

Warm-up as you log in

What are these linear shapes called for 1-D, 2-D, 3-D, M-D input?

$$\mathbf{x} \in \mathbb{R}$$

$$\mathbf{x} \in \mathbb{R}^2$$

$$\mathbf{x} \in \mathbb{R}^3$$

$$\mathbf{x} \in \mathbb{R}^M$$

$$y = \mathbf{w}^T \mathbf{x} + b$$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$$\mathbf{w}^T \mathbf{x} + b \geq 0$$

Announcements

Assignments

- HW2
 - Due Mon, 9/21, 11:59 pm
- HW3
 - Out tomorrow, due Mon, 9/28, 11:59 pm
 - Written, but in Gradescope

Midterm 1

- Mon, 10/5
- In lecture; Gradescope exam (like HW1 written); proctored via Zoom
- Content up to and including linear regression and optimization
- Stay tuned to Piazza for details and a few forms to fill out

Plan

Last time

- Model selection
 - Parameters, Hyperparameters
 - Train, Test, and Validation sets

Today

- A few more things on model selection
- Regression
- Linear regression
- Optimization for linear regression



Introduction to Machine Learning

Linear Regression and Optimization

Instructor: Pat Virtue

Model Selection

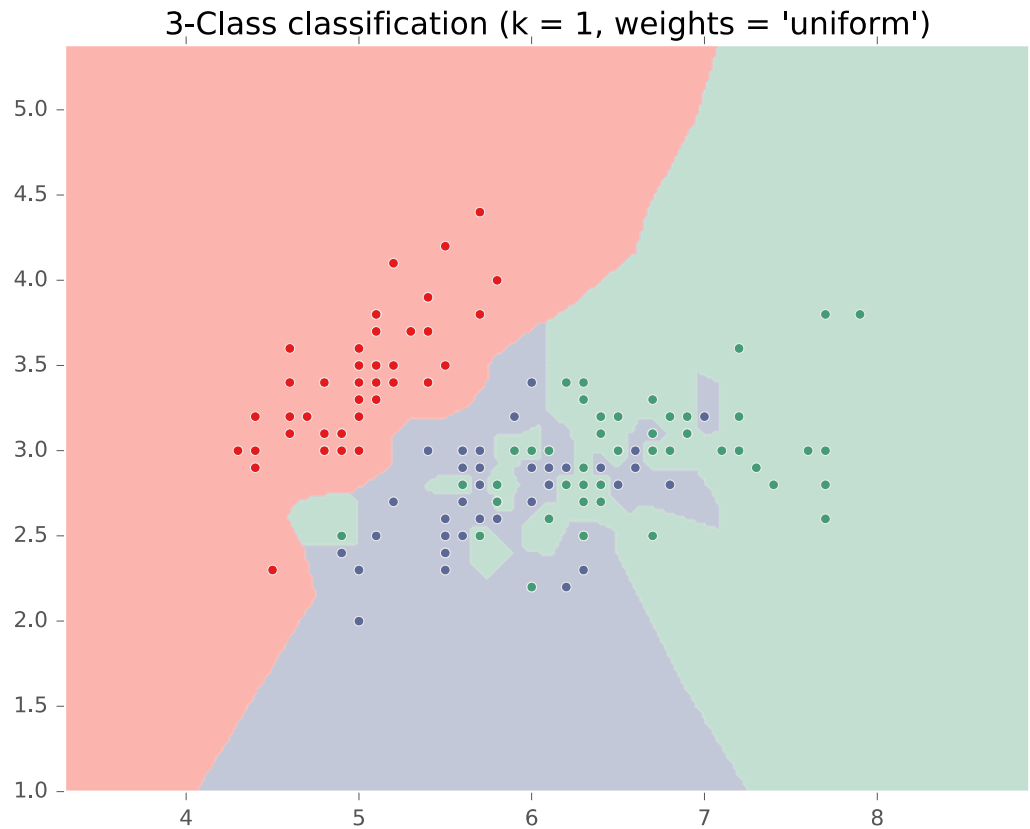
- Two very similar definitions:
 - Def: **model selection** is the process by which we choose the “best” model from among a set of candidates
 - Def: **hyperparameter optimization** is the process by which we choose the “best” hyperparameters from among a set of candidates (**could be called a special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

Experimental Design

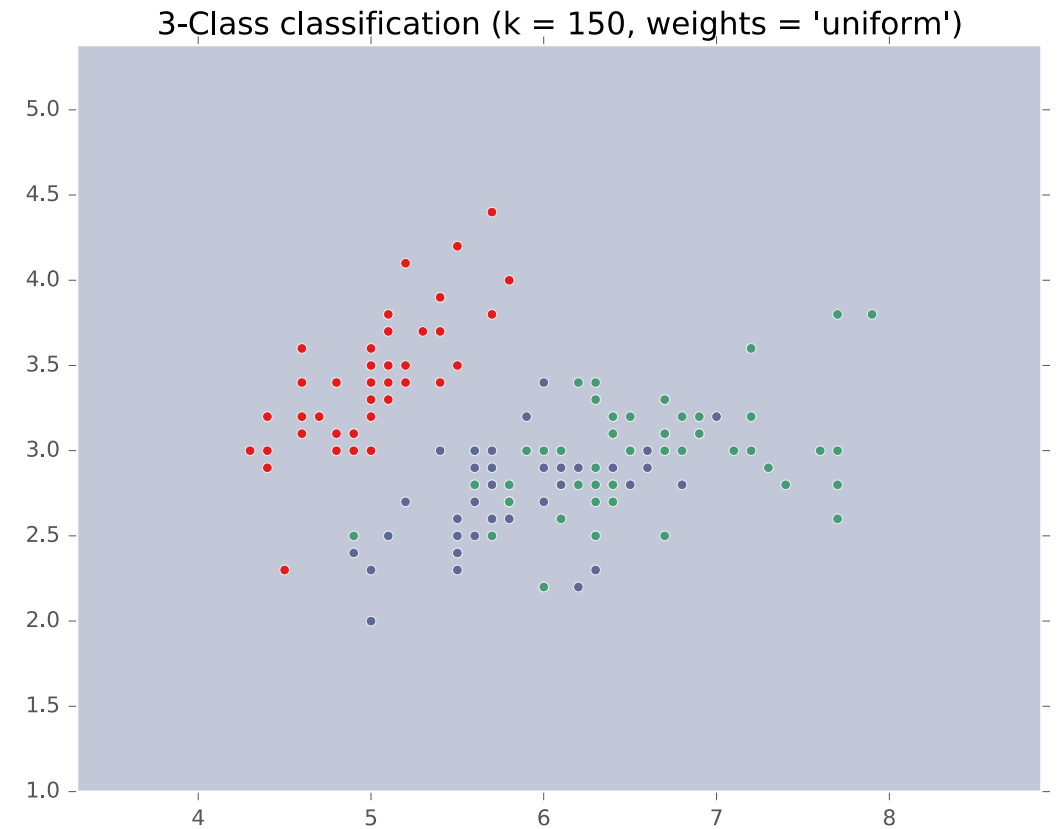
| | Input | Output | Notes |
|-----------------------------|---|---|--|
| Training | <ul style="list-style-type: none">• training dataset• hyperparameters | <ul style="list-style-type: none">• best model parameters | We pick the best model parameters by learning on the training dataset for a fixed set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

