

10-701: Introduction to Machine Learning Lecture 4 – Linear Regression

Henry Chai

1/29/24

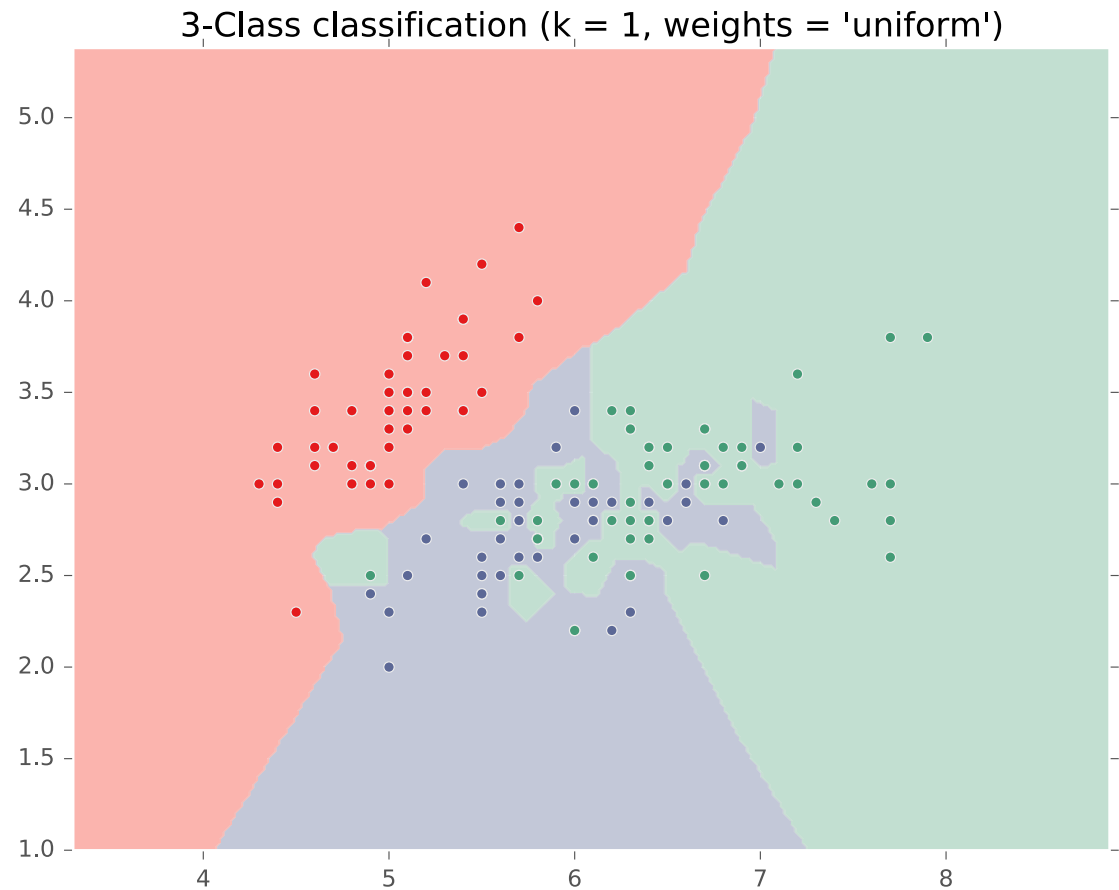
Front Matter

- Announcements:
 - HW1 released 1/24, due 2/2 at 11:59 PM
- Recommended Readings:
 - Murphy, [Sections 7.1-7.3](#)

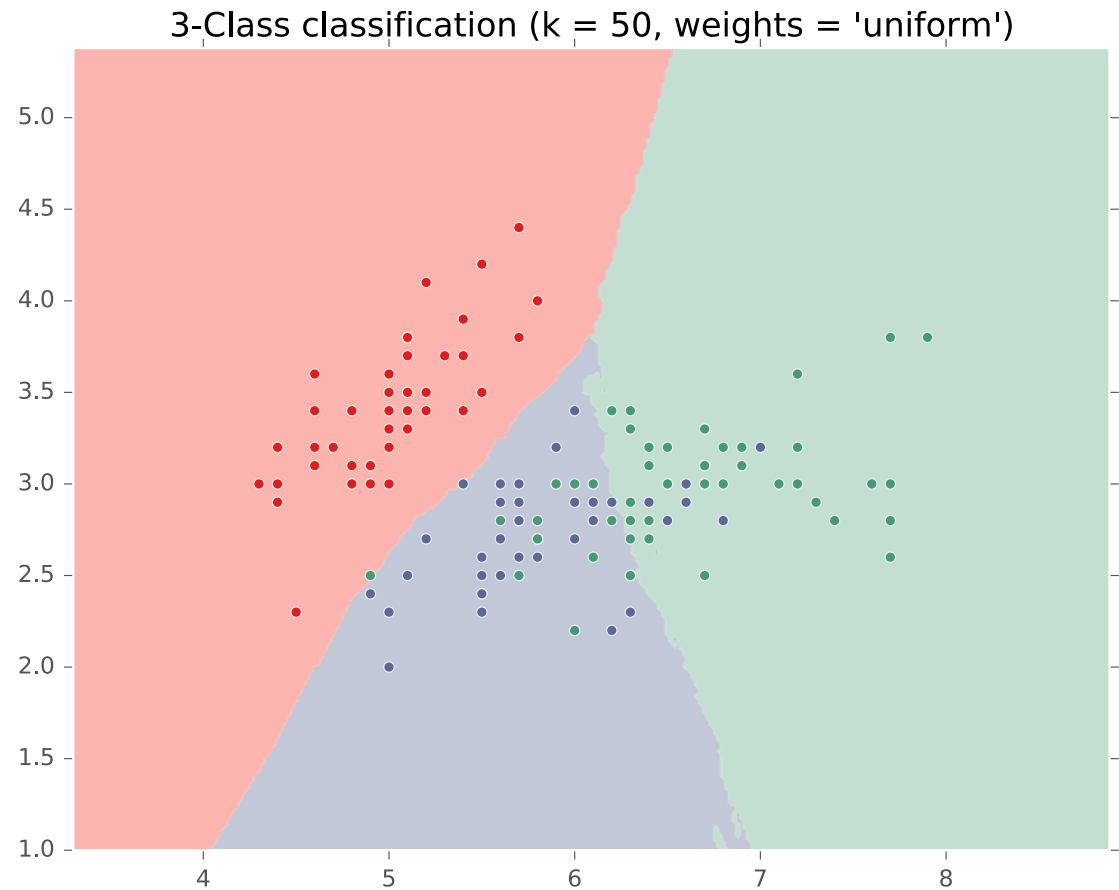
Recall: k -Nearest Neighbors (k NN)

