

RECITATION 3 CLASSIFICATION AND REGRESSION

10-301/10-601: INTRODUCTION TO MACHINE LEARNING
 02/08/2023

1 Decision Trees and Beyond

1. Decision Tree Classification with Continuous Attributes
 Given the dataset $D_1 = \{x^{(i)}, y^{(i)}\}_{i=1}^n$ where $x^{(i)} \in \mathbb{R}^2$, $y^{(i)} \in \{\text{Yellow, Purple, Green}\}$ as shown in Fig. 1, we wish to learn a decision tree for classifying such points. Provided with a possible tree structure in Fig. 1, what values should each leaf node predict to minimize the training Mean Squared Error (MSE) of our regression? Assume each leaf node just predicts a constant.

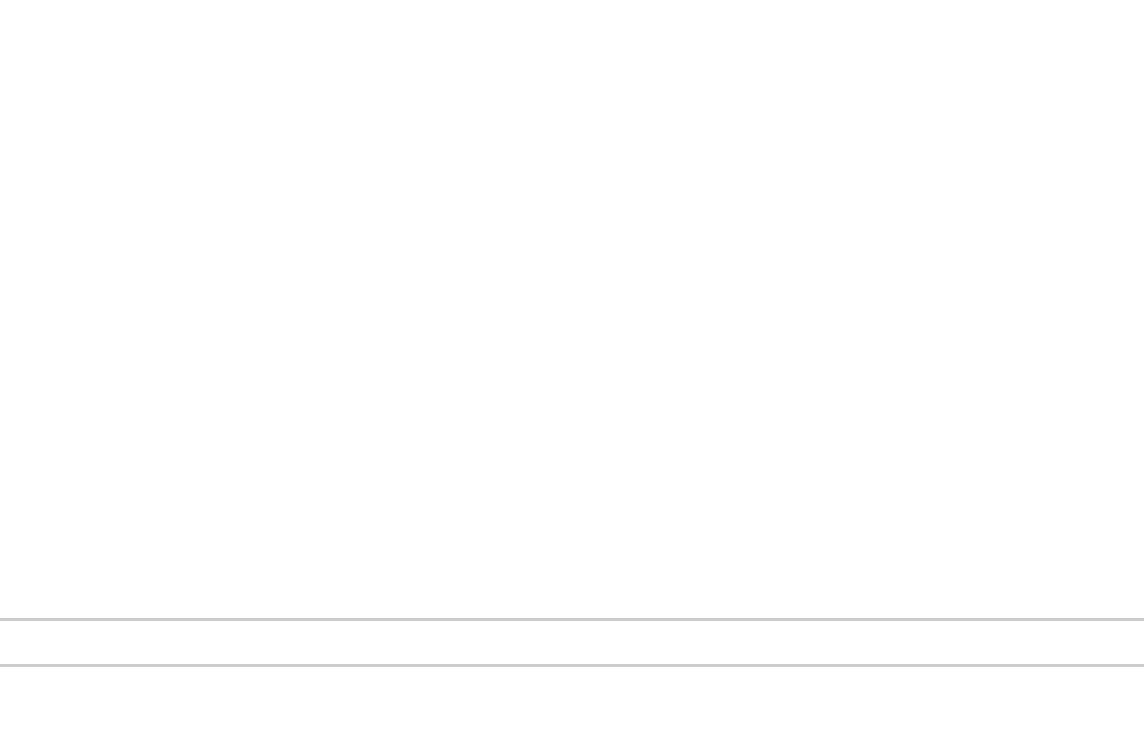


Figure 1: Classification of 2D points, with Decision Tree to fill in

2. Decision Tree Regression with Continuous Attributes
 Now instead if we had dataset $D_2 = \{x^{(i)}, y^{(i)}\}_{i=1}^n$ where $x^{(i)} \in \mathbb{R}^2$, $y^{(i)} \in \mathbb{R}$ as shown in Fig. 2, we wish to learn a decision tree for regression on such points. Using the same tree structure and values of α, β as before, what values should each leaf node predict to minimize the training Mean Squared Error (MSE) of our regression? Assume each leaf node just predicts a constant.

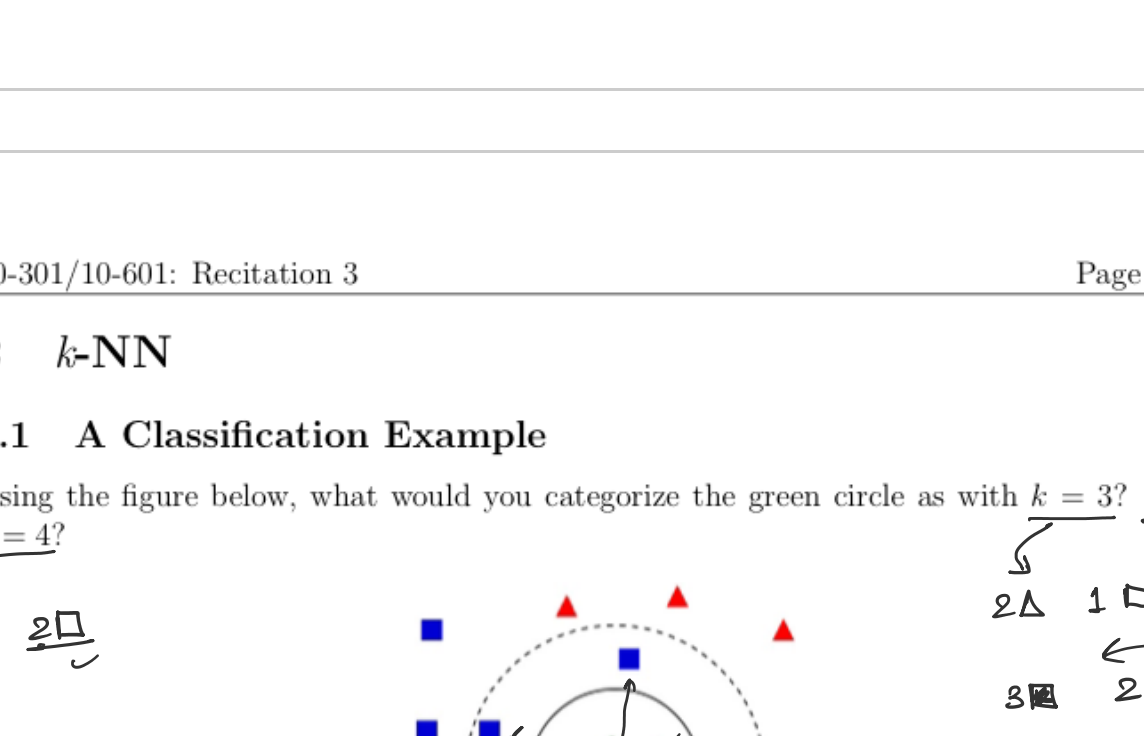


Figure 2: Regression on 2D points, with Decision Tree to fill in

$$J = \text{MSE} = \frac{1}{n} \sum_{i=1}^n \|y_i - \hat{y}_i\|^2 \Rightarrow \frac{\partial J}{\partial \hat{y}_i} = 0 \Rightarrow \frac{\partial}{\partial \hat{y}_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 0 \Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i) = 0$$

$$\text{min MSE} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\sum_{i=1}^n y_i - n\hat{y} = 0 \Rightarrow \hat{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

3. Choosing a Tree: What might happen if we increased the max-depth of the tree? When predicting on unseen data, would we prefer the depth-2 tree above or a very deep tree?