Special Cases of k-NN

k=1: Nearest Neighbor

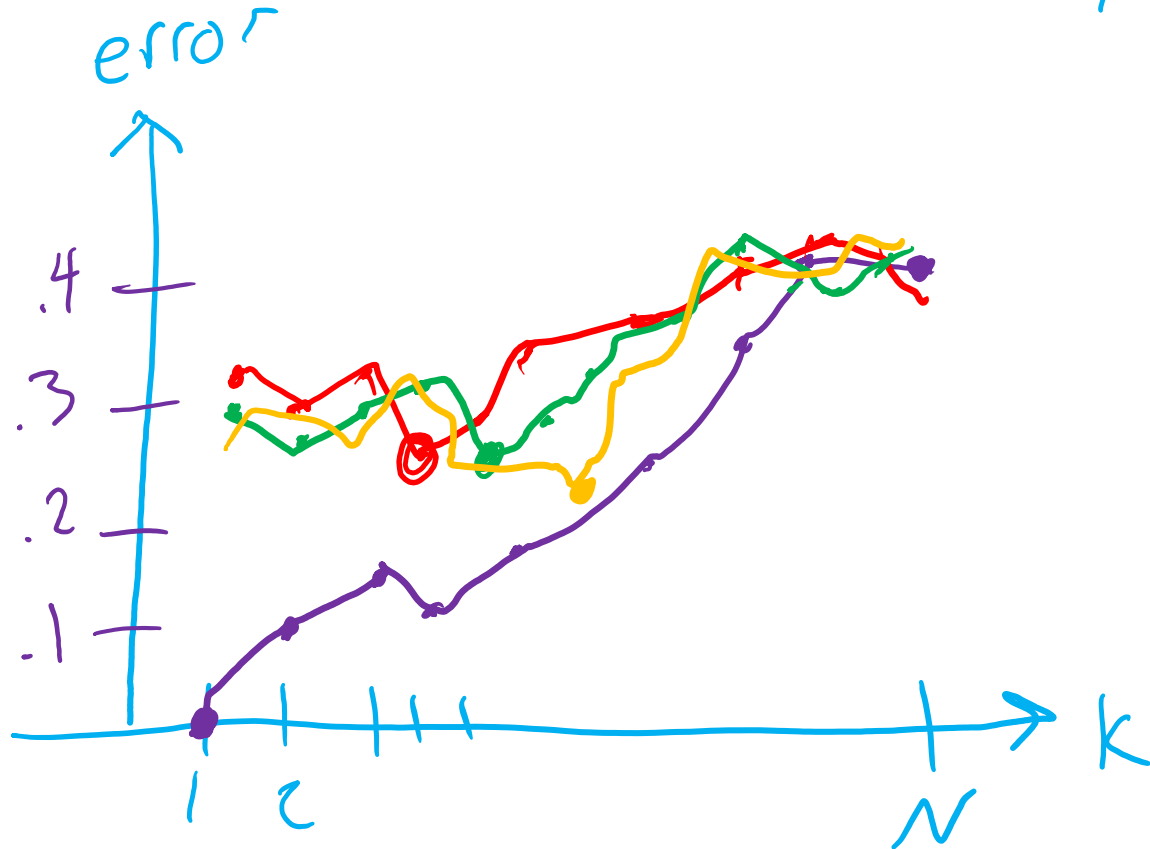


k=N: Majority Vote



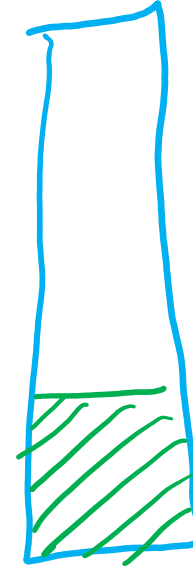
Example of Hyperparameter Optimization

Choosing k for k -NN



$\gamma = 0$ 40% $\gamma = 1$ 60%

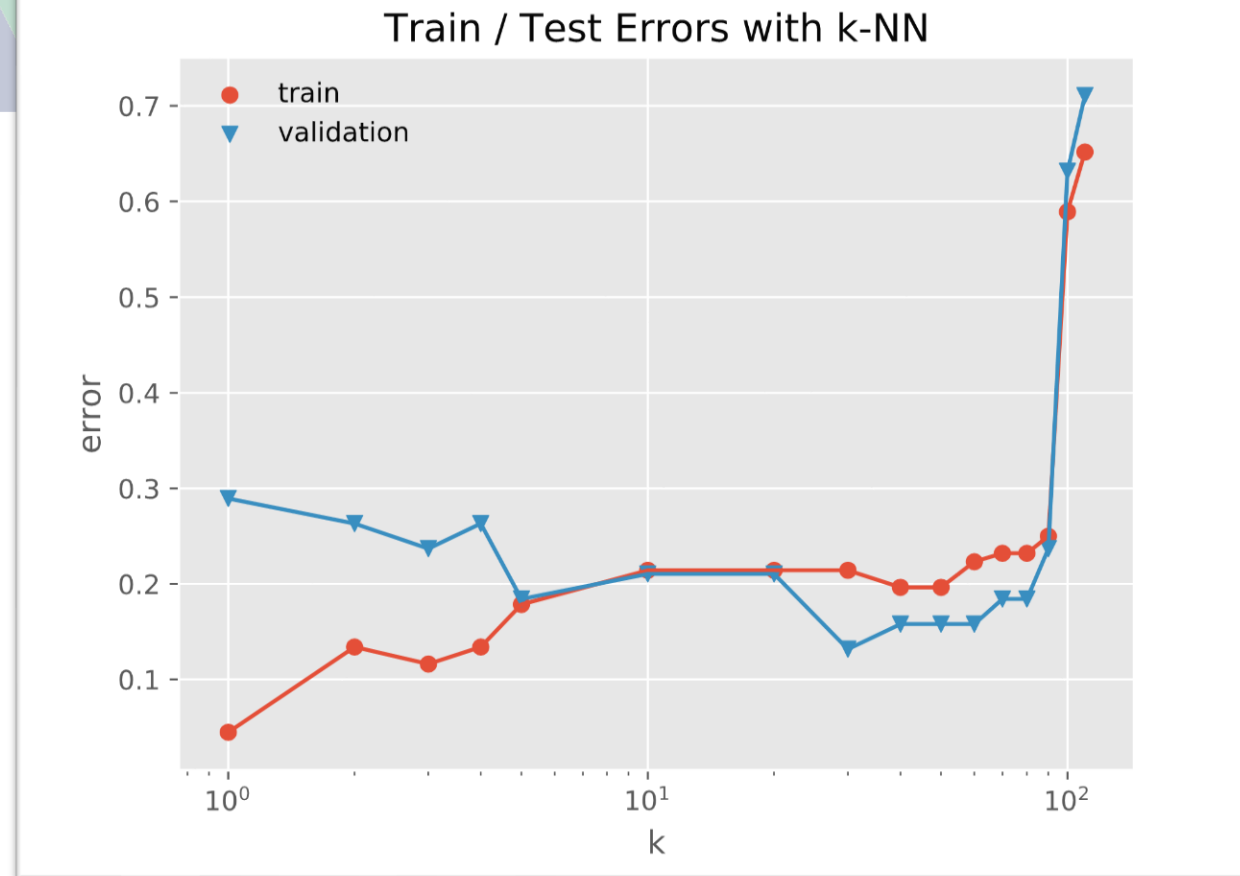
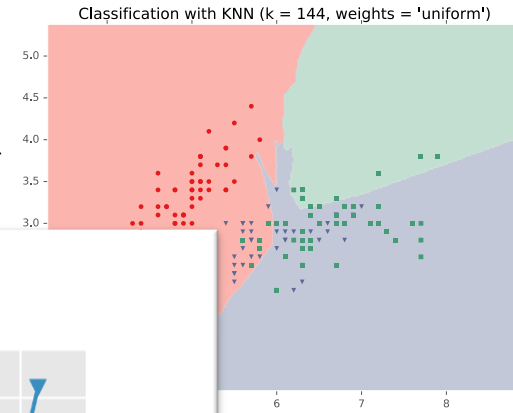
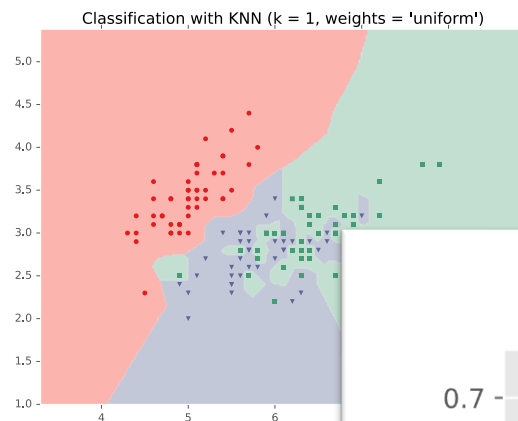
1200