- Classify a point as the most common label among the labels of the k nearest training points
- Tie-breaking (in case of even k and/or more than 2 classes)
 - Weight votes by distance
 - Remove furthest neighbor
 - Add next closest neighbor
 - Use a different distance metric

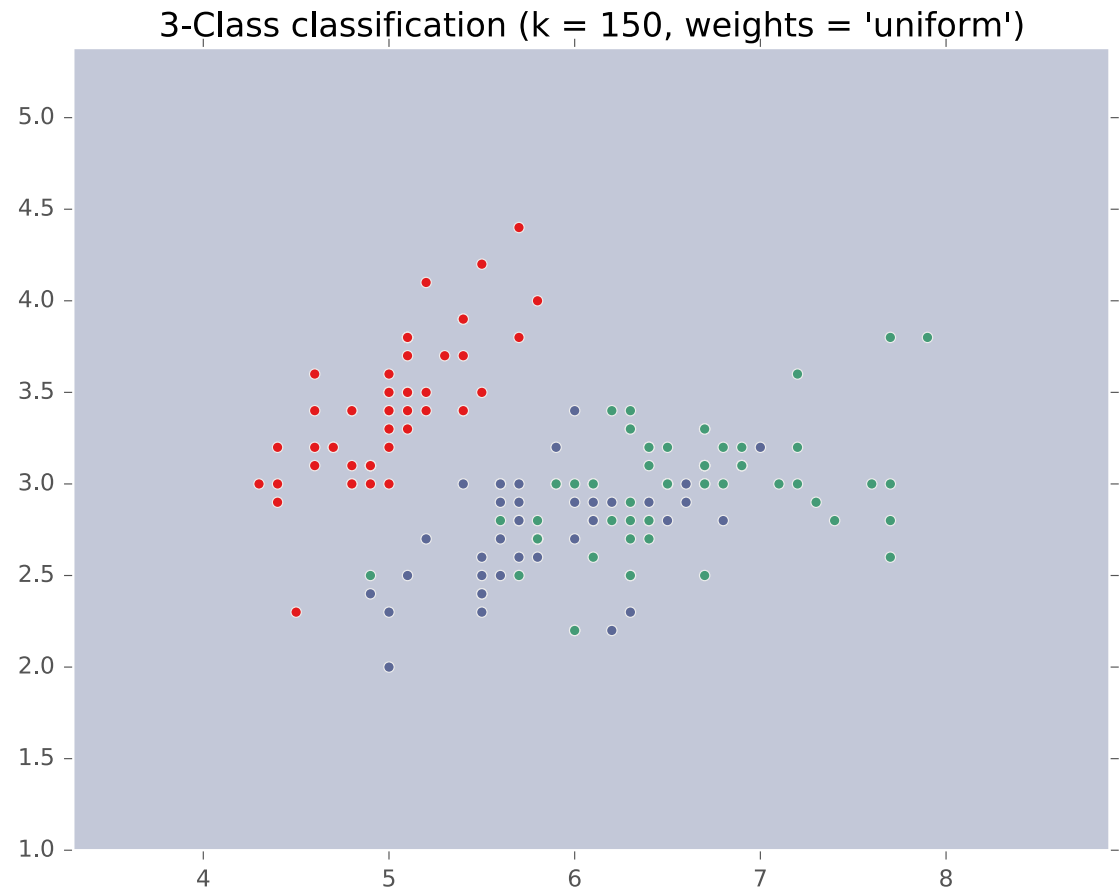
k NN on Fisher Iris Data



k NN on Fisher Iris Data



k NN on Fisher Iris Data



Setting k

- When $k = 1$:
 - many, complicated decision boundaries
 - may *overfit*
- When $k = N$:
 - no decision boundaries; always predicts the most common label in the training data
 - may *underfit*
- k controls the complexity of the hypothesis set $\implies k$ affects how well the learned hypothesis will generalize

Setting k

- Theorem:
 - If k is some function of N s.t. $k(N) \rightarrow \infty$ and $\frac{k(N)}{N} \rightarrow 0$ as $N \rightarrow \infty$...
 - ... then (under certain assumptions) the true error of a k NN model \rightarrow the Bayes error rate
- Heuristics:
 - $k = \lfloor \sqrt{N} \rfloor$
 - $k = 3$
- This is fundamentally a question of **model selection**: each value of k corresponds to a different “model”

Model Selection

- A **model** is a (typically infinite) set of classifiers that a learning algorithm searches through to find the best one (the "hypothesis space")
- **Model parameters** are the numeric values or structure that are selected by the learning algorithm
- **Hyperparameters** are the tunable aspects of the model that are not selected by the learning algorithm

Example: Decision Trees

- Model = set of all possible trees, potentially narrowed down according to the hyperparameters (see below)
- Model parameters = structure of a specific tree e.g., splits, split order, predictions at leaf nodes,
- Hyperparameters = splitting criterion, max-depth, tie-breaking procedures, etc...

Model Selection

- A **model** is a (typically infinite) set of classifiers that a learning algorithm searches through to find the best one (the "hypothesis space")
- **Model parameters** are the numeric values or structure that are selected by the learning algorithm
- **Hyperparameters** are the tunable aspects of the model that are not selected by the learning algorithm

Example: k NN

- Model = set of all possible nearest neighbors classifiers
- Model parameters = none! k NN is a "non-parametric model"
- Hyperparameters = k ,
distance metric,
tie-breaking, etc...