2 k-NN

2.1 A Classification Example

Using the figure below, what would you categorize the green circle as with $k=3$? $k=5$? $k=4$?

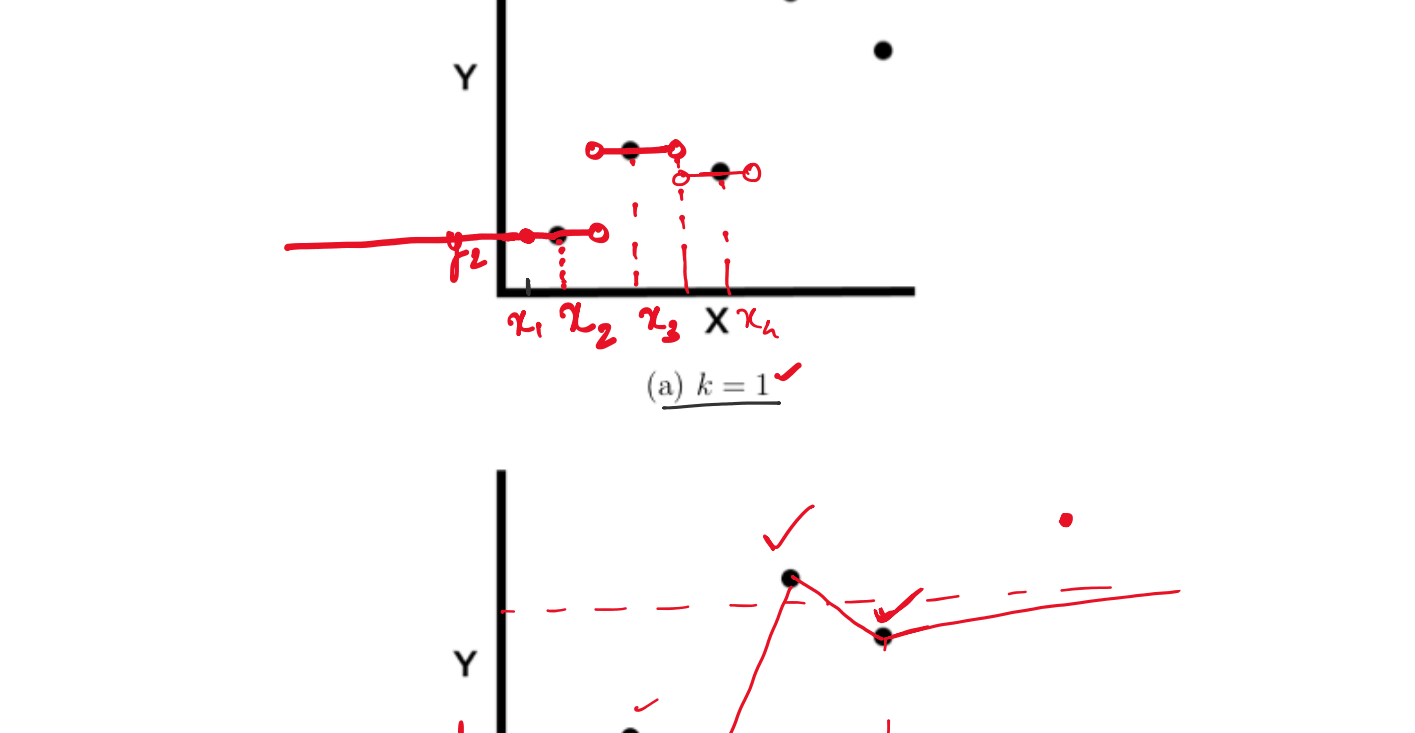
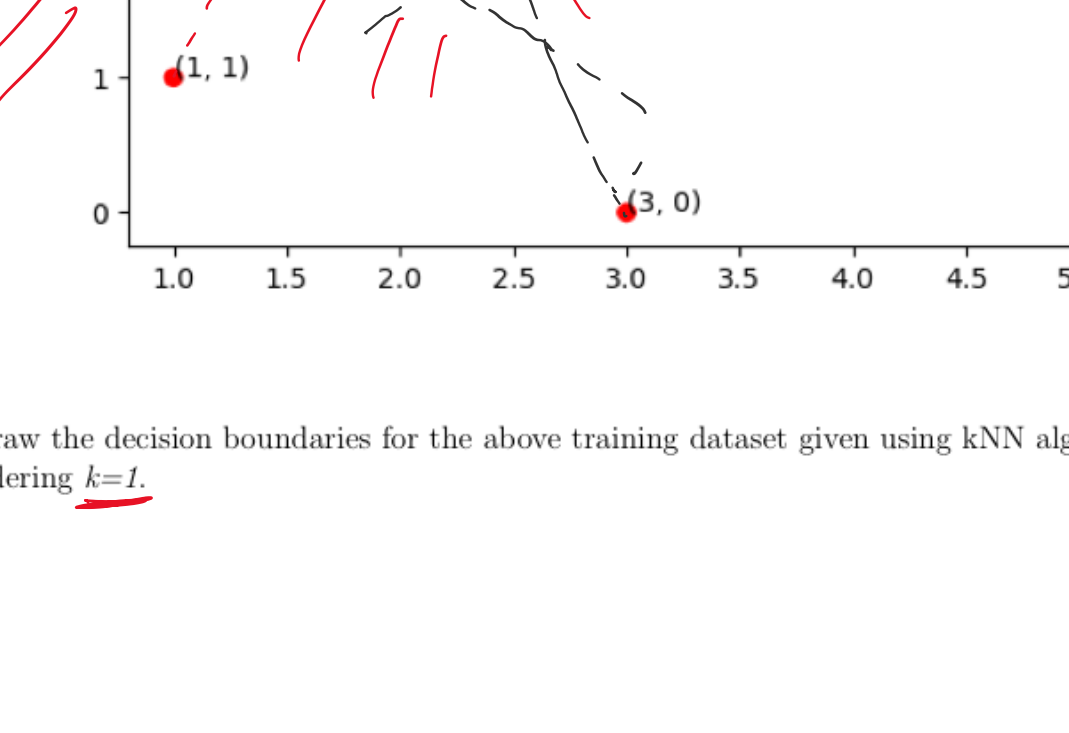
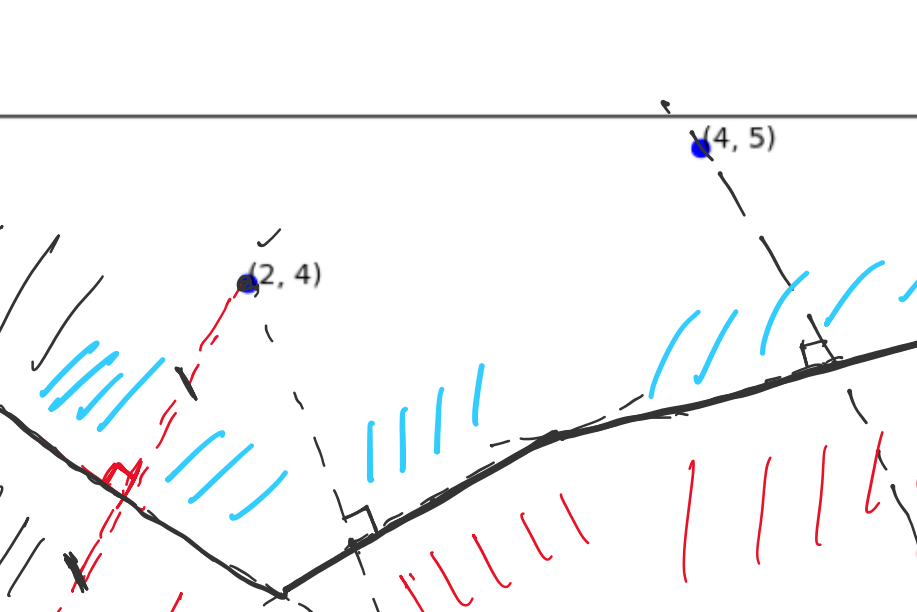