Val 300

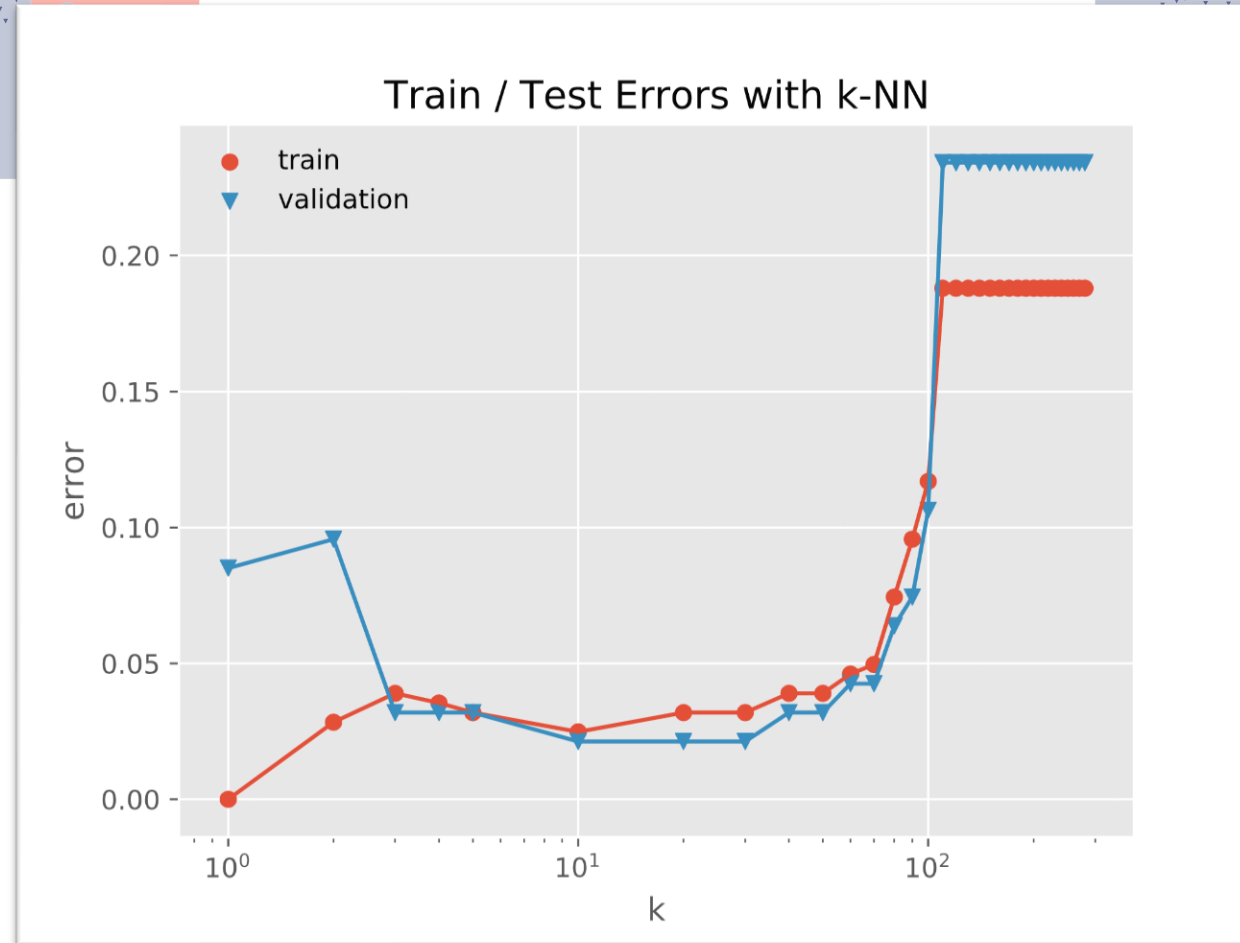
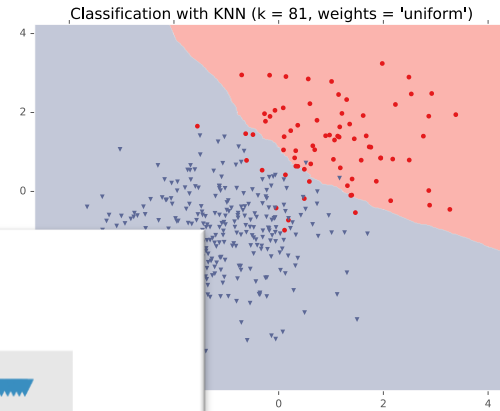
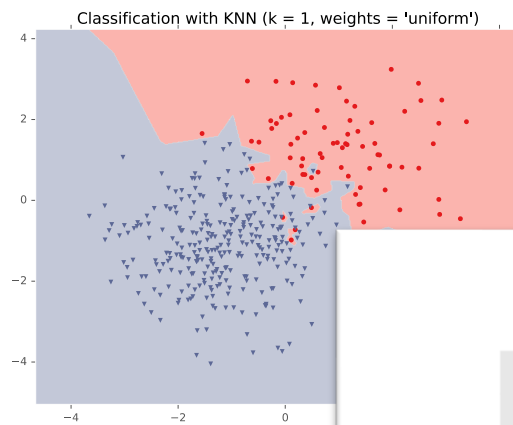


k-NN: Choosing k



Fisher Iris Data: varying the value of k

k-NN: Choosing k



Gaussian Data: varying the value of k

Validation

Why do we need validation?

- Choose hyperparameters
- Choose technique
- Help make any choices beyond our parameters

But now, we have another choice to make!

- How do we split training and validation?

Trade-offs

- More held-out data, better meaning behind validation numbers
- More held-out data, less data to train on!

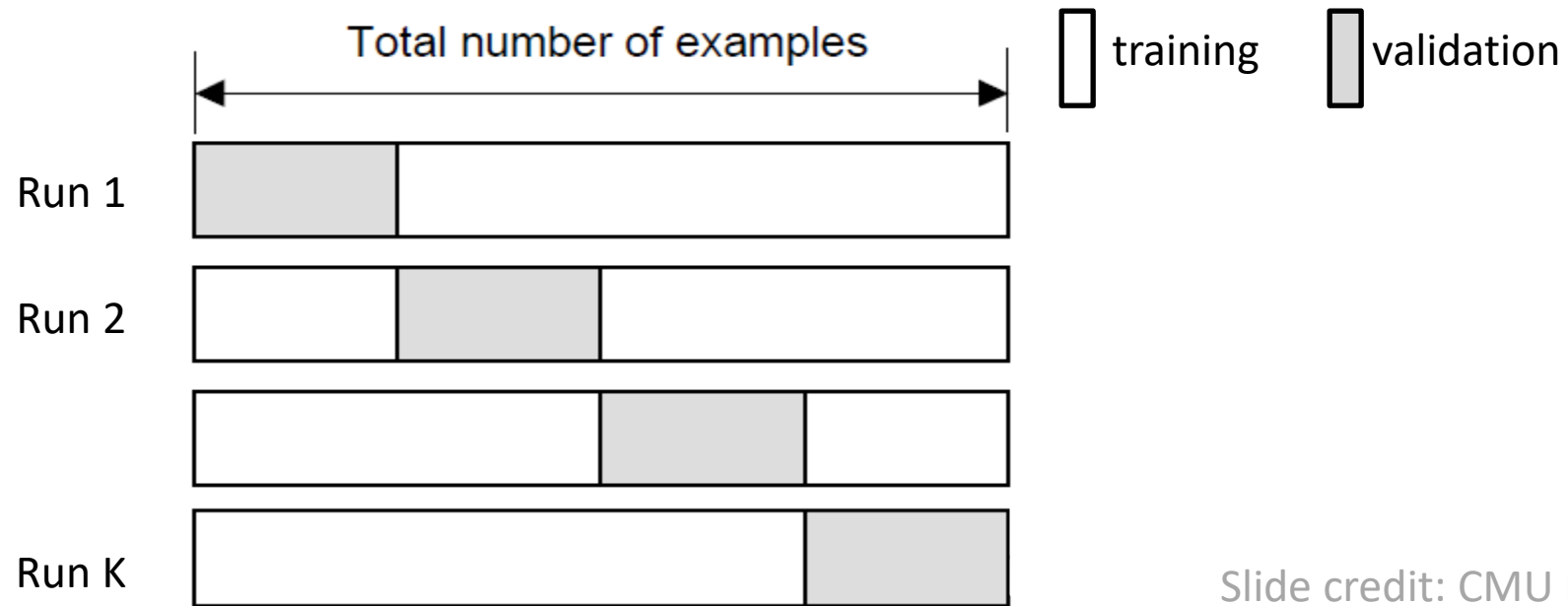
Cross-validation

K-fold cross-validation

Create K-fold partition of the dataset.

Do K runs: train using K-1 partitions and calculate validation error on remaining partition (rotating validation partition on each run).

Report average validation error

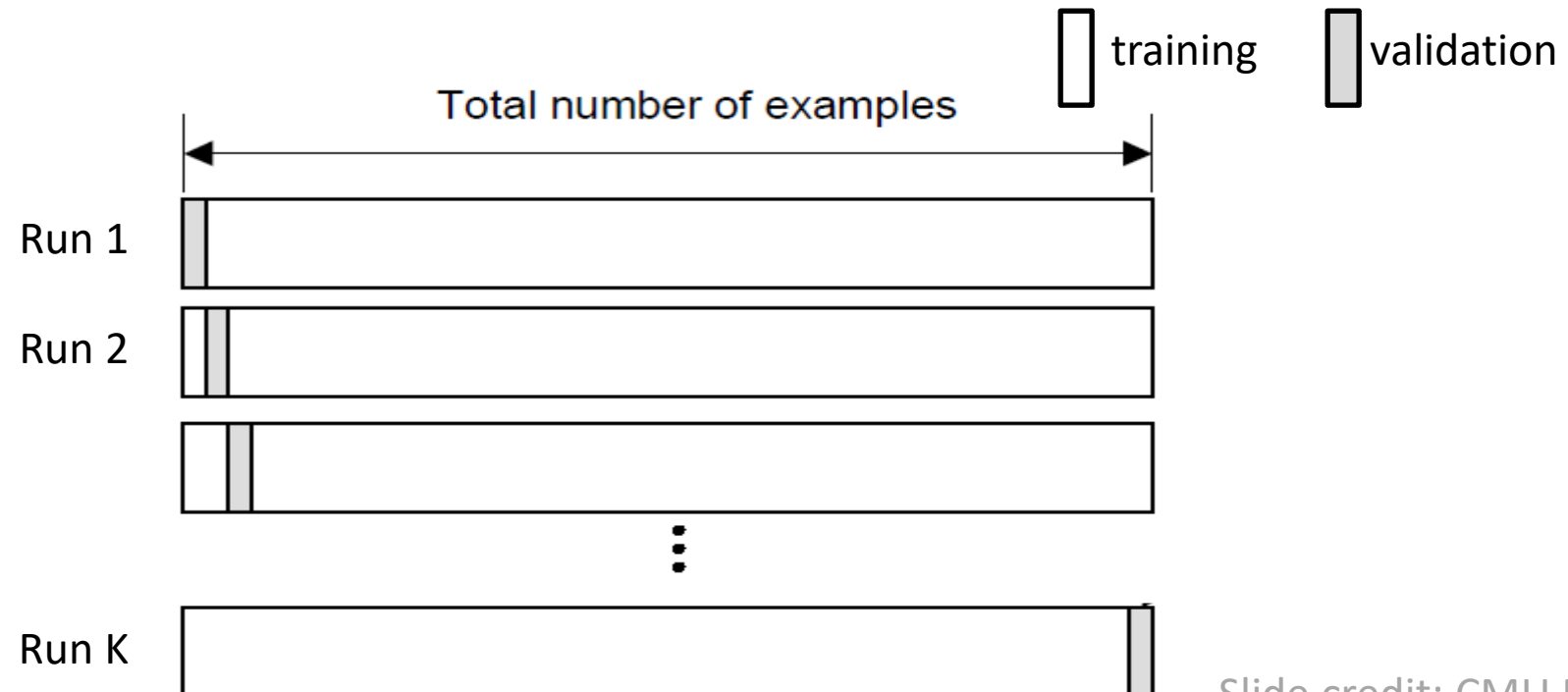


Cross-validation

Leave-one-out (LOO) cross-validation

Special case of K-fold with $K=N$ partitions

Equivalently, train on $N-1$ samples and validate on only one sample per run for N runs