Model Selection with Test Sets

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{test}$, suppose we have multiple candidate models:

$$\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_M$$

- Learn a classifier from each model using only \mathcal{D}_{train} :

$$h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2, \dots, h_M \in \mathcal{H}_M$$

- Evaluate each one using \mathcal{D}_{test} and choose the one with lowest test error:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{test})$$

Model Selection with Test Sets?

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{test}$, suppose we have multiple candidate models:

$$\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_M$$

- Learn a classifier from each model using only \mathcal{D}_{train} :

$$h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2, \dots, h_M \in \mathcal{H}_M$$

- Evaluate each one using \mathcal{D}_{test} and choose the one with lowest test error:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{test})$$

- Is $\operatorname{err}(h_{\hat{m}}, \mathcal{D}_{test})$ a good estimate of $\operatorname{err}(h_{\hat{m}})$?

Model Selection with Validation Sets

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$, suppose we have multiple candidate models:

$$\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_M$$

- Learn a classifier from each model using only \mathcal{D}_{train} :

$$h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2, \dots, h_M \in \mathcal{H}_M$$

- Evaluate each one using \mathcal{D}_{val} and choose the one with lowest *validation error*:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{val})$$

Hyperparameter Optimization with Validation Sets

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$, suppose we have multiple candidate hyperparameter settings:

$$\theta_1, \theta_2, \dots, \theta_M$$

- Learn a classifier for each setting using only \mathcal{D}_{train} :

$$h_1, h_2, \dots, h_M$$

- Evaluate each one using \mathcal{D}_{val} and choose the one with lowest *validation* error:

$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{val})$$

Pro tip: train your final model using *both* training and validation datasets

- Given $\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$, suppose we have multiple candidate hyperparameter settings:

$$\theta_1, \theta_2, \dots, \theta_M$$

- Learn a classifier for each setting using only \mathcal{D}_{train} :

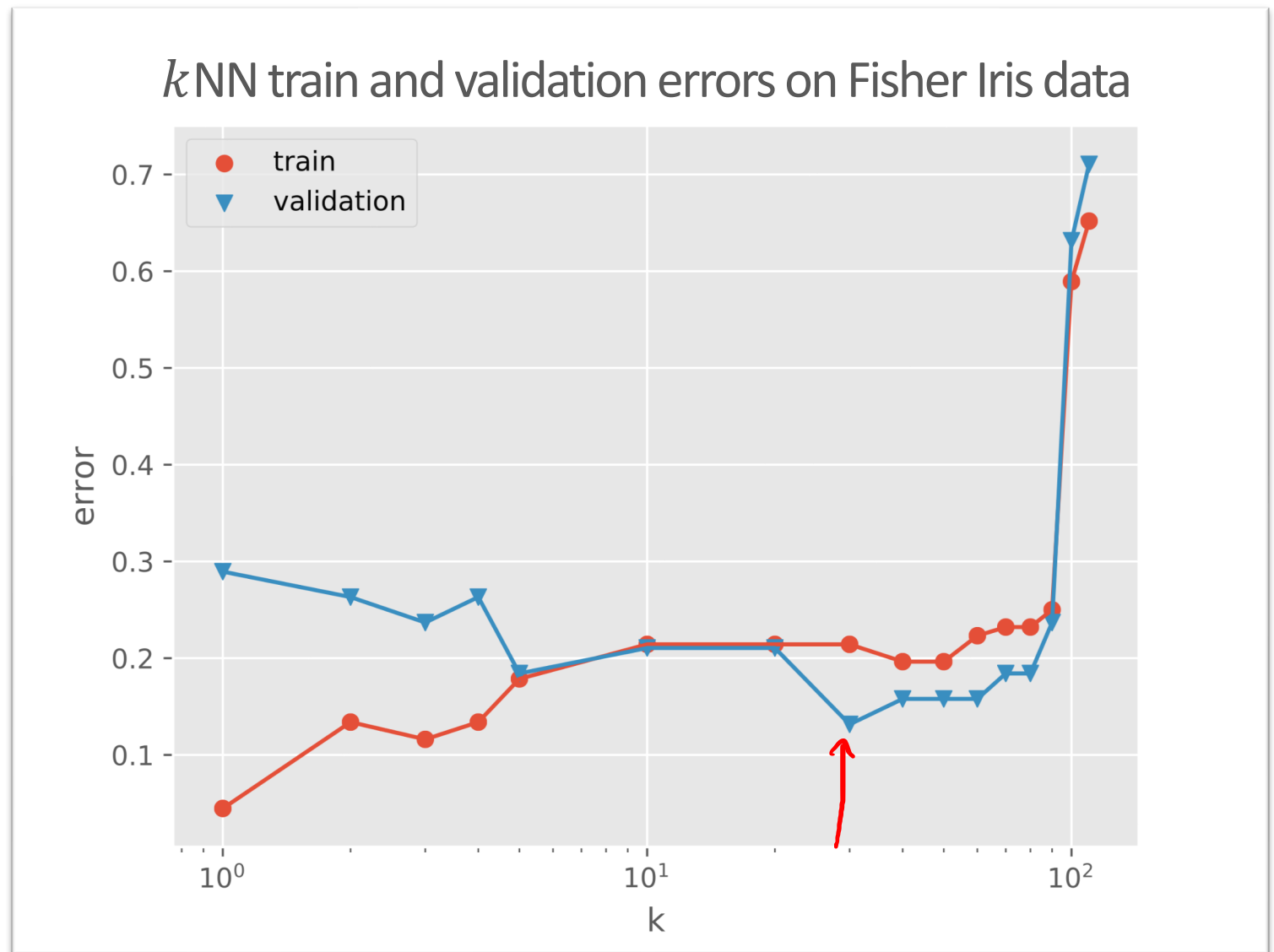
$$h_1, h_2, \dots, h_M$$

- Evaluate each one using \mathcal{D}_{val} and choose the one with lowest *validation* error:

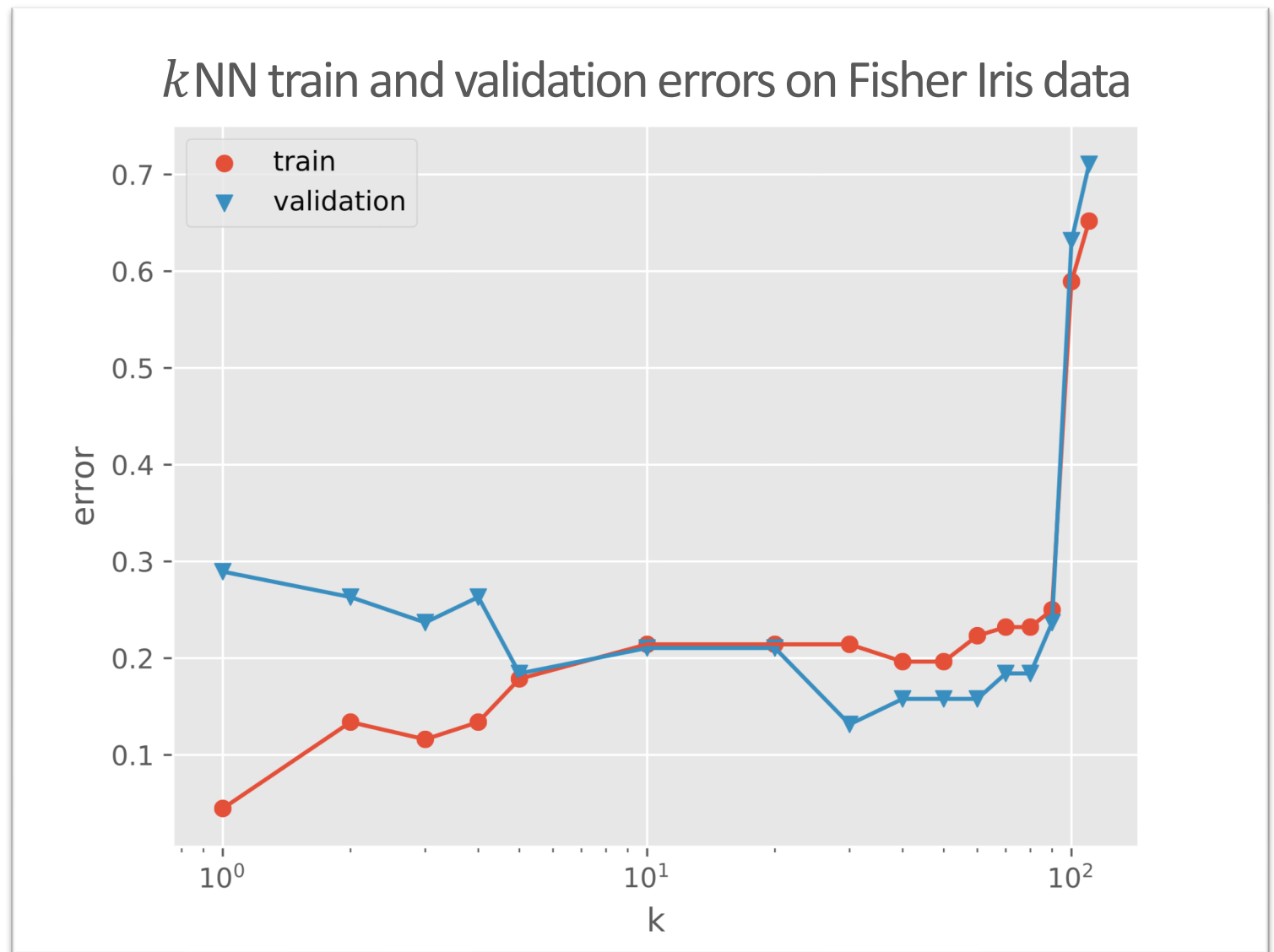
$$\hat{m} = \operatorname{argmin}_{m \in \{1, \dots, M\}} \operatorname{err}(h_m, \mathcal{D}_{val})$$

- Train a new model on $\mathcal{D}_{train} \cup \mathcal{D}_{val}$ using $\theta_{\hat{m}}, h_{\hat{m}}^+$
- $\operatorname{err}(h_{\hat{m}}^+, \mathcal{D}_{test})$ is still a good estimate of $\operatorname{err}(h_{\hat{m}}^+)$!

Setting k for k NN with Validation Sets



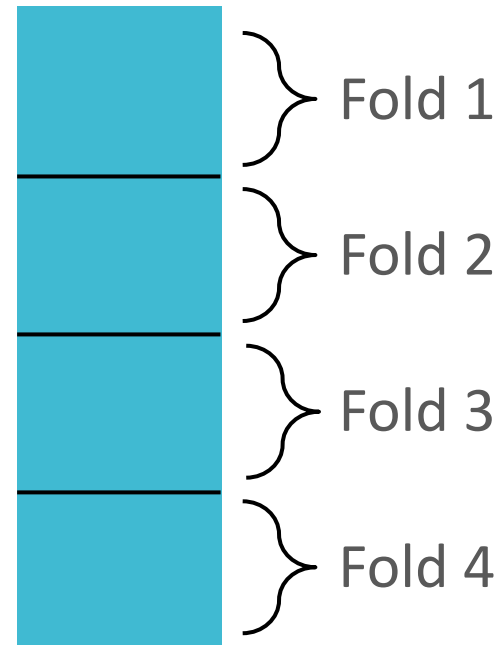
How should we partition our dataset?



K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

- Use each one as a validation set once:



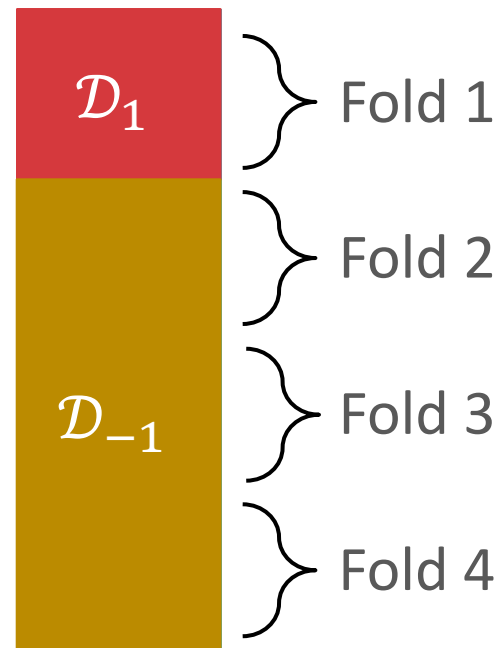
- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