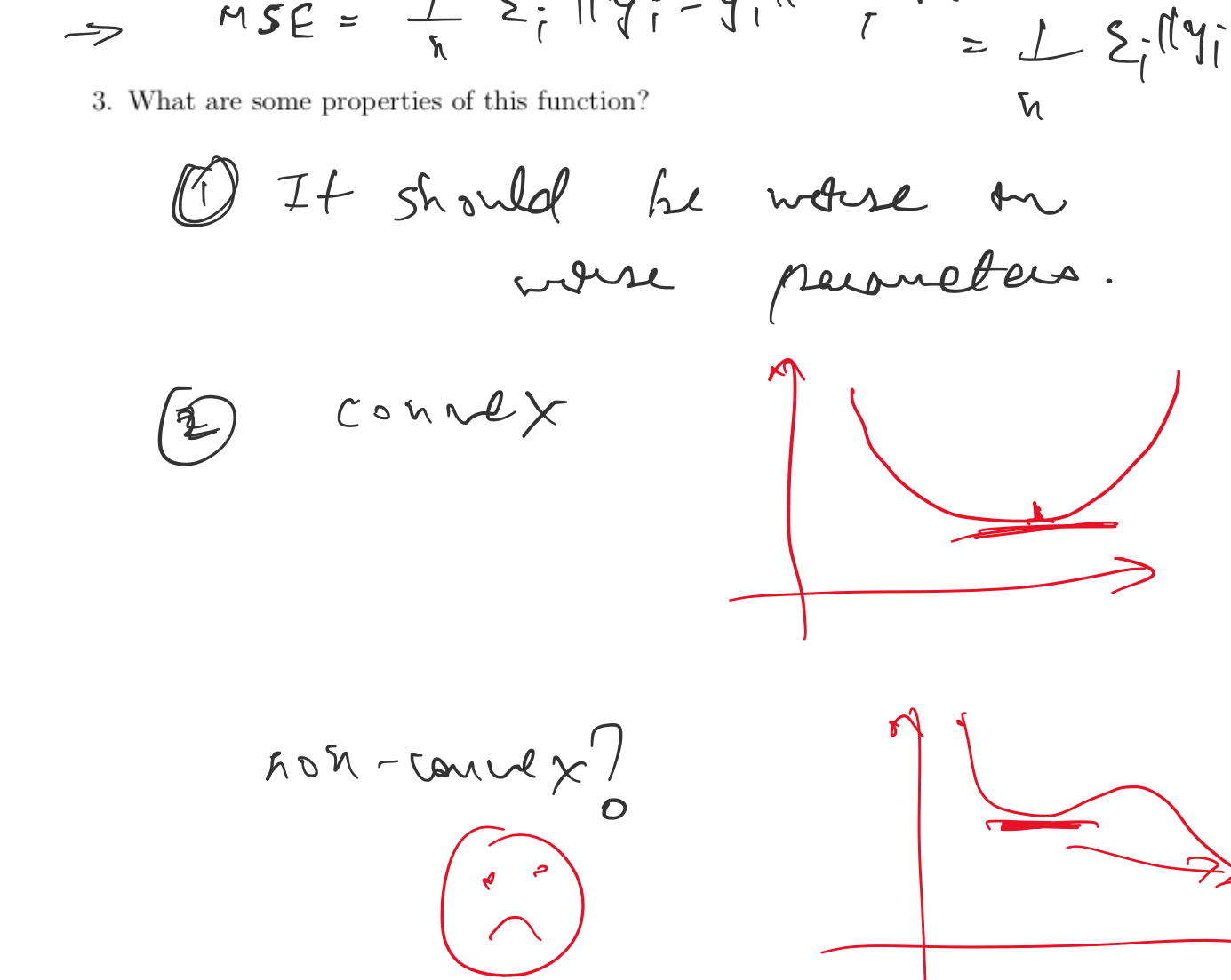
Figure 3: An example of k-NN on a small dataset; image source from Wikipedia

2.2 k-NN for Regression

You want to predict a continuous variable Y with a continuous variable X . Having just learned k -NN, you are super eager to try it out for regression. Given the data below, draw the regression lines (what k -NN would predict Y to be for every X value if it was trained for the given data) for k -NN regression with $k=1$, weighted $k=2$, and unweighted $k=2$. For weighted $k=2$, take the weighted average of the two nearest points. For unweighted $k=2$, take the unweighted average of the two nearest points. (Note: the points are equidistant along the x-axis)



2.3 k-NN Decision Boundary



Draw the decision boundaries for the above training dataset given using kNN algorithm considering $k=1$.

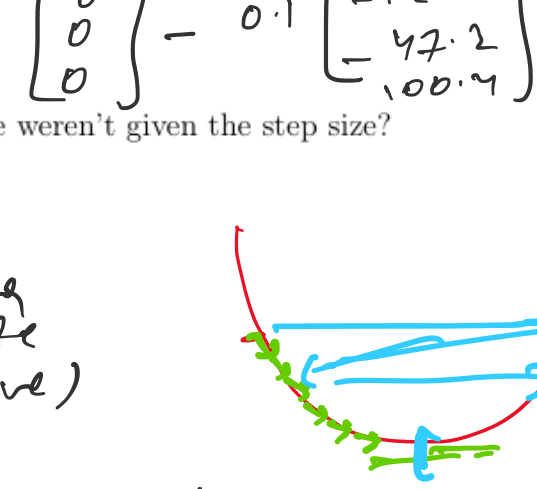
3 Linear Regression

3.1 Defining the Objective Function

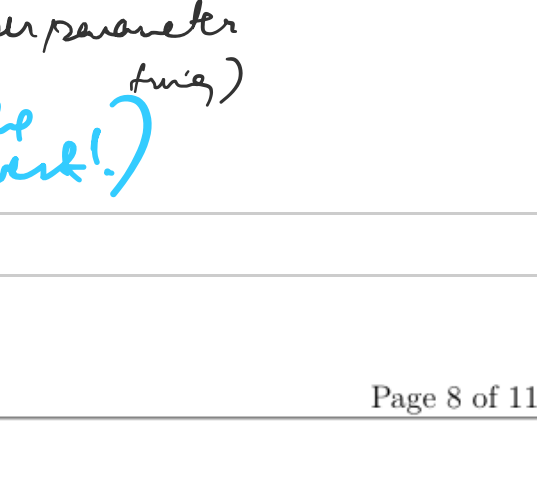
1. What does an objective function $J(\theta)$ do?
 → How bad your parameters are!