Cross-validation

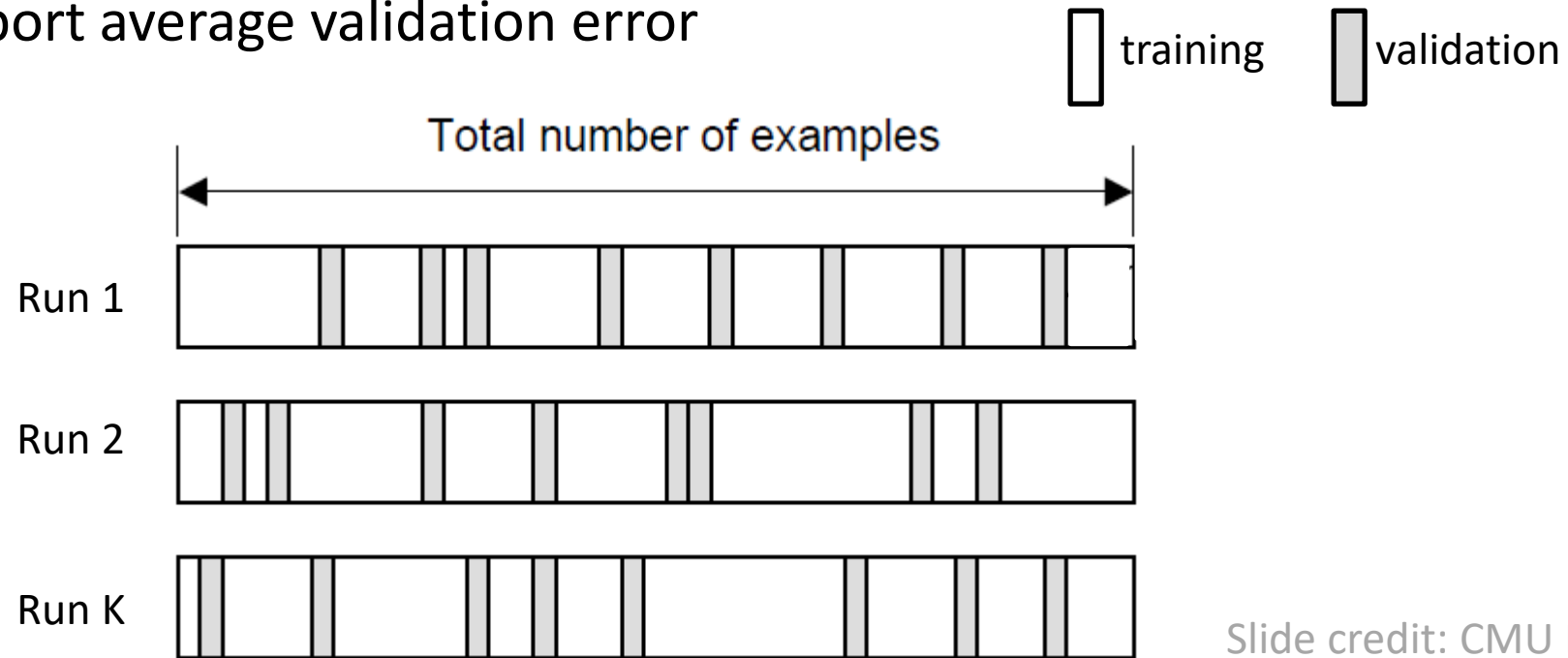
Random subsampling

Randomly subsample a fixed fraction αN ($0 < \alpha < 1$) of the dataset for validation.

Compute validation error with remaining data as training data.

Repeat K times

Report average validation error



Practical Issues in Cross-validation

How to decide the values for K and α ?

- Large K
 - + Validation error can approximate test error well
 - Observed validation error will be unstable (few validation pts)
 - The computational time will be very large as well (many experiments)
- Small K
 - + The # experiments and, therefore, computation time are reduced
 - + Observed validation error will be stable (many validation pts)
 - Validation error cannot approximate test error well

Common choice: $K = 10$, $\alpha = 0.1$ 😊

Piazza Poll 1

Say you are choosing amongst 10 discrete values of a decision tree *mutual information threshold*, and you want to do K=10-fold cross-validation.

How many times do I have to train my model?

- A. 0
- B. 1
- C. 10
- D. 20
- E. 100
- F. 10^{10}

Piazza Poll 1

Say you are choosing amongst 10 discrete values of a decision tree *mutual information threshold*, and you want to do K=10-fold cross-validation.

How many times do I have to train my model?

- A. 0
- B. 1
- C. 10
- D. 20
- E. 100**
- F. 10^{10}

Model Selection

WARNING (again):

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

Main Takeaway:

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

Model Selection Learning Objectives

You should be able to...

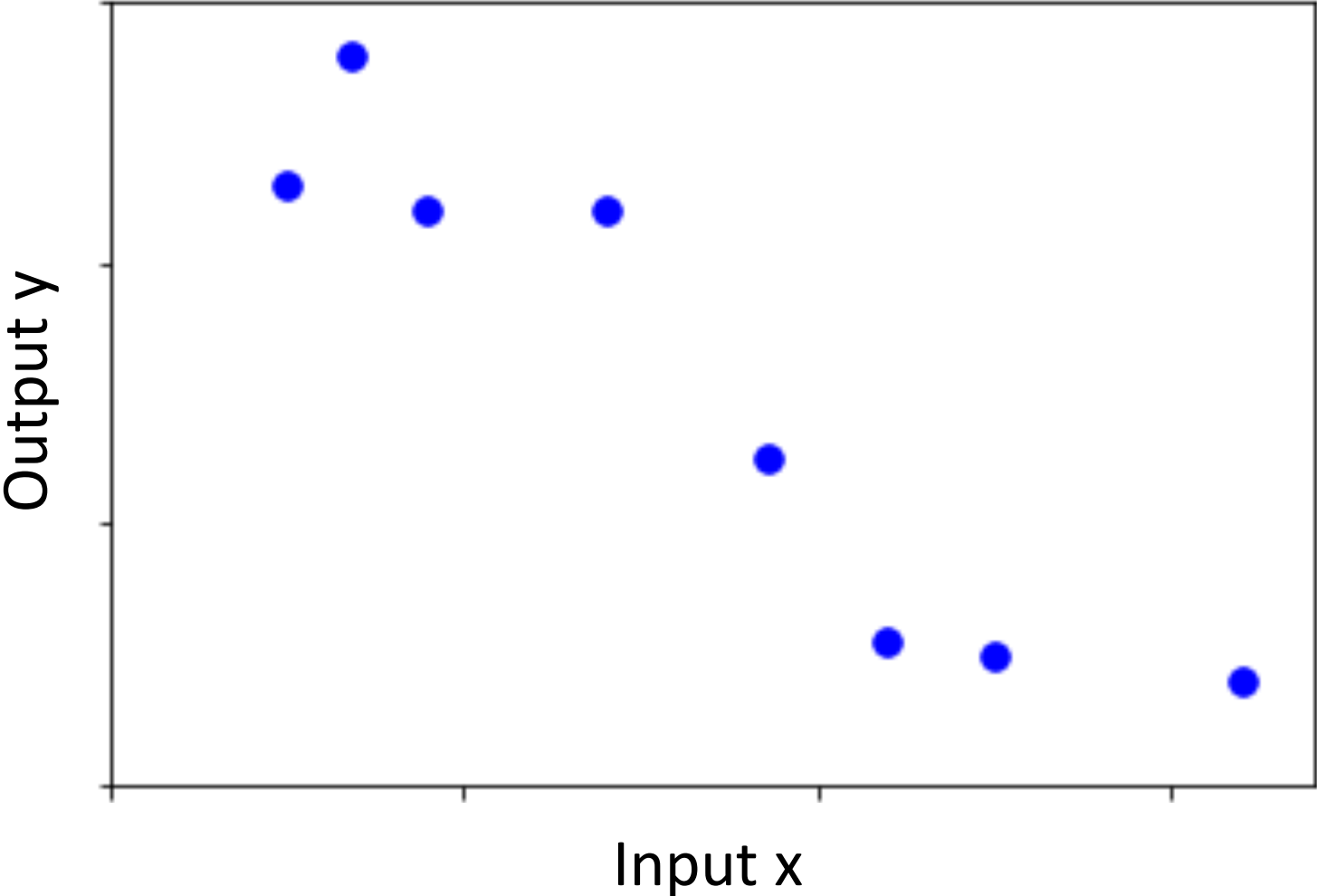
- Plan an experiment that uses training, validation, and test datasets to predict the performance of a classifier on unseen data (without cheating)
- Explain the difference between (1) training error, (2) validation error, (3) cross-validation error, (4) test error, and (5) true error
- For a given learning technique, identify the model, learning algorithm, parameters, and hyperparameters