- Use each one as a validation set once:



- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

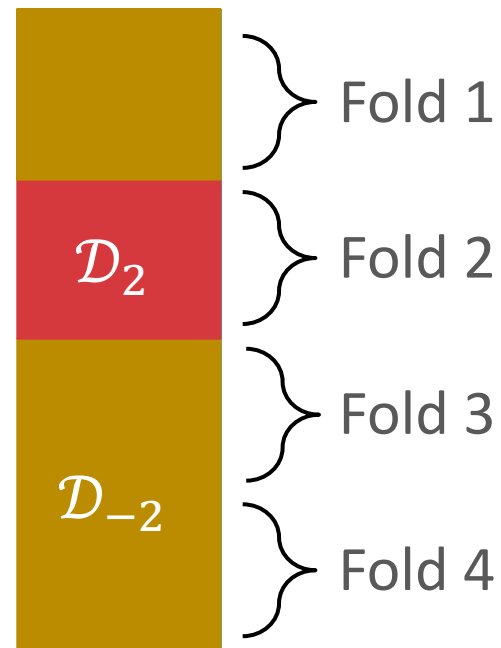
$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds:

$$\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$$

- Use each one as a validation set once:



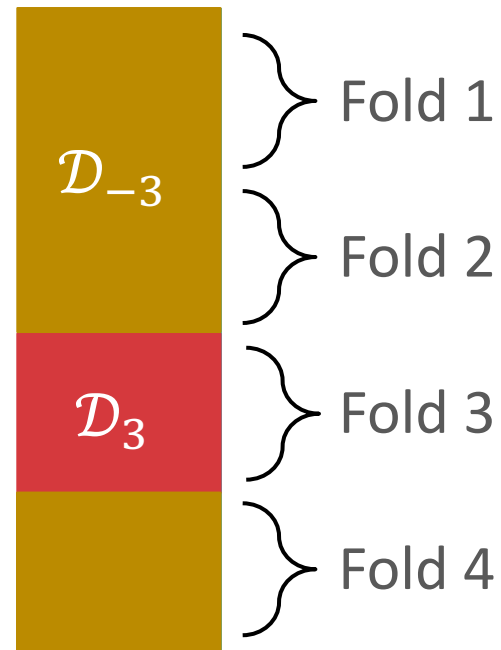
- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

- Use each one as a validation set once:



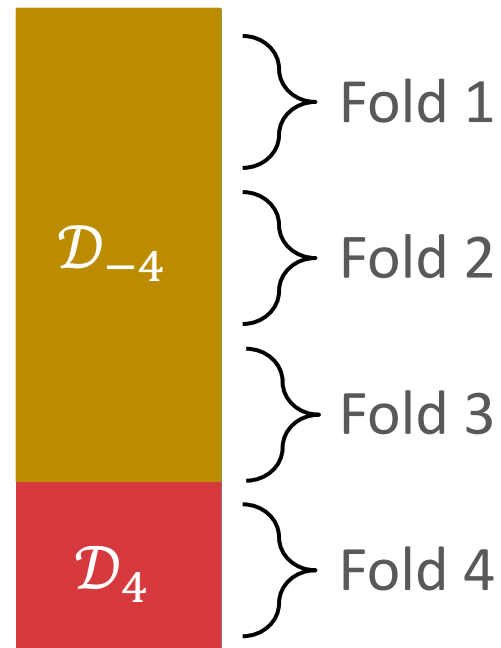
- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

- Use each one as a validation set once:



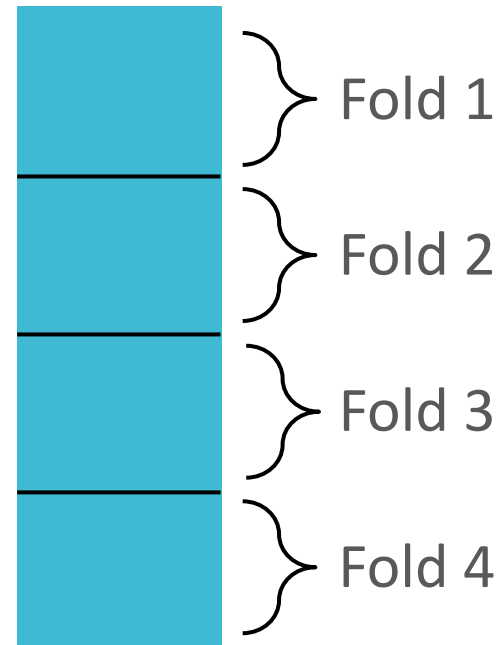
- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

K -fold cross-validation

- Given \mathcal{D} , split \mathcal{D} into K equally sized datasets or folds: $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$

- Use each one as a validation set once:



- Let h_{-i} be the classifier learned using $\mathcal{D}_{-i} = \mathcal{D} \setminus \mathcal{D}_i$ (all folds other than \mathcal{D}_i) and let $e_i = \text{err}(h_{-i}, \mathcal{D}_i)$
- The K -fold cross validation error is

$$\text{err}_{cv_K} = \frac{1}{K} \sum_{i=1}^K e_i$$

- Special case when $K = N$: Leave-one-out cross-validation
- Choosing between m candidates requires training mK times

Summary

	Input	Output
Training	<ul style="list-style-type: none">• training dataset• hyperparameters	<ul style="list-style-type: none">• best model parameters
Hyperparameter Optimization	<ul style="list-style-type: none">• training dataset• validation dataset	<ul style="list-style-type: none">• best hyperparameters
Cross-Validation	<ul style="list-style-type: none">• training dataset• validation dataset	<ul style="list-style-type: none">• cross-validation error
Testing	<ul style="list-style-type: none">• test dataset• classifier	<ul style="list-style-type: none">• test error

Key Takeaways

- Real-valued features and decision boundaries
- Nearest neighbor model and generalization guarantees
- k NN “training” and prediction
- Effect of k on model complexity
- k NN inductive bias
- Differences between training, validation and test datasets in the model selection process
- Cross-validation for model selection
- Relationship between training, hyperparameter optimization and model selection