2. What are some examples?
 → $MSE = \frac{1}{n} \sum_{i=1}^n \|y_i - \hat{y}_i\|^2$, $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

3. What are some properties of this function?
 (1) It should be convex or concave.
 (2) convex



non-convex?



(3) Differentiable w.r.t to the parameters (eq. for gradient descent)

3.2 Solving Linear Regression using Gradient Descent

	$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$x^{(4)}$	$x^{(5)}$
x_1	1.0	2.0	3.0	4.0	5.0
x_2	-2.0	-5.0	-6.0	-8.0	-11.0
y	2.0	4.0	7.0	8.0	11.0

$= 33$

Now, we want to implement the gradient descent method.
 Assuming that $\gamma = 0.1$ and θ has been initialized to $[0, 0, 0]^T$, perform one iteration of gradient descent:

1. What is the gradient of the objective function $J(\theta)$ with respect to θ : $\nabla_{\theta} J(\theta)$?

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2 \quad \hat{y}_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

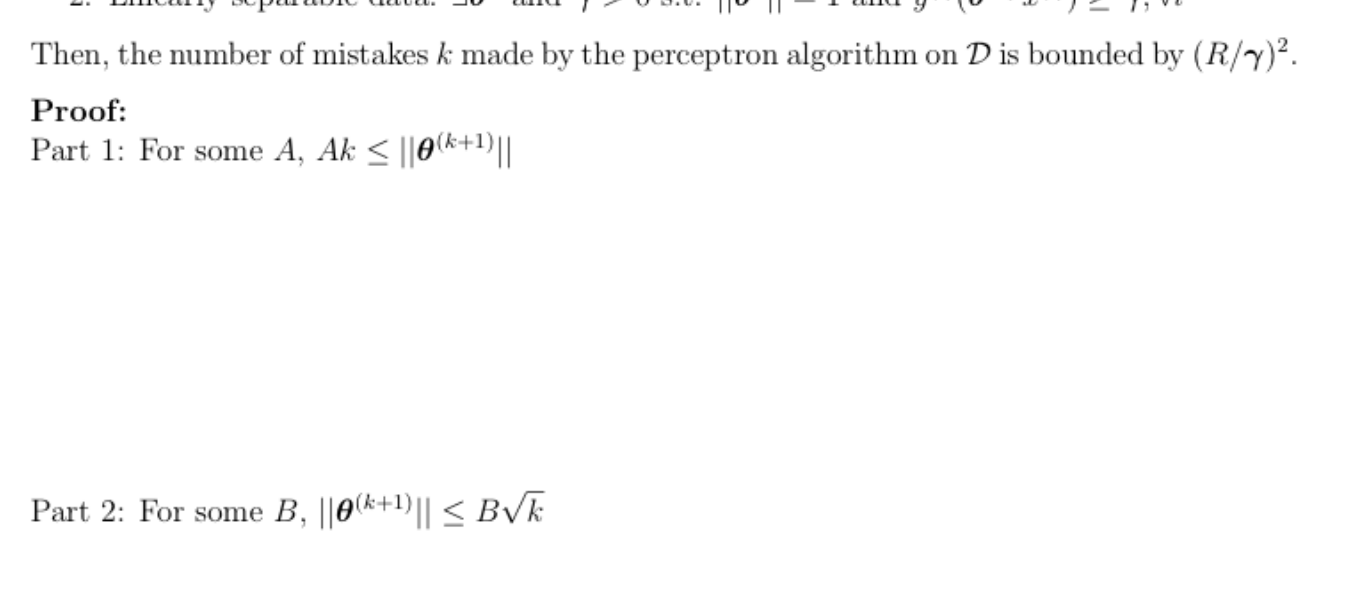
$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum_{i=1}^n (\theta^T x_i - y_i) \cdot 2 \cdot x_{ij} \quad \nabla_{\theta} J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \begin{bmatrix} x_{10} \\ x_{11} \\ x_{12} \end{bmatrix}$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{2}{n} \sum_{i=1}^n (y_i - \theta^T x_i) x_i \quad = \begin{bmatrix} -12.8 \\ -42.2 \\ 10.0 \end{bmatrix}$$

2. How do we carry out the update rule?
 $\theta_{new} = \theta - \gamma \nabla_{\theta} J(\theta)$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -12.8 \\ -42.2 \\ 10.0 \end{bmatrix} = \begin{bmatrix} 1.28 \\ 4.22 \\ -1.00 \end{bmatrix}$$

3. How could we pick which value of γ to use if we weren't given the step size?



4 Perceptron

4.1 Perceptron Mistake Bound Guarantee

If a dataset has margin γ and all points inside a ball of radius R , then the perceptron makes less than or equal to $(R/\gamma)^2$ mistakes.

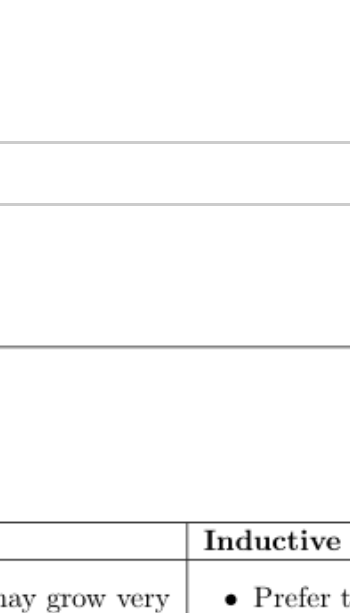


Figure 5: Perceptron Mistake Bound Setup

4.2 Definitions

- Margin:**
 - The margin of example x wrt a linear separator w is the (absolute) distance from x to the plane $w \cdot x = 0$.
 - The margin γ_w of a set of examples S wrt a linear separator w is the smallest margin over points $x \in S$.
 - The margin γ of a set of examples S is the maximum γ_w over all linear separators w .
- Linear Separability:** For a binary classification problem, a set of examples S is linearly separable if there exists a linear decision boundary that can separate the points.
- Update Rule:** When the k -th mistake is made on data point $x^{(k)}$, the parameter update is $\theta^{(k+1)} = \theta^{(k)} + y^{(k)} x^{(k)}$

We say the (batch) perceptron algorithm has converged when it stops making mistakes on the training data.