LINEAR REGRESSION AND OPTIMIZATION

Breakout Room

In your breakout room

Come up with a story for this data



Lecture 2: Problem Formulation

Experience $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ $x \in \mathcal{X}$ $y \in \mathcal{Y}$ } Train Data
Test Data

Regression: $y \in \mathbb{R}$

output, class label, class, label
input, features, attributes

Classification: y discrete (and not ordered) set

Hypothesis

function: $\hat{y} = h(x)$ $h: \mathcal{X} \rightarrow \mathcal{Y}$

Hypothesis space $h \in \mathcal{H}$

$$\hat{y} = \underline{m}x + \underline{b}$$

Performance measure

Error

Loss

Objective function

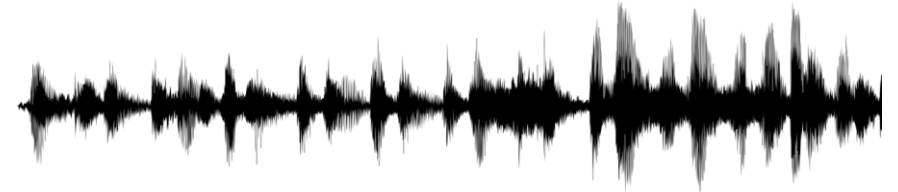
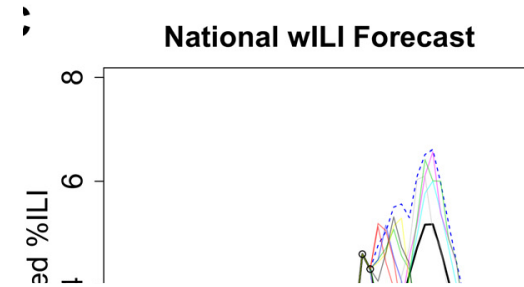
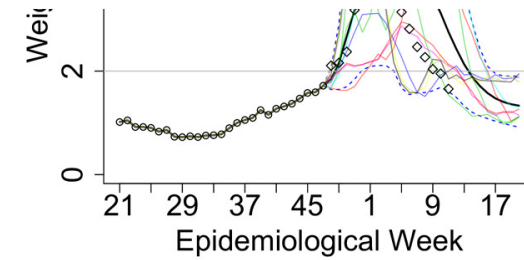
Regression

Goal:

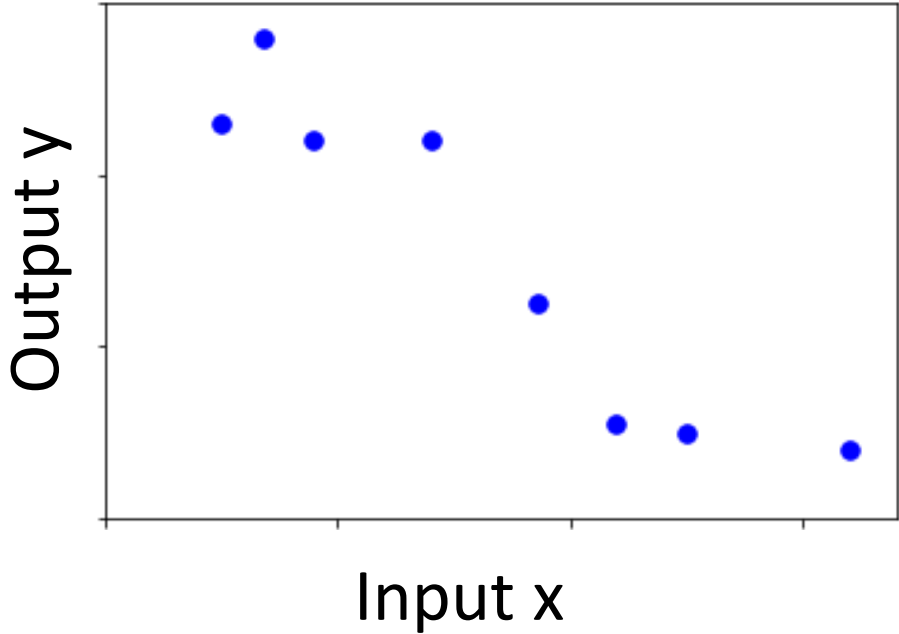
- Given a training dataset of pairs (x,y) where
 - y is a continuous, rather than a label
- Learn a function (aka. curve or line)
 $\hat{y} = h(x)$ that best fits the training data

Example Applications:

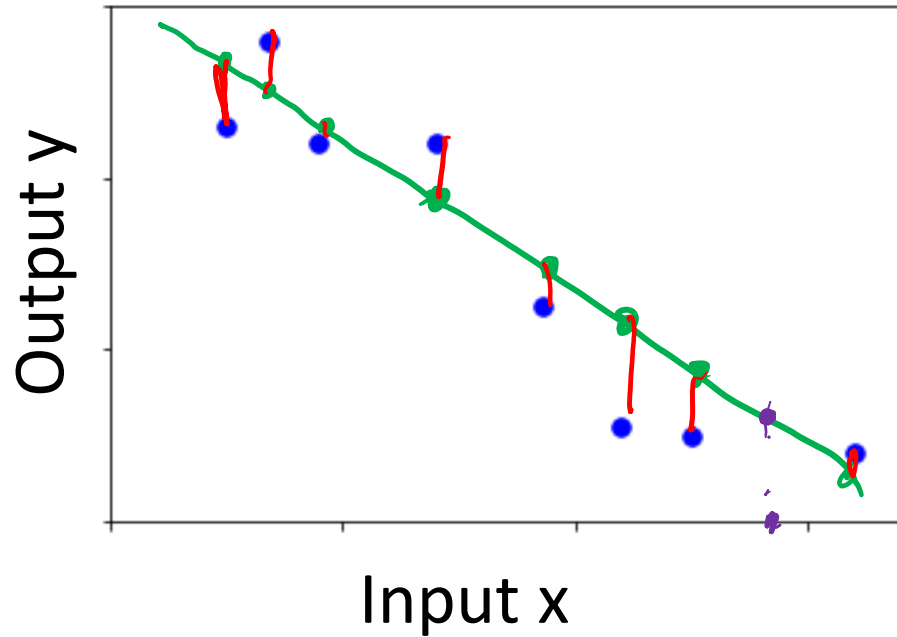
- Stock price prediction
- Forecasting epidemics
- Speech synthesis
- Generation of images (e.g. *Deep Fake*)