Recall: Regression

- Learning to diagnose heart disease

as a **(supervised)**

regression task

features

targets

data points

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000

Decision Tree Regression

- Learning to diagnose heart disease

as a **(supervised)**

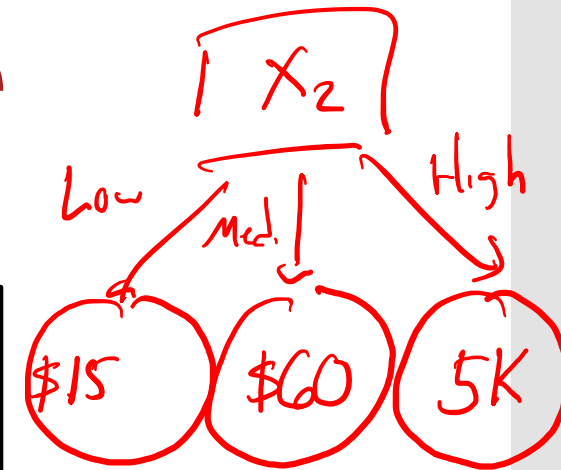
regression task

features

targets

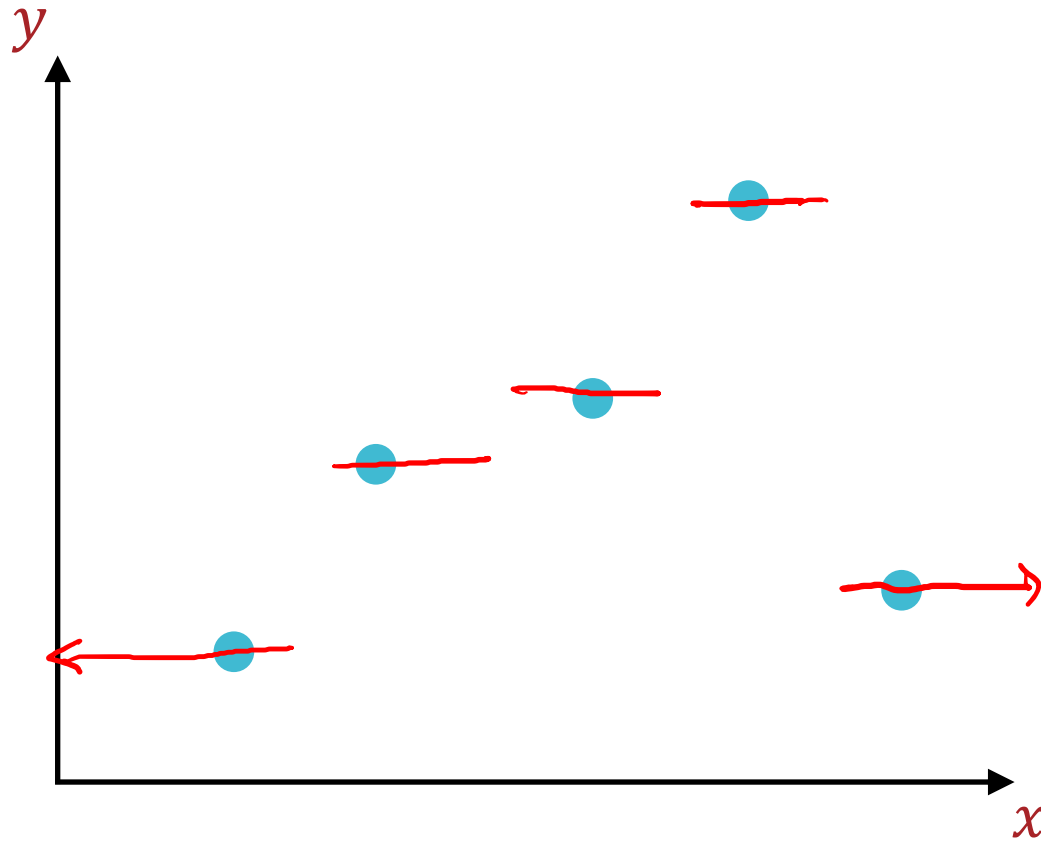
data points

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000



1-NN Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and one-dimensional inputs $x \in \mathbb{R}$



Linear Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and D -dimensional inputs $\mathbf{x} = [x_1, \dots, x_D]^T \in \mathbb{R}^D$
- **Assume**

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

Linear Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and D -dimensional inputs $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- Assume

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

$\curvearrowright [w_0, w_1, w_2, \dots, w_D]$

- Notation: given training data $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$

$$\bullet X = \begin{bmatrix} 1 & \mathbf{x}^{(1)T} \\ 1 & \mathbf{x}^{(2)T} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(N)T} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_D^{(1)} \\ 1 & x_1^{(2)} & \dots & x_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)}$$

is the *design matrix*

- $\mathbf{y} = [y^{(1)}, \dots, y^{(N)}]^T \in \mathbb{R}^N$ is the *target vector*

General Recipe for Machine Learning

1. Define a model and model parameters
2. Write down an objective function
3. Optimize the objective w.r.t. the model parameters

Recipe for Linear Regression

1. Define a model and model parameters

- Assume linear dependence

- Parameters: $w = [w_0, w_1, \dots, w_D]^T$

2. Write down an objective function

minimize the mean squared error

$$l_D(w) = \frac{1}{N} \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2$$

3. Optimize the objective w.r.t. the model parameters

Solve in closed-form: take gradients, set equal to 0 and solve

$$X = \begin{bmatrix} 1 & x^{(1)T} \\ 1 & x^{(2)T} \\ \vdots & \vdots \end{bmatrix}$$