4.3 Perceptron Mistake Bound: Example

Given dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, suppose:

- Finite size inputs: $\|x^{(i)}\| \leq R$
- Linearly separable data: $\exists \theta^*$ and $\gamma > 0$ s.t. $\|\theta^*\| = 1$ and $y^{(i)} (\theta^* \cdot x^{(i)}) \geq \gamma, \forall i$

Then, the number of mistakes k made by the perceptron algorithm on \mathcal{D} is bounded by $(R/\gamma)^2$. The following table shows a dataset of linearly separable datapoints.

x_1	x_2	y
1	-1	1
0	2	-1
4	0	1

Assuming that the linear separator with the largest margin is given by:

$\theta^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$, where $\theta = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

Calculate the theoretical mistake bound for the perceptron.

$$R = \frac{\sqrt{(4-0)^2 + (0-0)^2}}{\gamma} = 4$$

$$\gamma = \frac{|[-1 \ 1] \cdot [-1 \ 1]^T|}{\sqrt{(-1)^2 + 1^2}} = \frac{|(-1) \cdot (-1) + 1 \cdot (-1)|}{\sqrt{2}} = \frac{2}{\sqrt{2}} = \sqrt{2}$$

$$\left(\frac{R}{\gamma}\right)^2 = \left(\frac{4}{\sqrt{2}}\right)^2 = \frac{16}{2} = 8$$

$$\theta^T x + b \quad \theta = [b \ \theta_1 \ \theta_2] \quad \theta^T x = b + \theta_1 x_1 + \theta_2 x_2$$

$$x = [1 \ x_1 \ x_2]$$

4.4 Theorem: Block, Novikoff

Given dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, suppose:

- Finite size inputs: $\|x^{(i)}\| \leq R$
- Linearly separable data: $\exists \theta^*$ and $\gamma > 0$ s.t. $\|\theta^*\| = 1$ and $y^{(i)} (\theta^* \cdot x^{(i)}) \geq \gamma, \forall i$

Then, the number of mistakes k made by the perceptron algorithm on \mathcal{D} is bounded by $(R/\gamma)^2$.

Proof:
 Part 1: For some $A, Ak \leq \|\theta^{(k+1)}\|$

Part 2: For some $B, \|\theta^{(k+1)}\| \leq B\sqrt{k}$

Part 3: Combine the bounds

Main Takeaway:

5 Summary

5.1 Decision Tree

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and interpret Very fast for inference 	<ul style="list-style-type: none"> Tree may grow very large and tend to overfit. Greedy behaviour may be sub-optimal 	<ul style="list-style-type: none"> Prefer the smallest tree consistent w/ the training data (i.e. 0 error rate) 	<ul style="list-style-type: none"> Most cases. Random forests are widely used in industry.

5.2 k-NN

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> No training of parameters Can apply to multi-class problems and use different metrics 	<ul style="list-style-type: none"> Slow for large datasets Must select good k Imbalanced data and outliers can lead to misleading results 	<ul style="list-style-type: none"> Similar (i.e. nearby) points should have similar labels All label dimensions are created equal 	<ul style="list-style-type: none"> Small dataset Small dimensionality Data is clean (no missing data) Inductive bias is strong for dataset

5.3 Linear regression

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and train Closed form solution 	<ul style="list-style-type: none"> Sensitive to noise (other than zero-mean Gaussian noise) 	<ul style="list-style-type: none"> The true relationship between the inputs and output is linear. 	<ul style="list-style-type: none"> Most cases (can be extended by adding non-linear feature transformations)

5.4 Perceptron

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and works for online learning. Provable guarantees on mistakes made for linearly separable data. 	<ul style="list-style-type: none"> No guarantees on finding best (maximum-margin) hyperplane. Output is sensitive to noise in the training data. 	<ul style="list-style-type: none"> The binary classes are separable in the feature space by a line. 	<ul style="list-style-type: none"> Not used much anymore, but variants (kernel perceptron) may have more success.

RECITATION 3 CLASSIFICATION AND REGRESSION

10-301/10-601: INTRODUCTION TO MACHINE LEARNING
 02/05/2023

1 Decision Trees and Beyond

1. Decision Tree Classification with Continuous Attributes
 Given the dataset $D_1 = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ where $x^{(i)} \in \mathbb{R}^2, y^{(i)} \in \{Yellow, Purple, Green\}$ as shown in Fig. 1, we wish to learn a decision tree for classifying such points. Provided with a possible tree structure in Fig. 1, what values should each leaf node predict to minimize the training Mean Squared Error (MSE) of our regression? Assume each leaf node just predicts a constant.

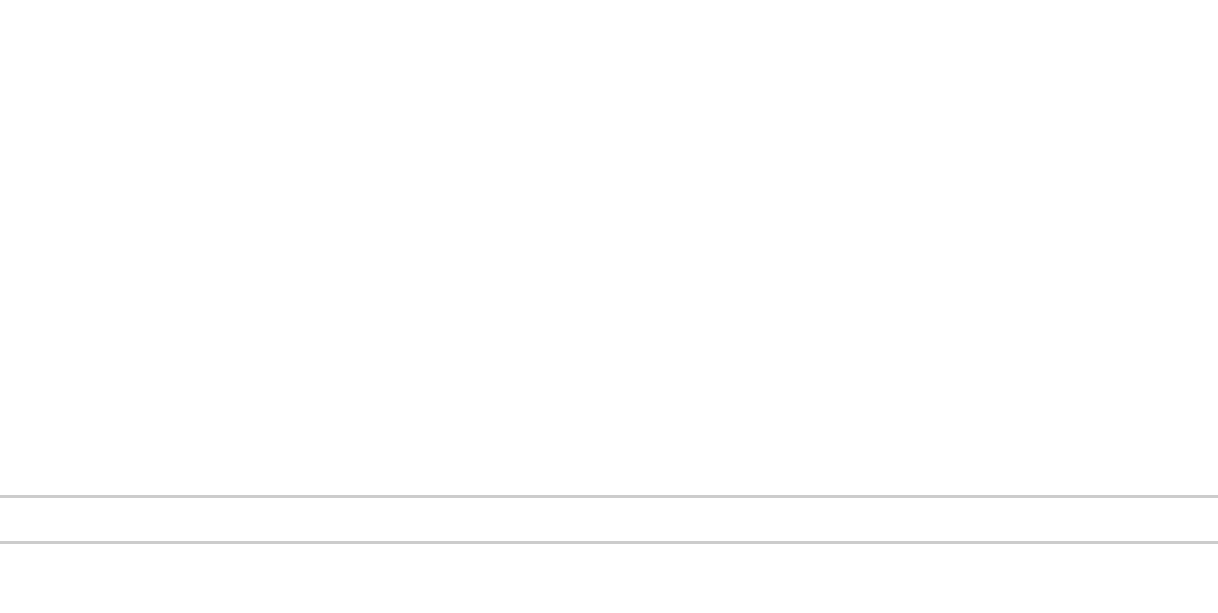


Figure 1: Classification of 2D points, with Decision Tree to fill in

10-301/10-601: Recitation 3 Page 2 of 11

2. Decision Tree Regression with Continuous Attributes
 Now instead if we had dataset $D_2 = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$ where $x^{(i)} \in \mathbb{R}^2, y^{(i)} \in \mathbb{R}$ as shown in Fig. 2, we wish to learn a decision tree for regression on such points. Using the same tree structure and values of α, β as before, what values should each leaf node predict to minimize the training Mean Squared Error (MSE) of our regression? Assume each leaf node just predicts a constant.