Lecture 2: Problem Formulation



Lecture 2: Problem Formulation



Regression

$$x \in \mathbb{R}$$

$$y \in \mathbb{R}$$

$$\hat{y} = h(x) = \underline{w}x + \underline{b}$$

$$\hat{y}^{(i)} = h(x^{(i)})$$

$$\begin{aligned} 3x + 2 \\ 4x - 5 \\ -32x + 7 \end{aligned}$$

Error $y - \hat{y}$

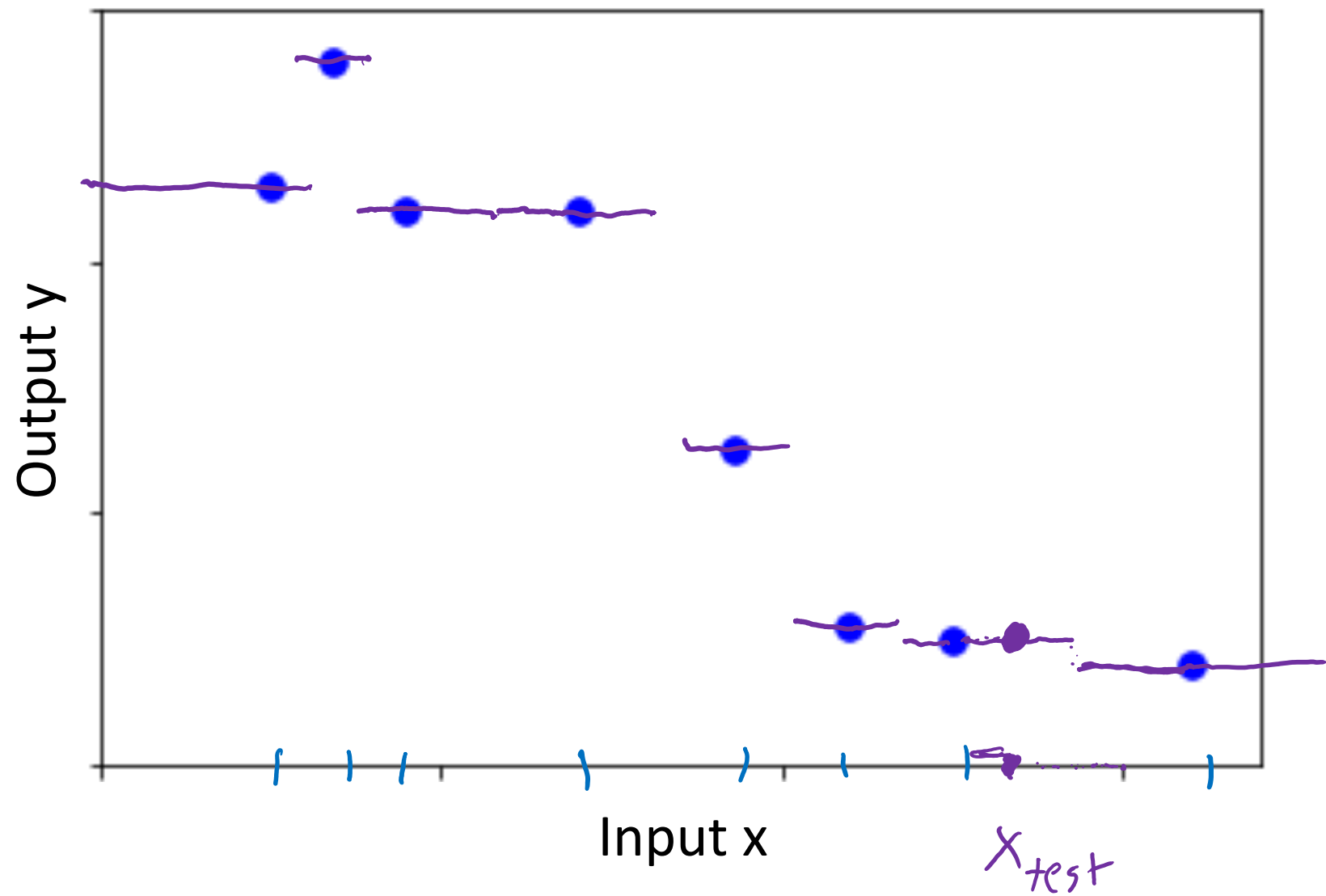
$$\sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)}) \rightarrow \sum |y^{(i)} - \hat{y}^{(i)}| \rightarrow \sum (y^{(i)} - \hat{y}^{(i)})^2$$

