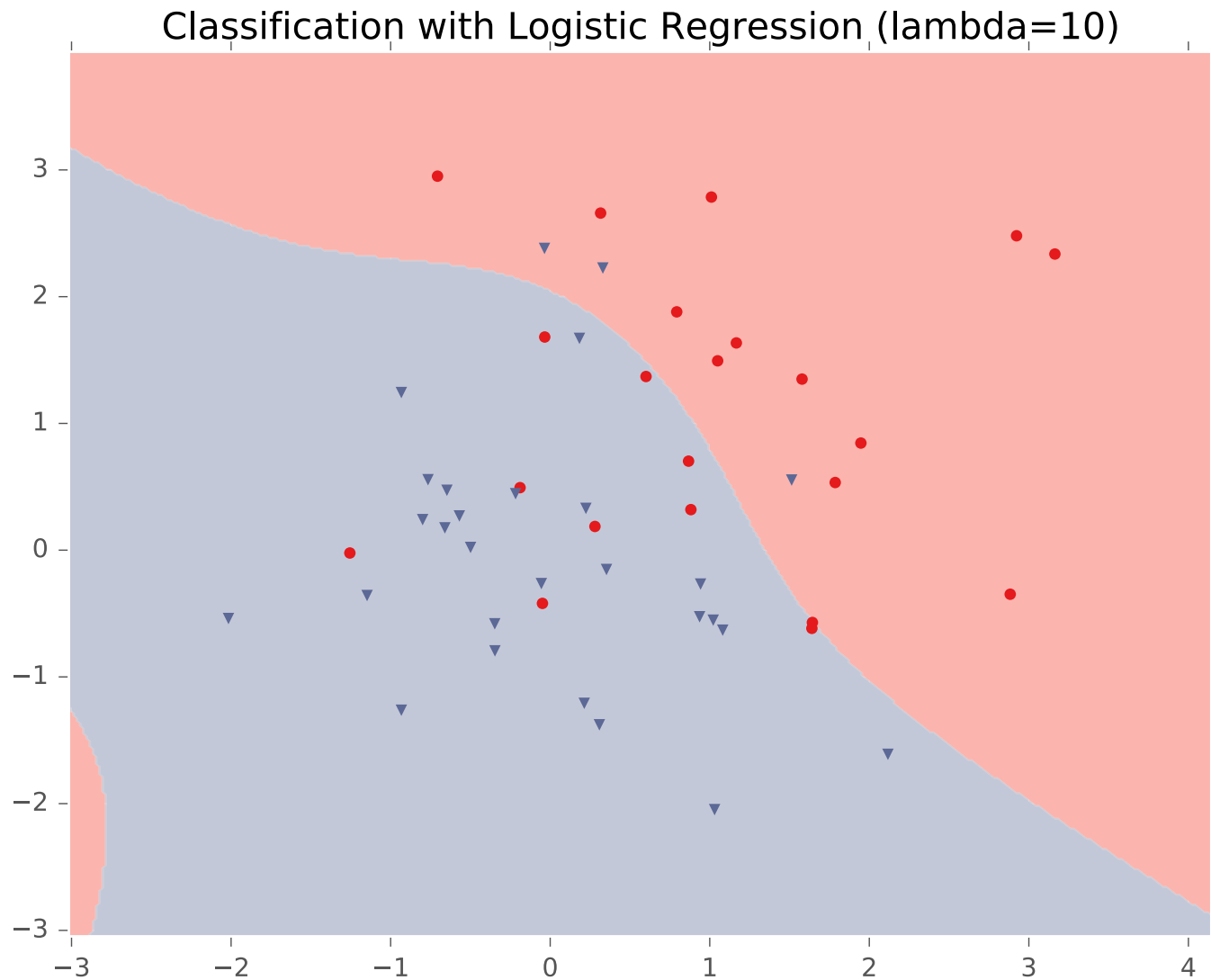
Example: Logistic Regression



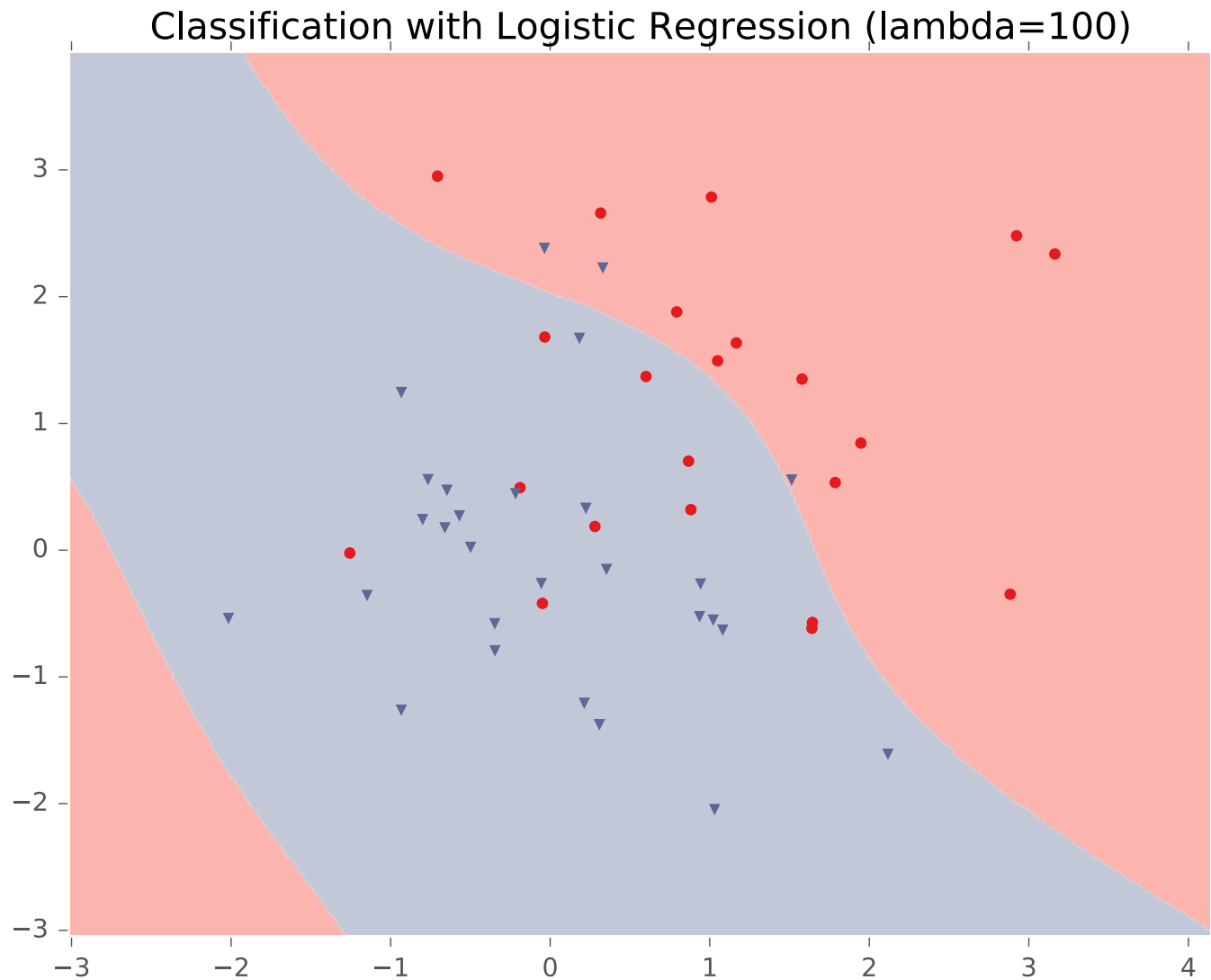
Example: Logistic Regression



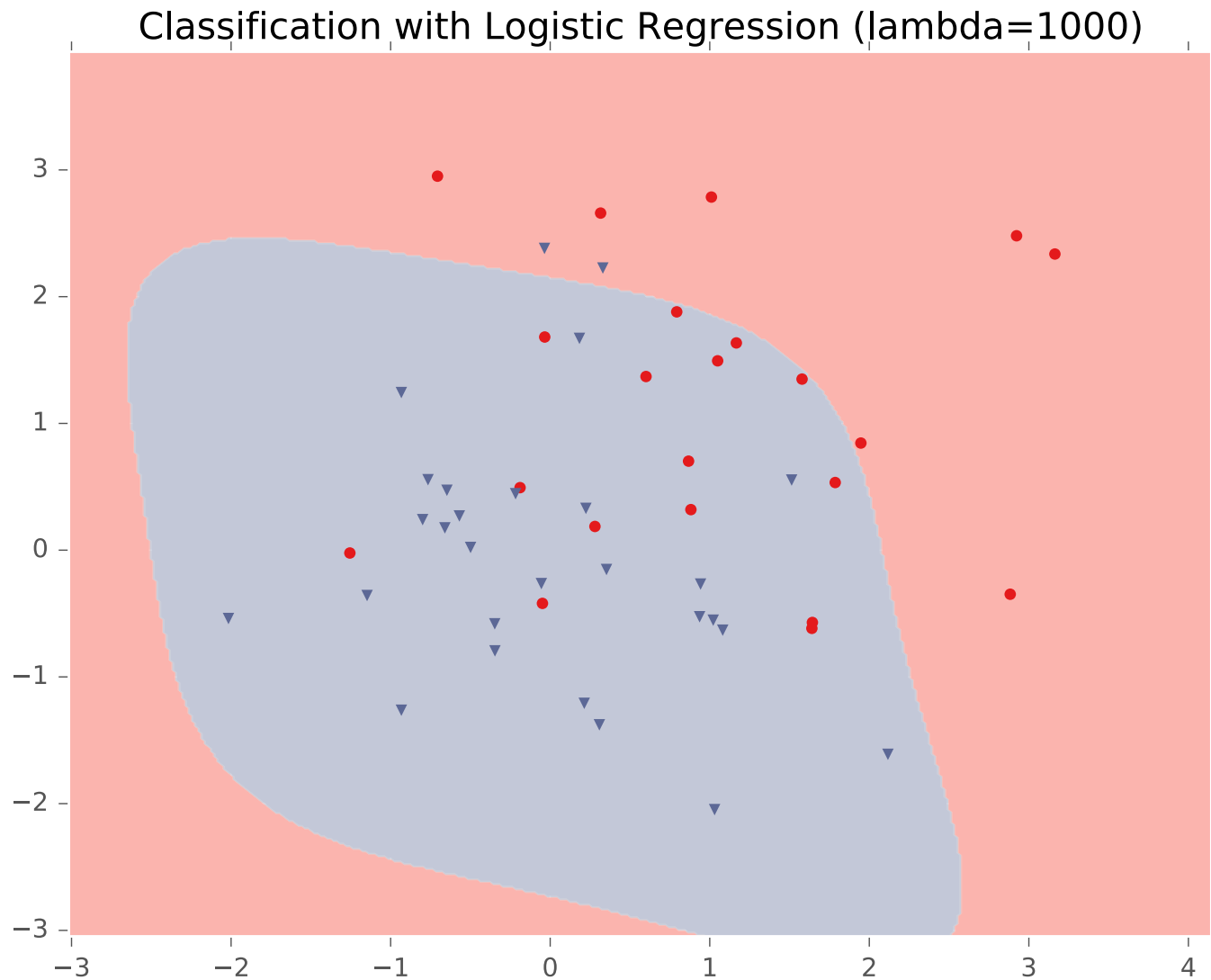
Example: Logistic Regression



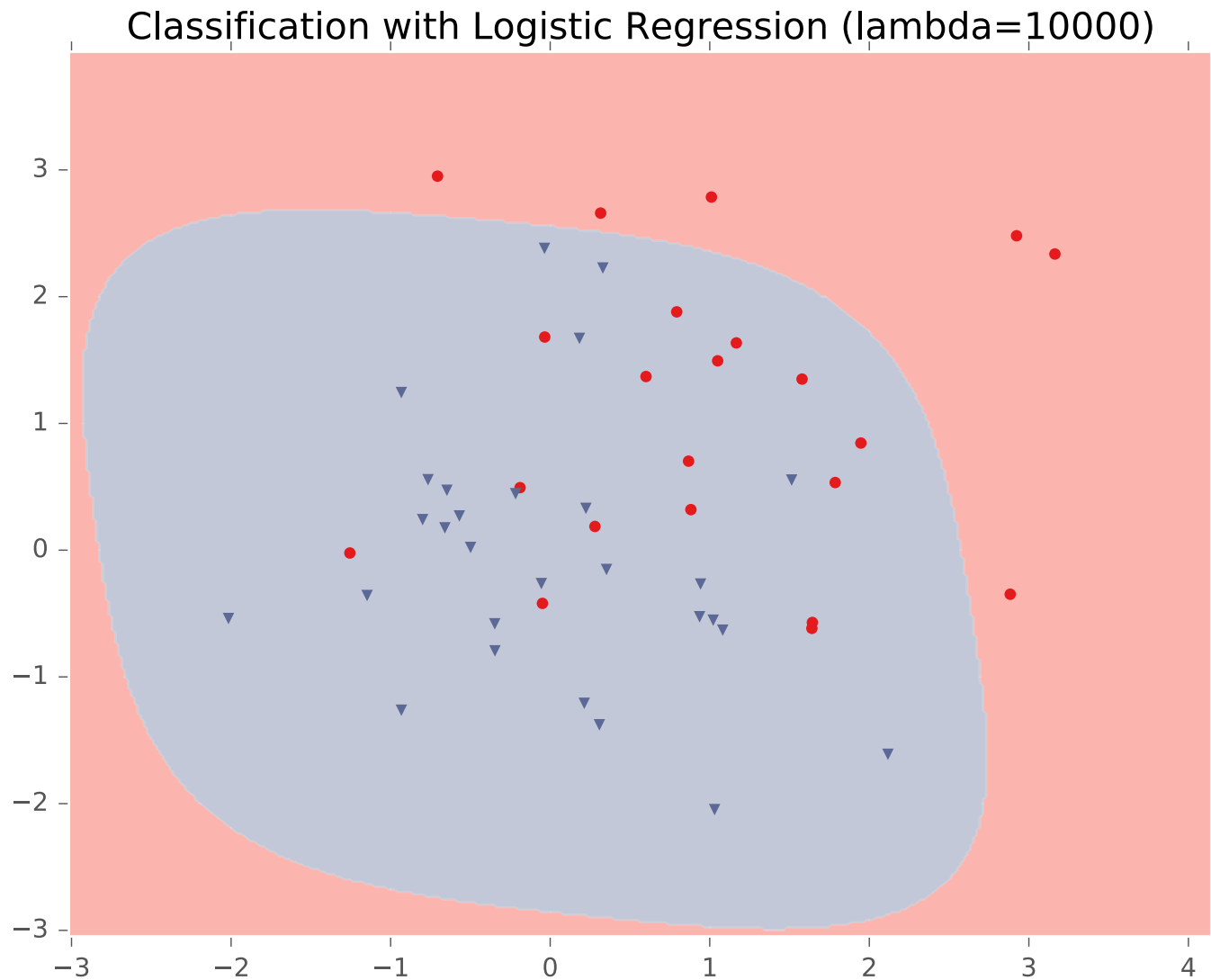