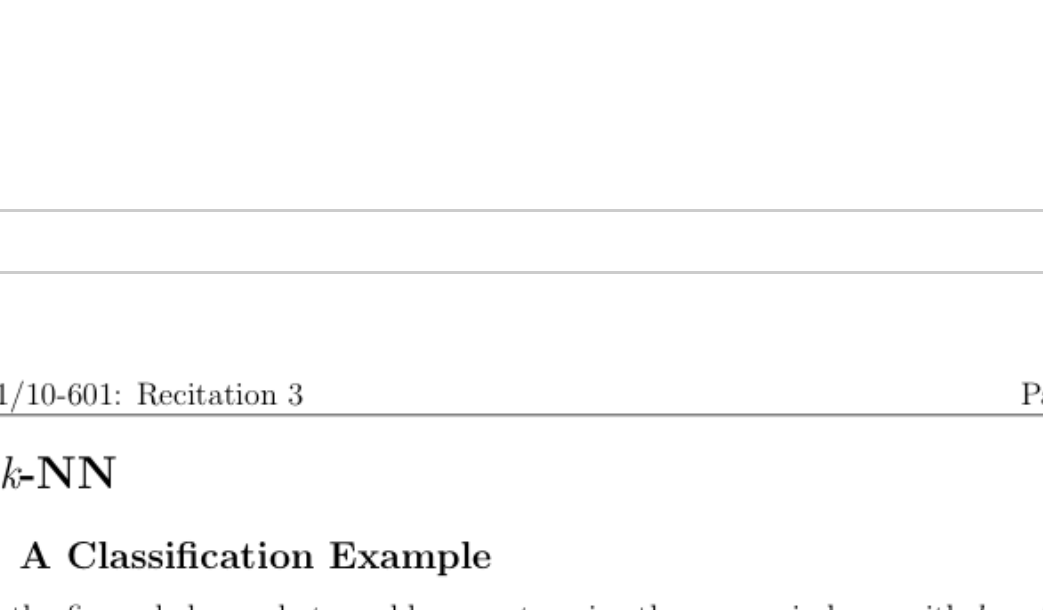


Figure 2: Regression on 2D points, with Decision Tree to fill in

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 \rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - q)^2 = 0$$

$$\rightarrow \sum_{i=1}^n y_i = n \cdot q \rightarrow q = \frac{\sum_{i=1}^n y_i}{n}$$

3. Choosing a Tree: What might happen if we increased the max-depth of the tree? When predicting on unseen data, would we prefer the depth-2 tree above or a very deep tree?

Overfitting!

10-301/10-601: Recitation 3 Page 3 of 11

2 k-NN

2.1 A Classification Example

Using the figure below, what would you predict the green circle as with $k=3$? $k=5$? $k=4$?

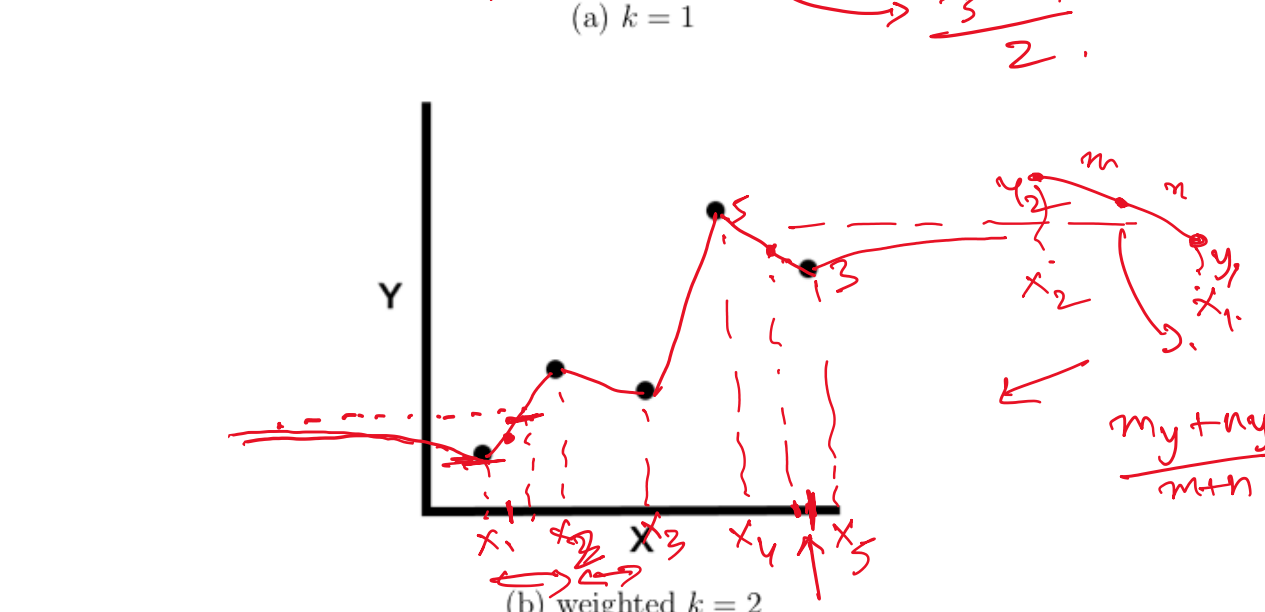
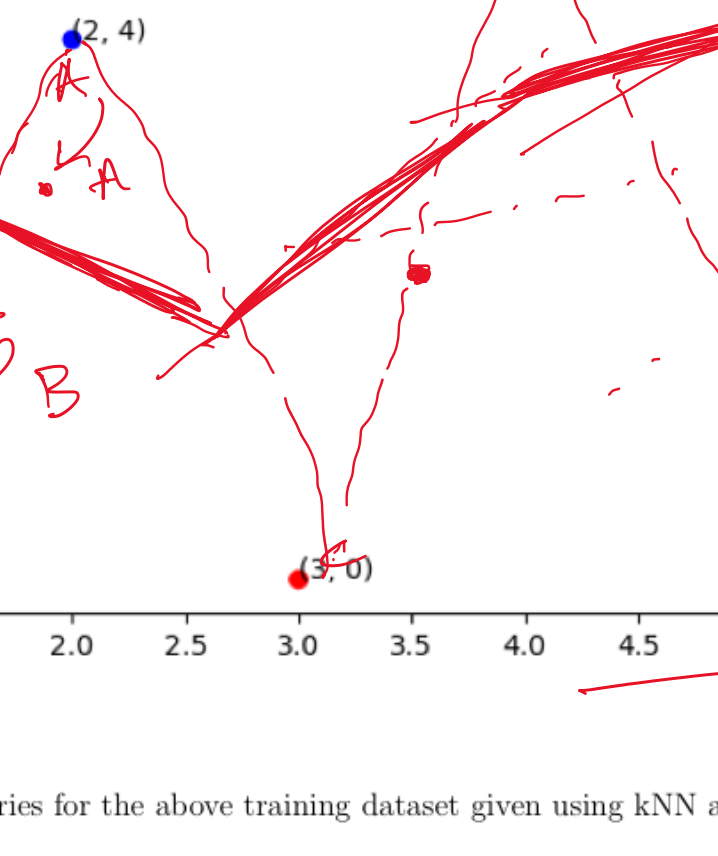