Sum sq error

$$\frac{1}{N} \sum (y^i - \hat{y}^i)^2$$

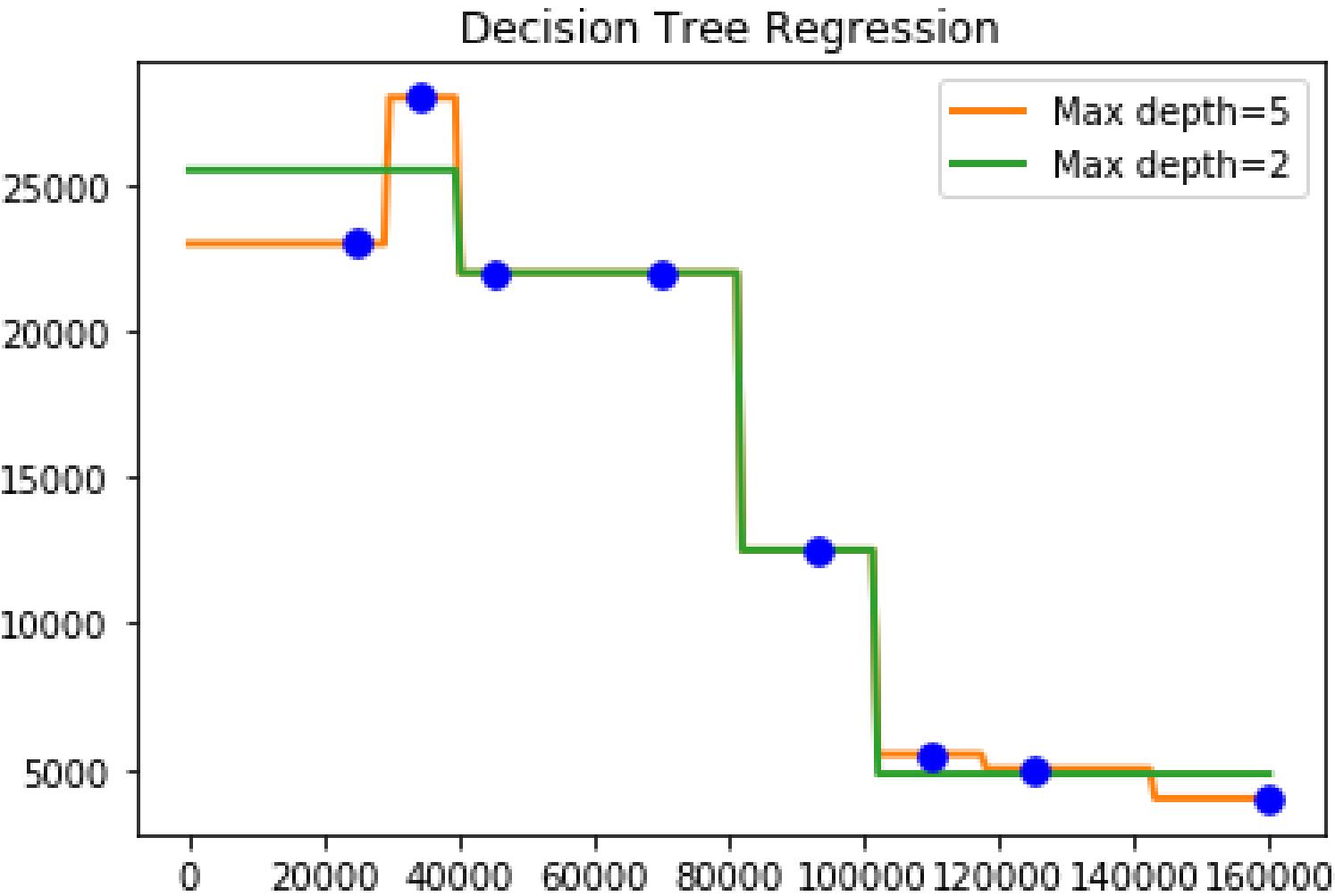
mean sq error

Regression: Nearest Neighbor



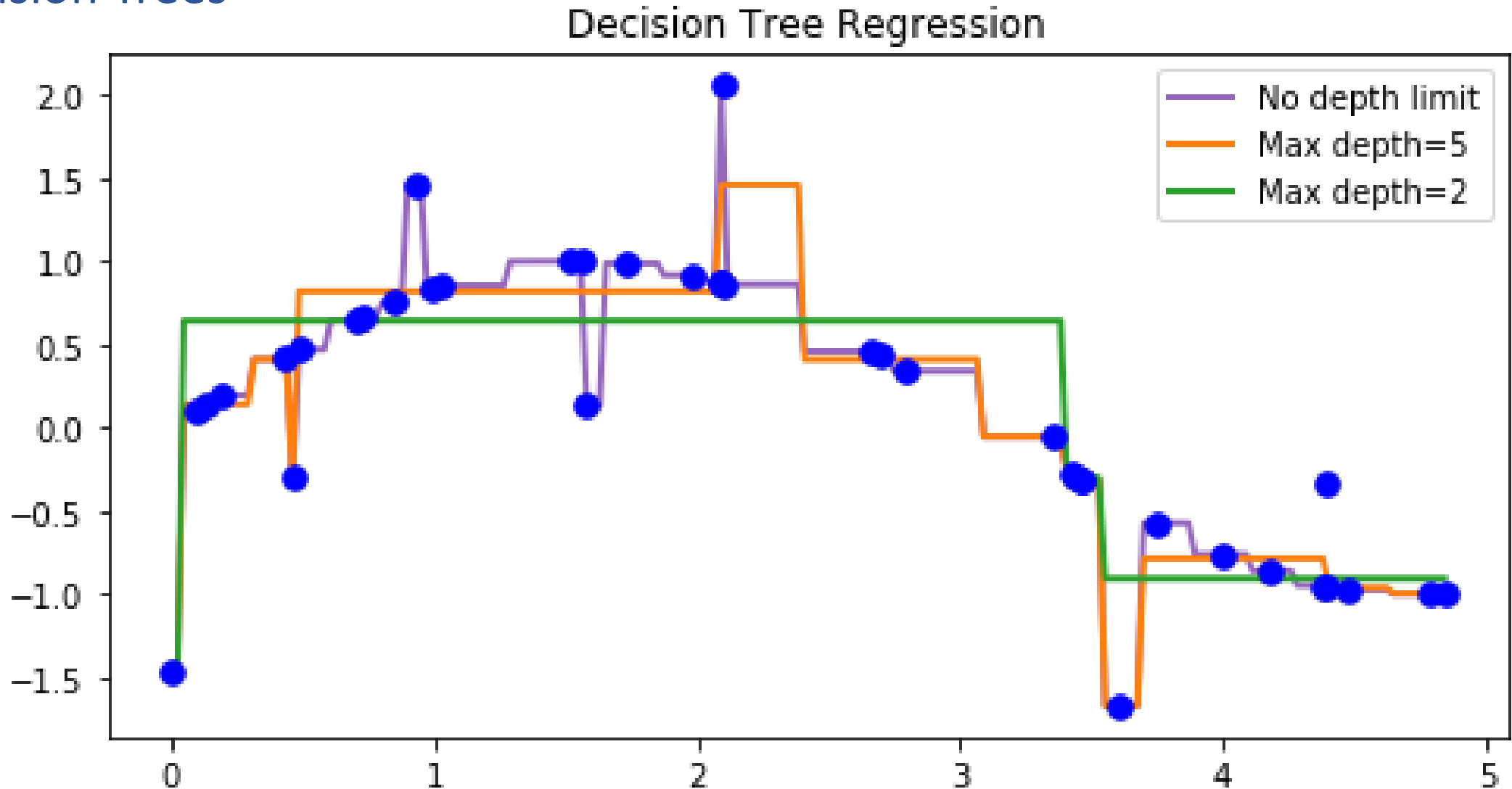
Regression

Decision Trees



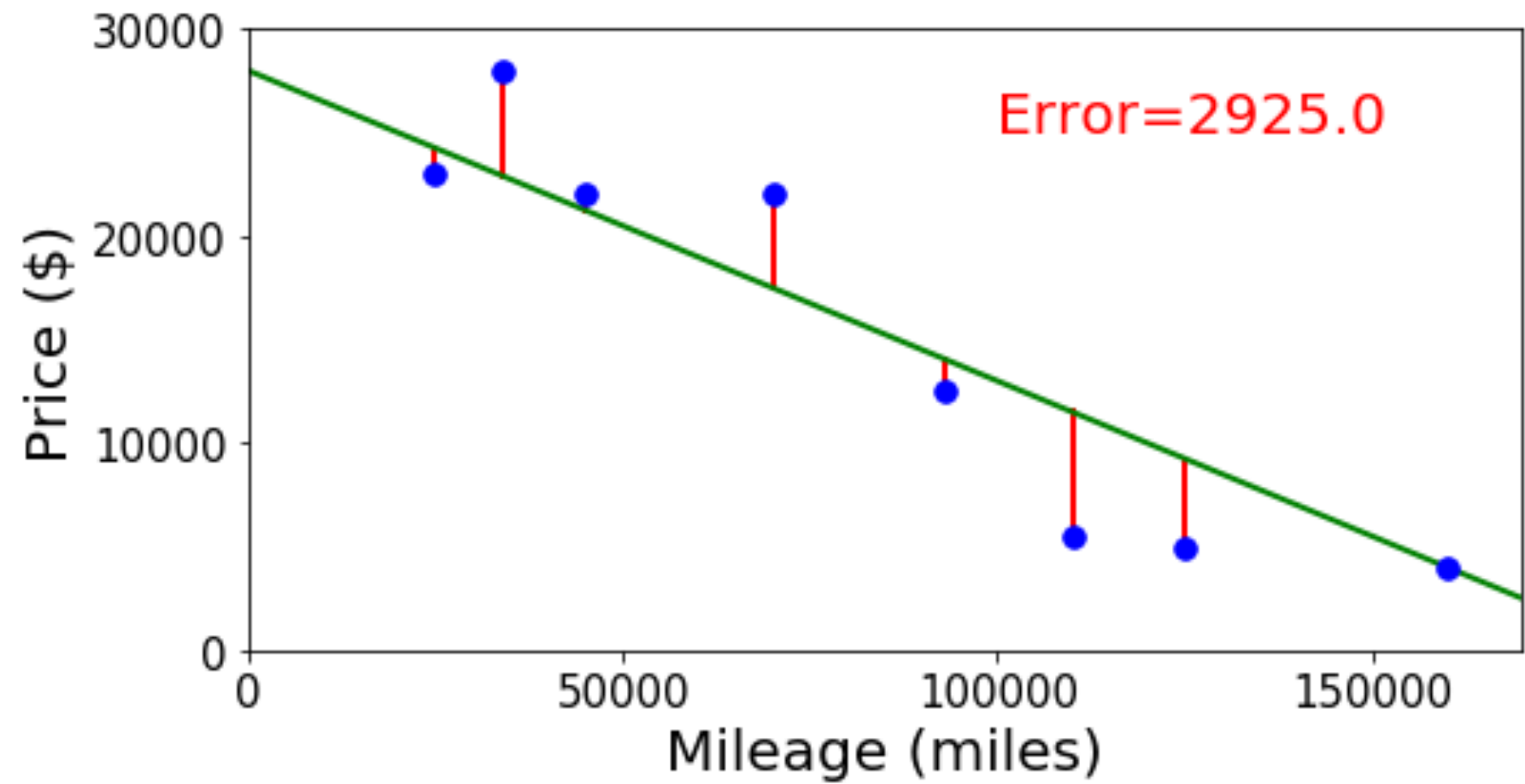
Nonparametric Regression

Decision Trees



Linear Regression

Selling my car



Linear Function

Linear function

If $f(\mathbf{x})$ is linear, then:

- $f(\mathbf{x} + \mathbf{z}) = f(\mathbf{x}) + f(\mathbf{z})$
- $f(\alpha \mathbf{x}) = \alpha f(\mathbf{x}) \quad \forall \alpha$
- $f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{z}) = \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z}) \quad \forall \alpha$

Linear in Higher Dimensions

What are these linear shapes called for 1-D, 2-D, 3-D, M-D input?

$$\mathbf{x} \in \mathbb{R}$$

$$\mathbf{x} \in \mathbb{R}^2$$

$$\mathbf{x} \in \mathbb{R}^3$$

$$\mathbf{x} \in \mathbb{R}^M$$

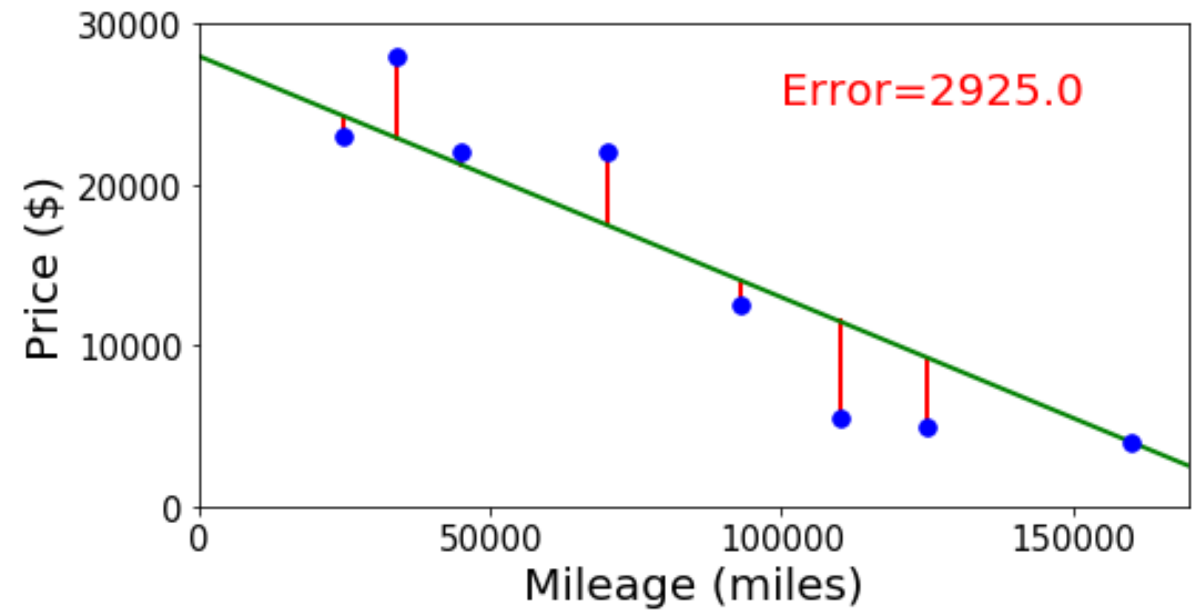
$$y = \mathbf{w}^T \mathbf{x} + b$$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$$\mathbf{w}^T \mathbf{x} + b \geq 0$$