Minimizing the Squared Error

$$H_w l_D(w)$$

$$= \frac{1}{N} (2X^T X)$$

which is positive semi-definite

\approx the second-order derivative ≥ 0

$$l_D(w) = \frac{1}{N} \sum_{n=1}^N \left(\underline{x^{(n)T} w} - y^{(n)} \right)^2$$

$$= \frac{1}{N} \left(\underline{Xw - y} \right)^T \left(\underline{Xw - y} \right) = \frac{1}{N} \|Xw - y\|_2^2$$

$$= \frac{1}{N} \left(\underline{w^T X^T X w} - \underline{2w^T X^T y} + \underline{y^T y} \right)$$

$$\nabla_w l_D(w) = \frac{1}{N} \left(\underline{2X^T X w} - \underline{2X^T y} + \underline{0} \right)$$

$$\Rightarrow \nabla_w l_D(\hat{w}) = 0 = \frac{1}{N} \left(2X^T X \hat{w} - 2X^T y \right)$$

$$\Rightarrow 2X^T X \hat{w} = 2X^T y \Rightarrow \hat{w} = \left(X^T X \right)^{-1} X^T y$$

Closed Form Solution

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

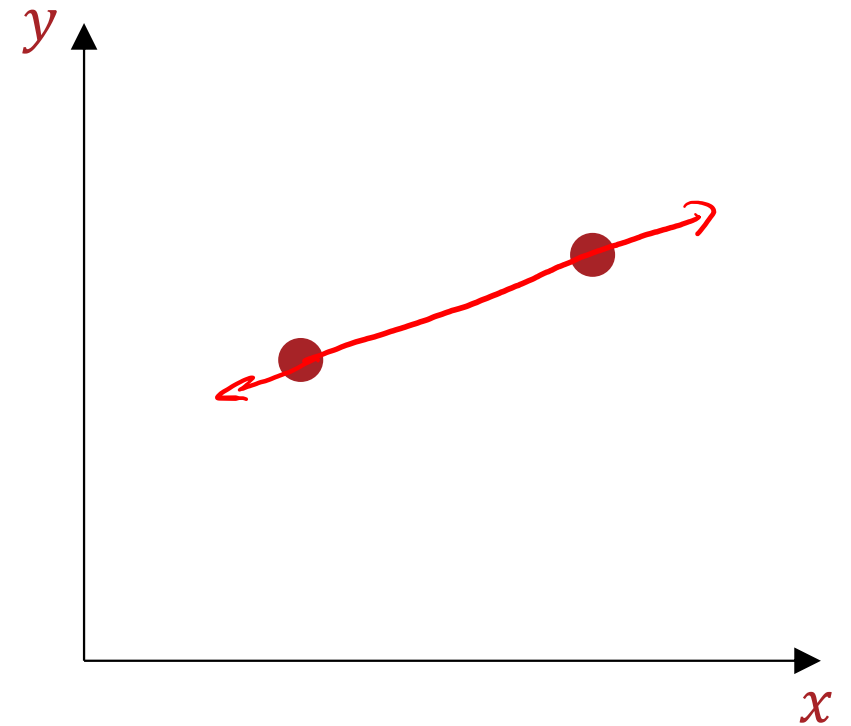
1. Is $\mathbf{X}^T \mathbf{X}$ invertible?

Almost always yes!

2. If so, how computationally expensive is inverting $\mathbf{X}^T \mathbf{X}$?

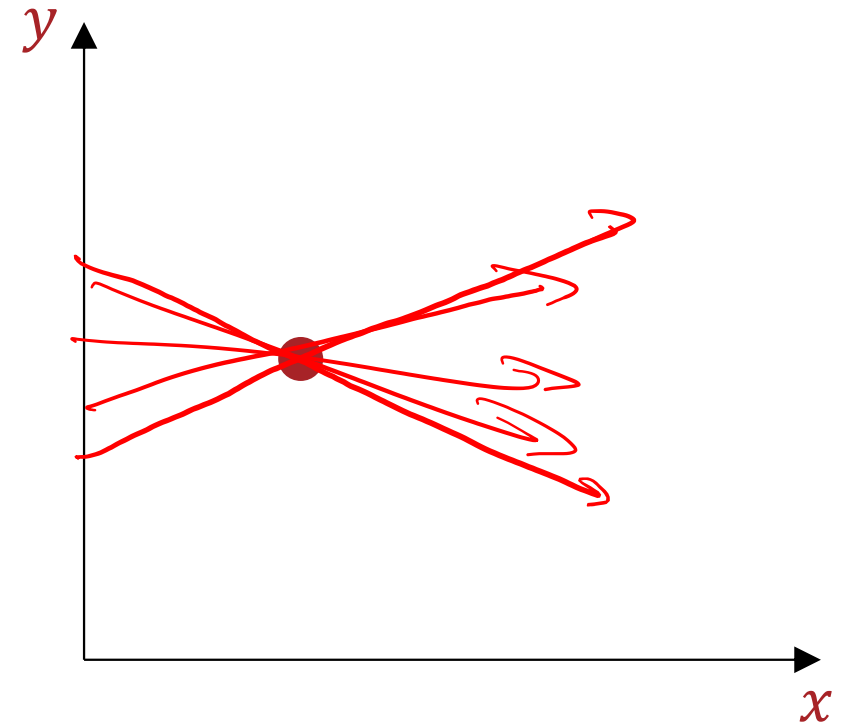
Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights \mathbf{w}) are there for the given dataset?



Linear Regression: Uniqueness

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights \mathbf{w}) are there for the given dataset?