Figure 3: An example of k-NN on a small dataset; image source from Wikipedia

2.2 k-NN for Regression

You want to predict a continuous variable Y with a continuous variable X . Having just learned k -NN, you are super eager to try it out for regression. Given the data below, draw the regression lines (what k -NN would predict Y to be for every X value if it was trained for the given data) for k -NN regression with $k=1$, weighted $k=2$, and unweighted $k=2$. For weighted $k=2$, take the weighted average of the two nearest points. For unweighted $k=2$, take the unweighted average of the two nearest points. (Note: the points are equidistant along the x-axis)

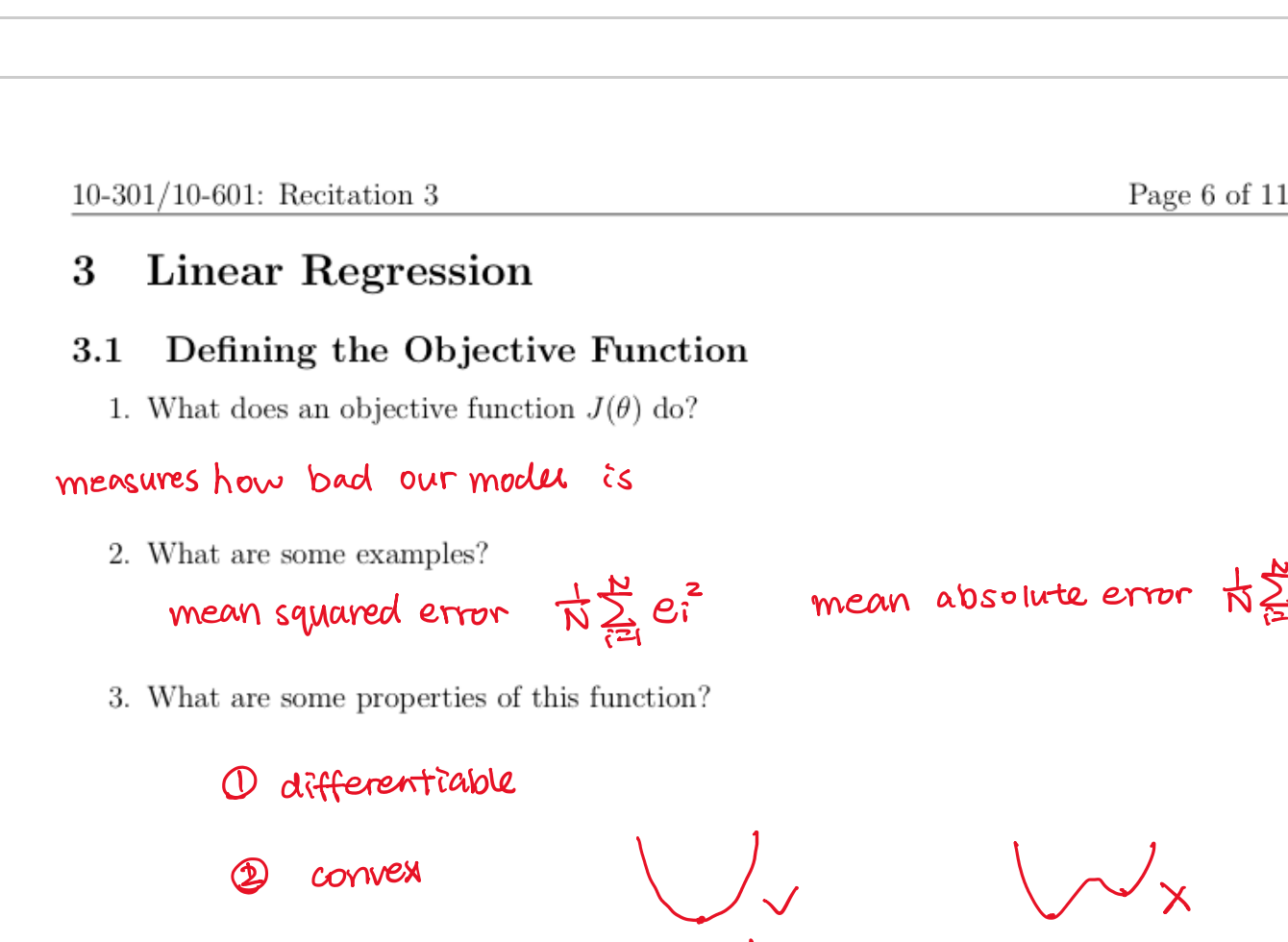
10-301/10-601: Recitation 3 Page 4 of 11



(a) $k=1$, (b) weighted $k=2$, (c) unweighted $k=2$

10-301/10-601: Recitation 3 Page 5 of 11

2.3 k-NN Decision Boundary



Draw the decision boundaries for the above training dataset given using kNN algorithm considering $k=1$.

10-301/10-601: Recitation 3 Page 6 of 11

3 Linear Regression

3.1 Defining the Objective Function

1. What does an objective function $J(\theta)$ do?
 measures how bad our model is

2. What are some examples?
 mean squared error $\frac{1}{n} \sum_{i=1}^n e_i^2$ mean absolute error $\frac{1}{n} \sum_{i=1}^n |e_i|$

3. What are some properties of this function?
 ① differentiable ② convex global minimum

10-301/10-601: Recitation 3 Page 7 of 11

3.2 Solving Linear Regression using Gradient Descent

x_0	1	1	1	1	1
x_1	1.0	2.0	3.0	4.0	5.0
x_2	-2.0	-5.0	-6.0	-8.0	-11.0
y	2.0	4.0	7.0	8.0	11.0

Now, we want to implement the gradient descent method. Assuming that $\gamma = 0.1$ and θ has been initialized to $[0, 0, 0]^T$, perform one iteration of gradient descent:

1. What is the gradient of the objective function $J(\theta)$ with respect to θ : $\nabla_{\theta} J(\theta)$?
 $J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad \frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)}) x_0^{(i)}$$

$$= \frac{2}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} - y^{(i)}) x_0^{(i)}$$

$$= \frac{2}{n} \sum_{i=1}^n (\theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} - y^{(i)})$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{2}{n} (2 \cdot (-2) \cdot 1 + 2 \cdot (0-4) \cdot 1 + 2 \cdot (0-7) \cdot 1 + 2 \cdot (0-8) \cdot 1 + 2 \cdot (0-11) \cdot 1) = -12.8$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{2}{n} (2 \cdot (-2) \cdot 1 + 2 \cdot (0-4) \cdot 2 + 2 \cdot (0-7) \cdot 3 + 2 \cdot (0-8) \cdot 4 + 2 \cdot (0-11) \cdot 5) = -47.2$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = \frac{2}{n} (2 \cdot (-2) \cdot 1 + 2 \cdot (0-4) \cdot 1 + 2 \cdot (0-7) \cdot 3 + 2 \cdot (0-8) \cdot 4 + 2 \cdot (0-11) \cdot 5) = -47.2$$

3. How could we pick which value of γ to use if we weren't given the step size?
 ① γ is small: slow to converge ② γ is large: faster, might overstep global minimum

fail to converge cross-validation hold-out validation

10-301/10-601: Recitation 3 Page 8 of 11

4 Perceptron

4.1 Perceptron Mistake Bound Guarantee

If a dataset has margin γ and all points inside a ball of radius R , then the perceptron makes less than or equal to $(R/\gamma)^2$ mistakes.