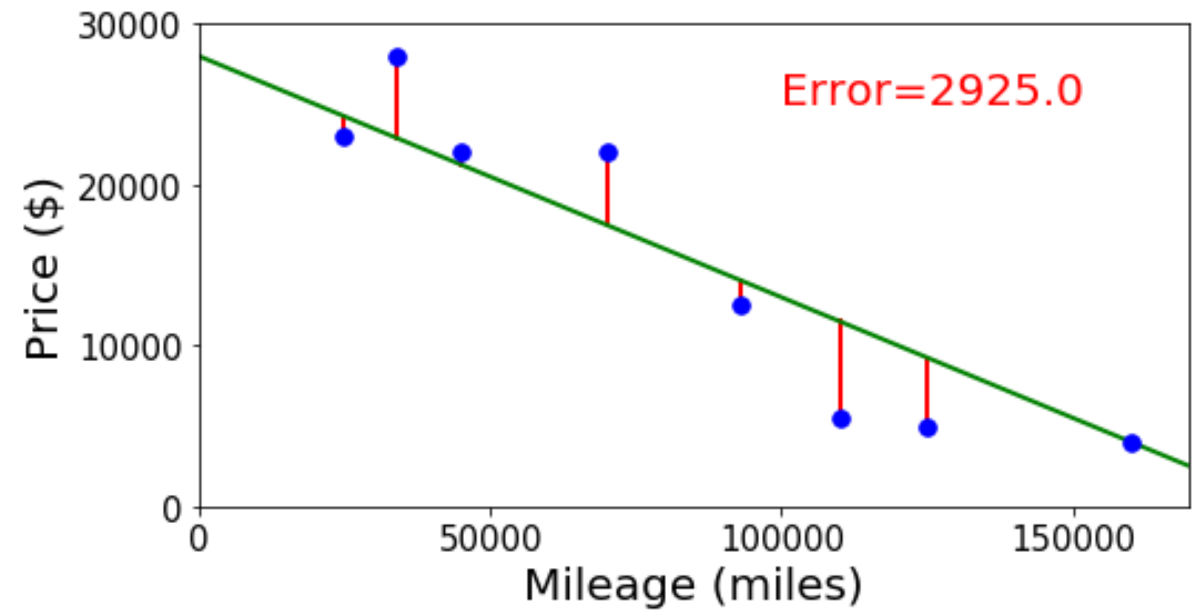
Linear Regression

Linear algebra formulation



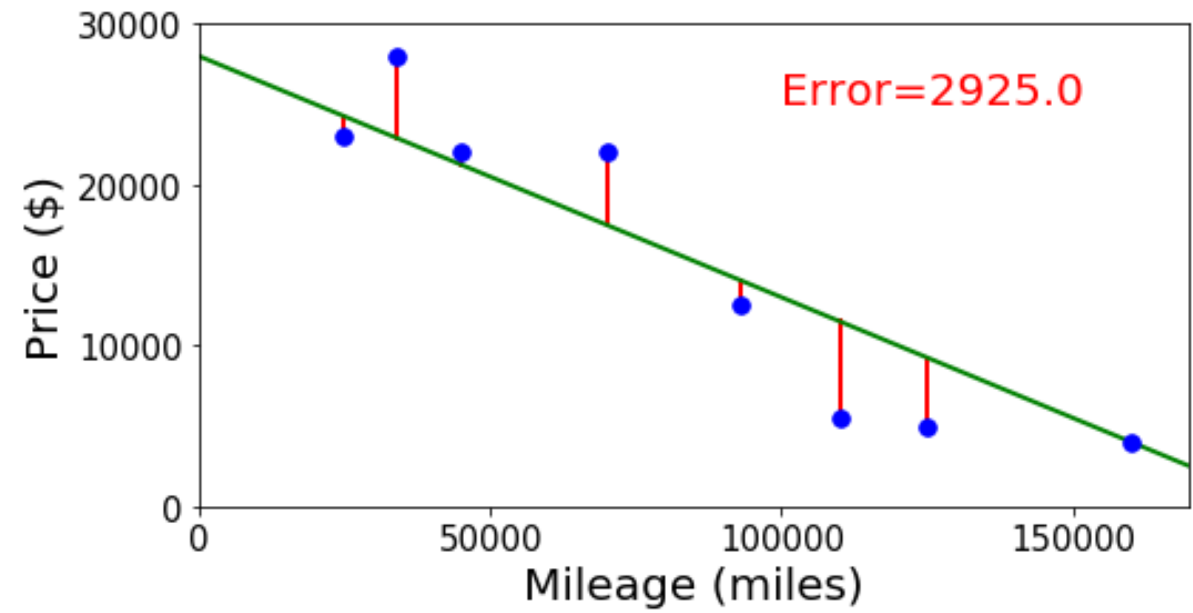
Linear Regression

Error and objectives



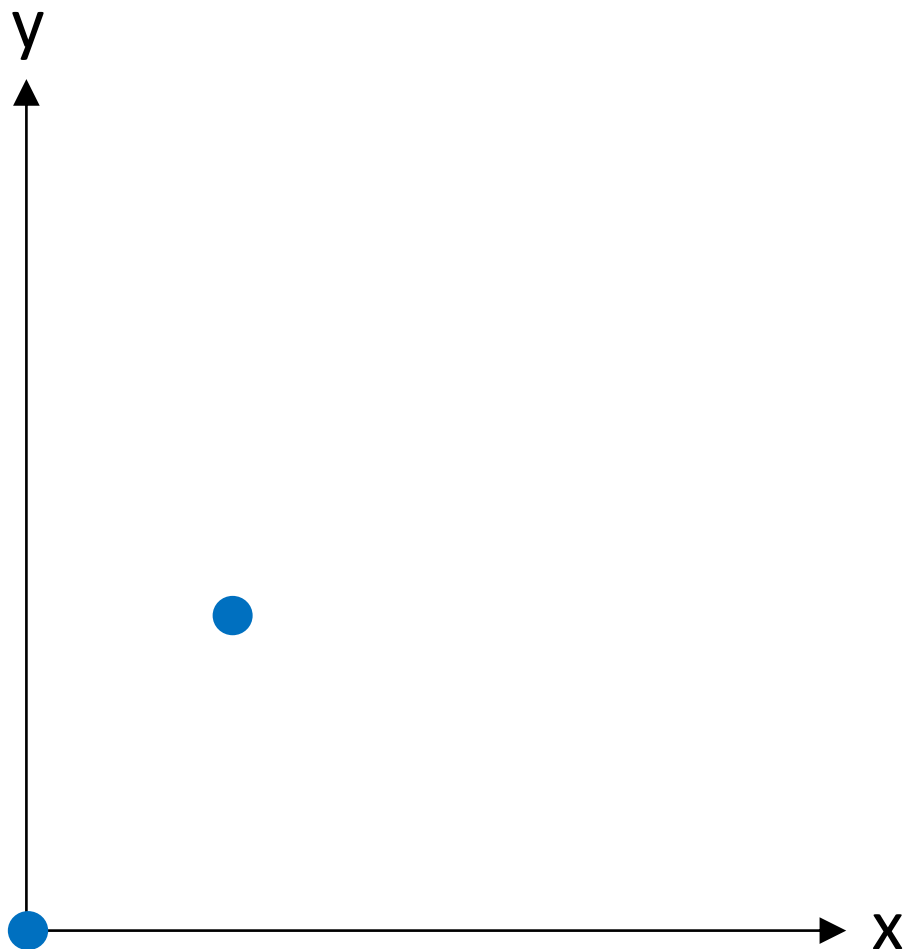
Linear Regression

Linear algebra formulation



Linear Regression

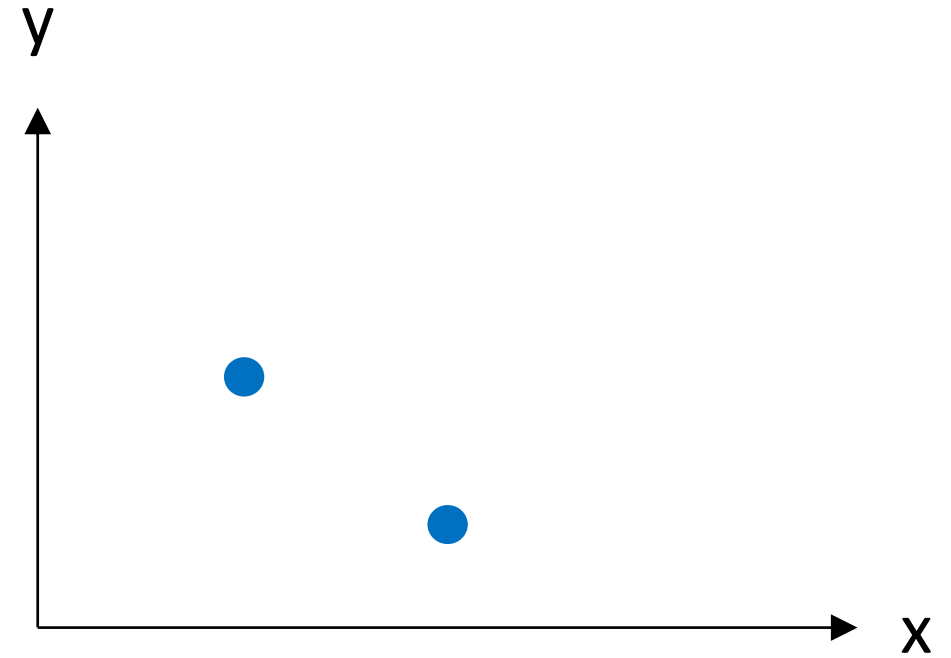
Optimizing the objective



Piazza Poll 2

For fixed data and fixed slope, w , what shape do we get by plotting MSE objective vs intercept, b ?

- A. Line
- B. Plane
- C. Half-plane
- D. Convex Parabola (U-shape)
- E. Concave parabola (up-side-down U)
- F. None of the above



Piazza Poll 2

For fixed data and fixed slope, w , what shape do we get by plotting MSE objective vs intercept, b ?

- A. Line
- B. Plane
- C. Half-plane
- D. Convex Parabola (U-shape)
- E. Concave parabola (up-side-down U)
- F. None of the above