Example: Logistic Regression



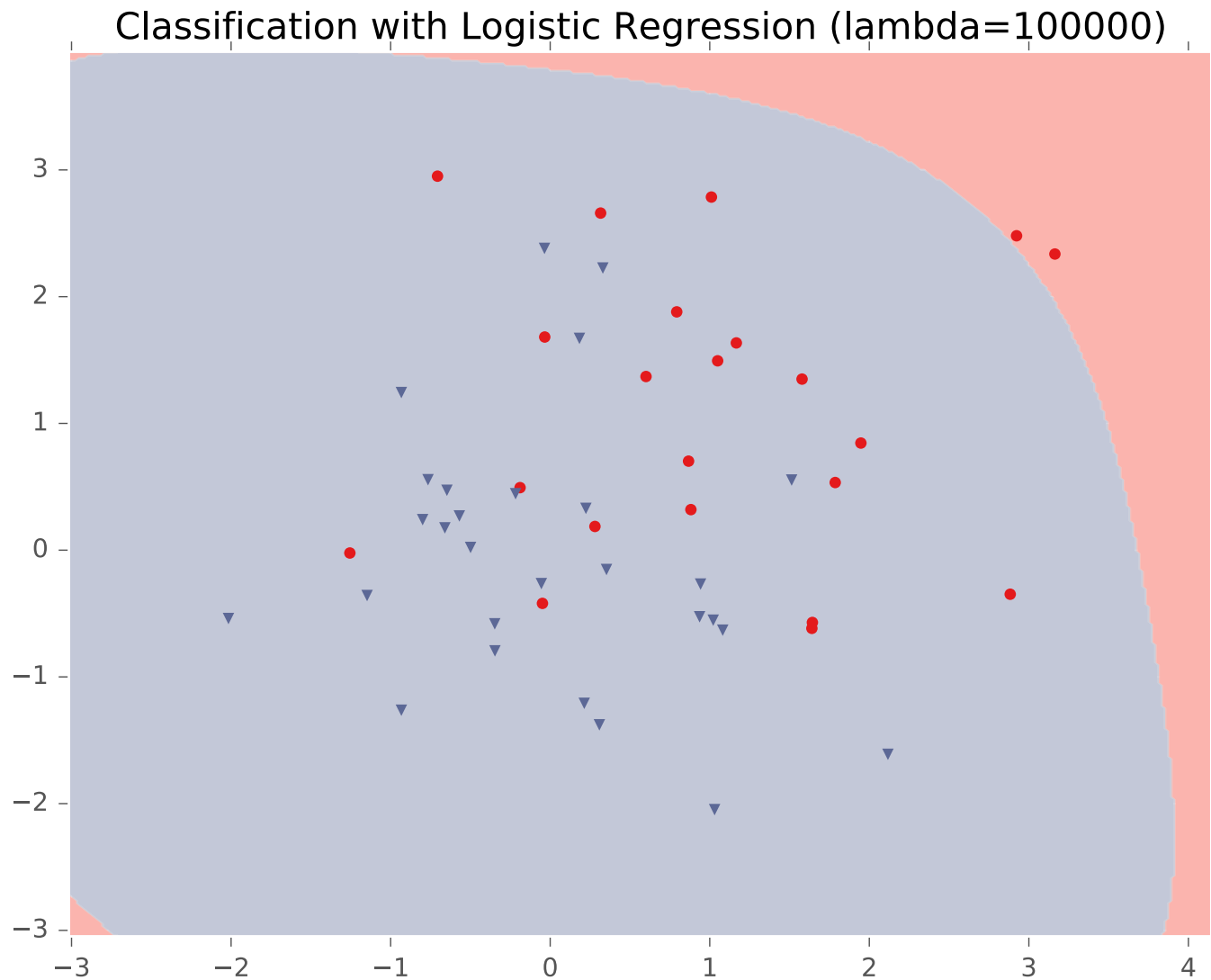
Example: Logistic Regression



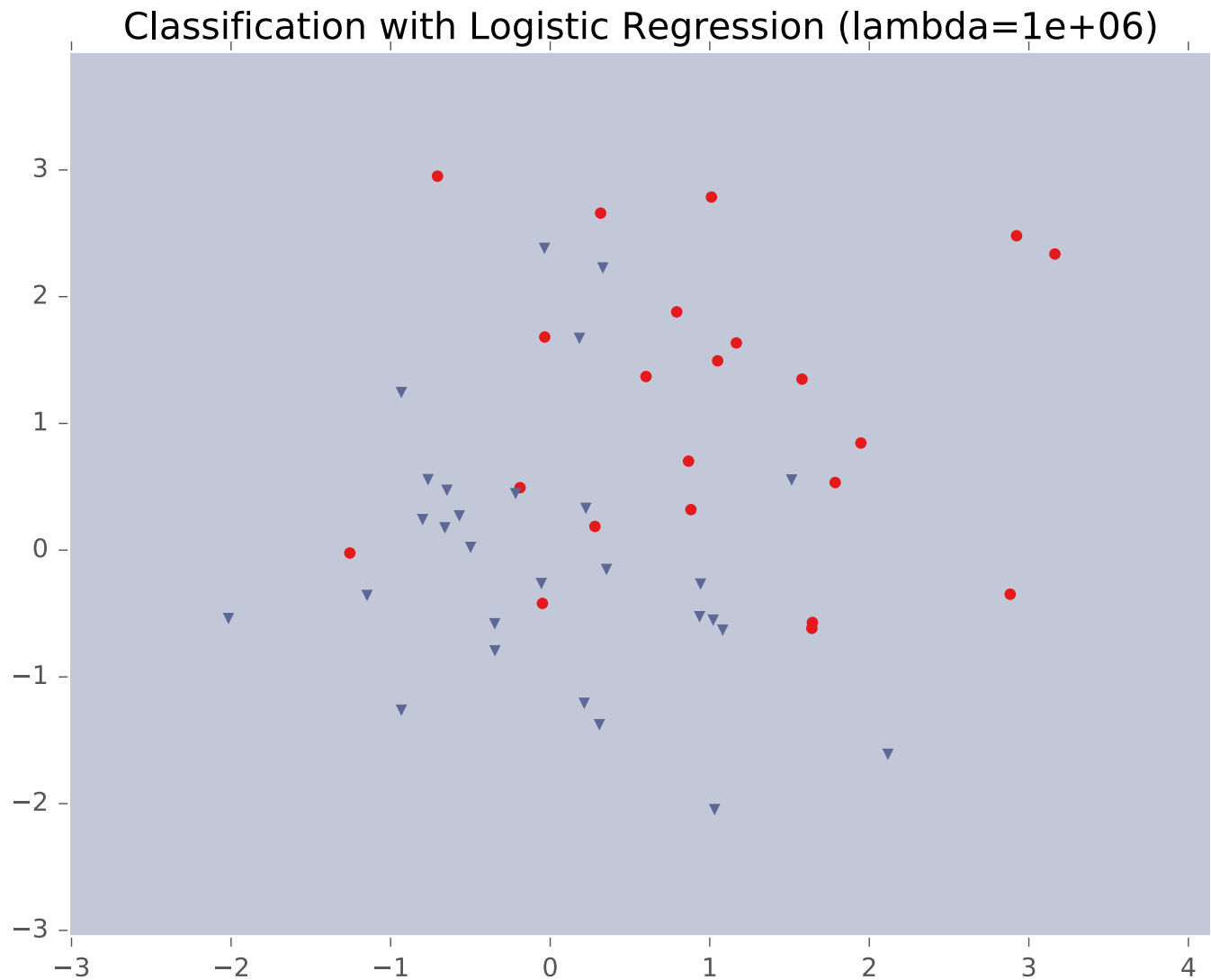
Example: Logistic Regression



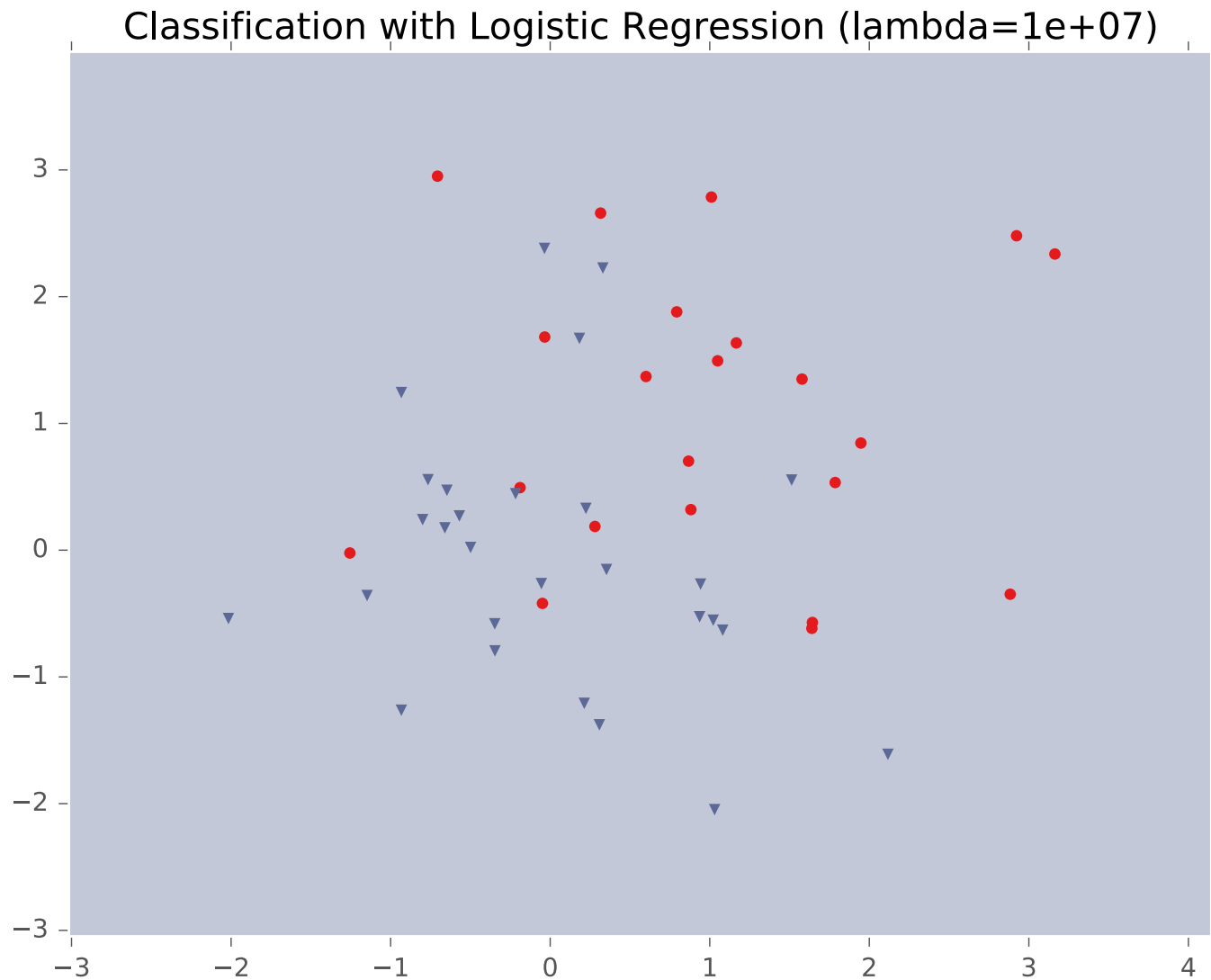
Example: Logistic Regression