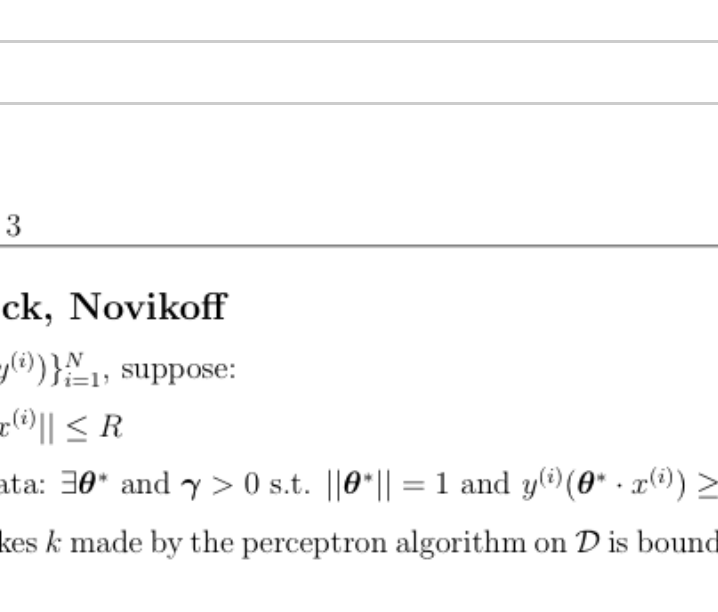


Figure 5: Perceptron Mistake Bound Setup

4.2 Definitions

- Margin:**
 - The margin of example x wrt a linear separator w is the (absolute) distance from x to the plane $w \cdot x = 0$.
 - The margin γ_w of a set of examples S wrt a linear separator w is the smallest margin over points $x \in S$.
 - The margin γ of a set of examples S is the maximum γ_w over all linear separators w .
- Linear Separability:** For a binary classification problem, a set of examples S is linearly separable if there exists a linear decision boundary that can separate the points.
- Update Rule:** When the k -th mistake is made on data point $x^{(k)}$, the parameter update is $\theta^{(k+1)} = \theta^{(k)} + y^{(k)} x^{(k)}$

We say the (batch) perceptron algorithm has *converged* when it stops making mistakes on the training data.

10-301/10-601: Recitation 3 Page 9 of 11

4.3 Perceptron Mistake Bound: Example

Given dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, suppose:
 1. Finite size inputs: $\|x^{(i)}\| \leq R$
 2. Linearly separable data: $\exists \theta^*$ and $\gamma > 0$ s.t. $\|\theta^*\| = 1$ and $y^{(i)}(\theta^* \cdot x^{(i)}) \geq \gamma, \forall i$

Then, the number of mistakes k made by the perceptron algorithm on \mathcal{D} is bounded by $(R/\gamma)^2$. The following table shows a dataset of linearly separable datapoints.

x_1	x_2	y
1	-1	1
0	2	-1
4	0	1

Assuming that the linear separator with the largest margin is given by:
 $\theta^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$, where $\theta = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

Calculate the theoretical mistake bound for the perceptron.
 $R \rightarrow \sqrt{(4-0)^2 + (0-0)^2} = 4$
 $\gamma \rightarrow \frac{|\theta^T x|}{\|\theta\|_2} = \frac{|[-1 \ 1] \cdot [1 \ -1]^T|}{\sqrt{(-1)^2 + 1^2}} = \frac{|-2|}{\sqrt{2}} = \frac{2}{\sqrt{2}} = \sqrt{2}$
 $(\frac{R}{\gamma})^2 = (\frac{4}{\sqrt{2}})^2 = \frac{16}{2} = 8$
 $\theta = [b, \theta_1, \theta_2]$
 $x = [x_1, x_2]$

10-301/10-601: Recitation 3 Page 10 of 11

4.4 Theorem: Block, Novikoff

Given dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n$, suppose:
 1. Finite size inputs: $\|x^{(i)}\| \leq R$
 2. Linearly separable data: $\exists \theta^*$ and $\gamma > 0$ s.t. $\|\theta^*\| = 1$ and $y^{(i)}(\theta^* \cdot x^{(i)}) \geq \gamma, \forall i$

Then, the number of mistakes k made by the perceptron algorithm on \mathcal{D} is bounded by $(R/\gamma)^2$.
Proof:
 Part 1: For some $A, Ak \leq \|\theta^{(k+1)}\|$

Part 2: For some $B, \|\theta^{(k+1)}\| \leq B\sqrt{k}$

Part 3: Combine the bounds

Main Takeaway:

10-301/10-601: Recitation 3 Page 11 of 11

5 Summary

5.1 Decision Tree

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and interpret Very fast for inference 	<ul style="list-style-type: none"> Tree may grow very large and tend to overfit. Greedy behaviour may be sub-optimal 	<ul style="list-style-type: none"> Prefer the smallest tree consistent w/ the training data (i.e. 0 error rate) 	<ul style="list-style-type: none"> Most cases. Random forests are widely used in industry.

5.2 k-NN

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> No training of parameters Can apply to multi-class problems and use different metrics 	<ul style="list-style-type: none"> Slow for large datasets Must select good k Imbalanced data and outliers can lead to misleading results 	<ul style="list-style-type: none"> Similar (i.e. nearby) points should have similar labels All label dimensions are created equal 	<ul style="list-style-type: none"> Small dataset Small dimensionality Data is clean (no missing data) Inductive bias is strong for dataset

5.3 Linear regression

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and train Closed form solution 	<ul style="list-style-type: none"> Sensitive to noise (other than zero-mean Gaussian noise) 	<ul style="list-style-type: none"> The true relationship between the inputs and output is linear. 	<ul style="list-style-type: none"> Most cases (can be extended by adding non-linear feature transformations)

5.4 Perceptron

Pros	Cons	Inductive bias	When to use
<ul style="list-style-type: none"> Easy to understand and works for online learning. Provable guarantees on mistakes made for linearly separable data. 	<ul style="list-style-type: none"> No guarantees on finding best (maximum-margin) hyperplane. Output is sensitive to noise in the training data. 	<ul style="list-style-type: none"> The binary classes are separable in the feature space by a line. 	<ul style="list-style-type: none"> Not used much anymore, but variants (kernel perceptron, structured perceptron) may have more success.