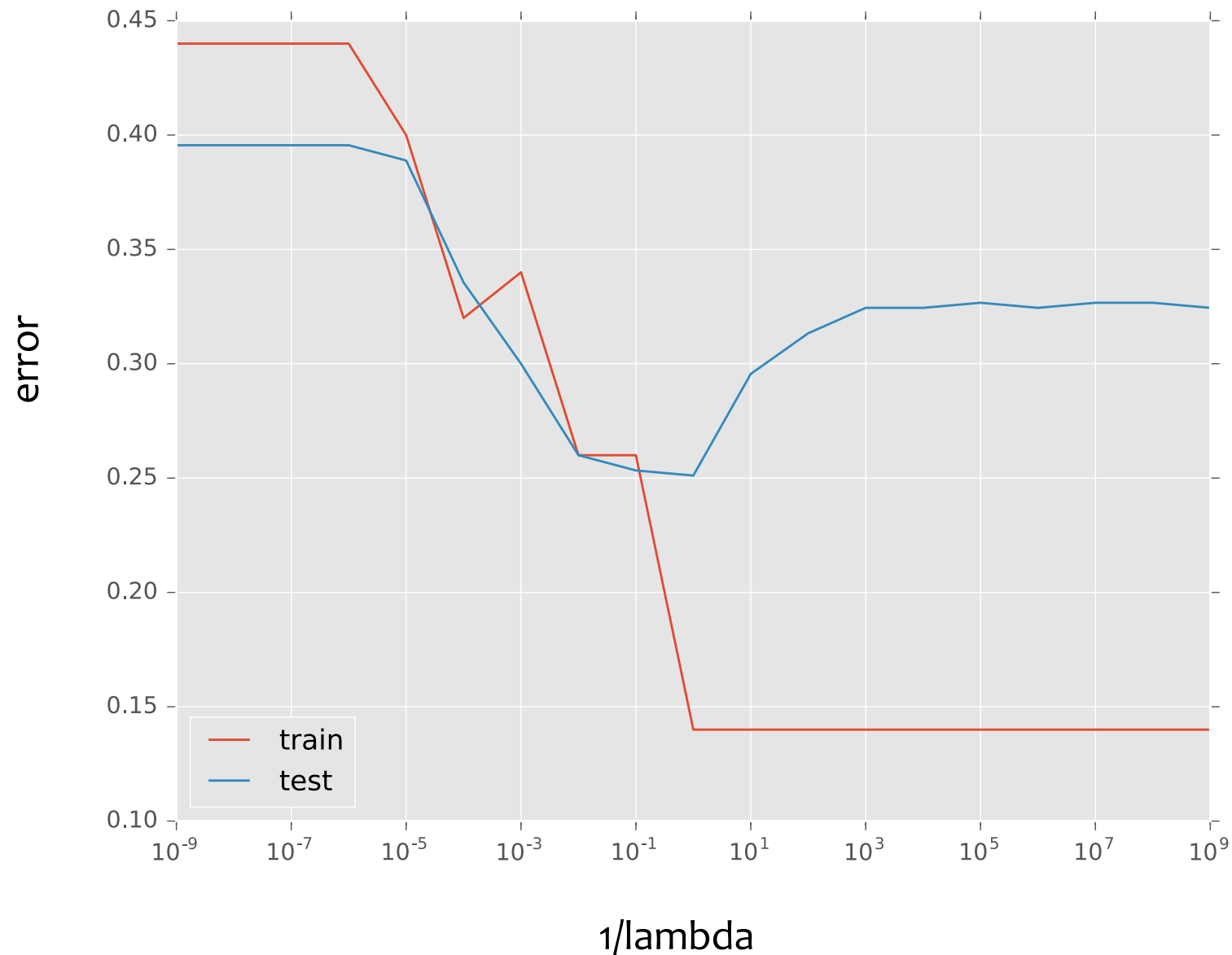
Example: Logistic Regression



Example: Logistic Regression



Example: Logistic Regression



Regularization as MAP

- L1 and L2 regularization can be interpreted as **maximum a-posteriori (MAP) estimation** of the parameters
- To be discussed later in the course...

Takeaways

1. **Nonlinear basis functions** allow **linear models** (e.g. Linear Regression, Logistic Regression) to capture **nonlinear** aspects of the original input
2. Nonlinear features **require no changes to the model** (i.e. just preprocessing)
3. **Regularization** helps to avoid **overfitting**
4. **Regularization** and **MAP estimation** are equivalent for appropriately chosen priors

Feature Engineering / Regularization Objectives

You should be able to...

- Engineer appropriate features for a new task
- Use feature selection techniques to identify and remove irrelevant features
- Identify when a model is overfitting
- Add a regularizer to an existing objective in order to combat overfitting
- Explain why we should **not** regularize the bias term
- Convert linearly inseparable dataset to a linearly separable dataset in higher dimensions
- Describe feature engineering in common application areas

ONLINE LEARNING LAB

Online Learning Lab

- Meet at the Gates-Hillman Center (GHC), room 6115
 - Use the 5th floor entrance near Cyert Hall
 - Take the elevator to the 6th floor
 - GHC 6115 is on your left as you exit the elevator
 - Email if you have trouble