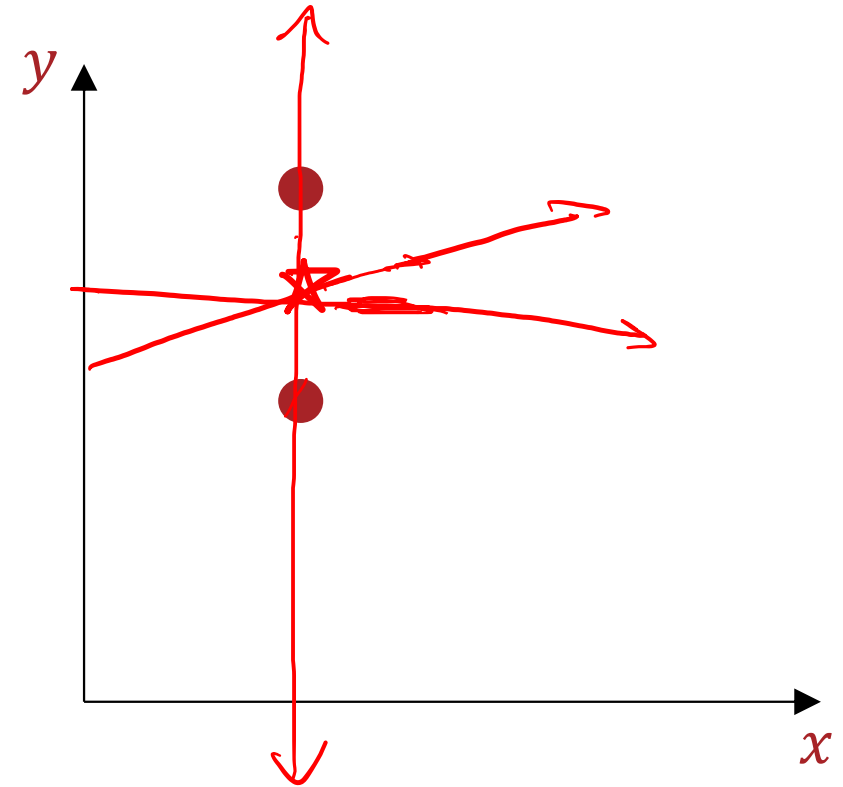


Linear Regression: Uniqueness

$$l_D(w) = \frac{1}{N} \sum_{n=1}^N |w^T x^{(n)} - y^{(n)}|$$

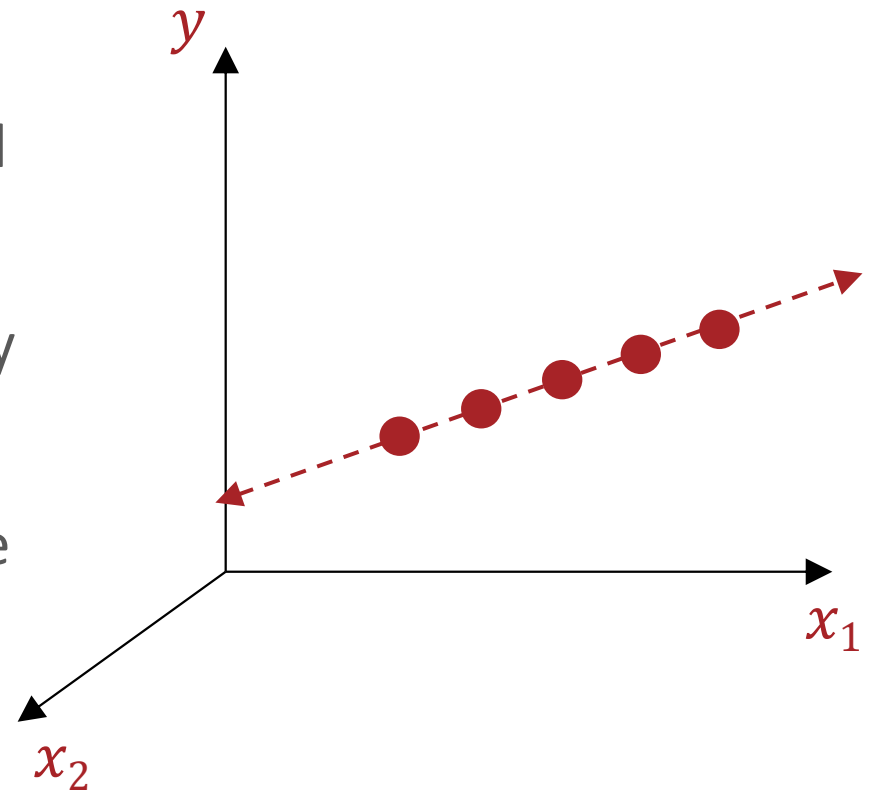
$x = 5$

- Consider a 1D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights w) are there for the given dataset?



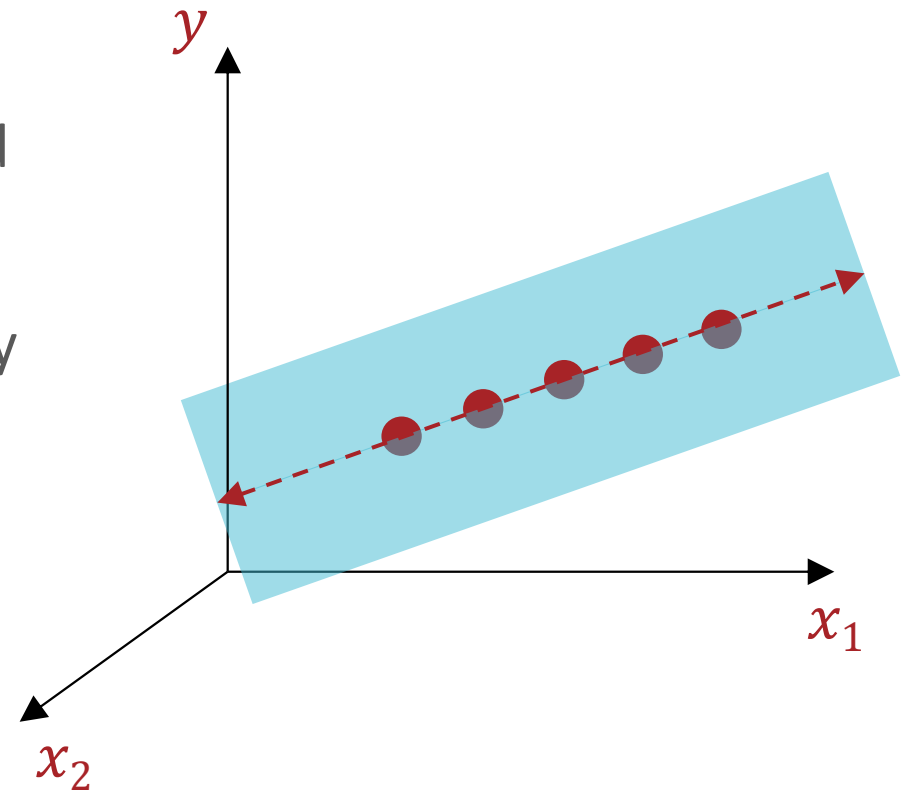
Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of parameters θ) are there for the given dataset?



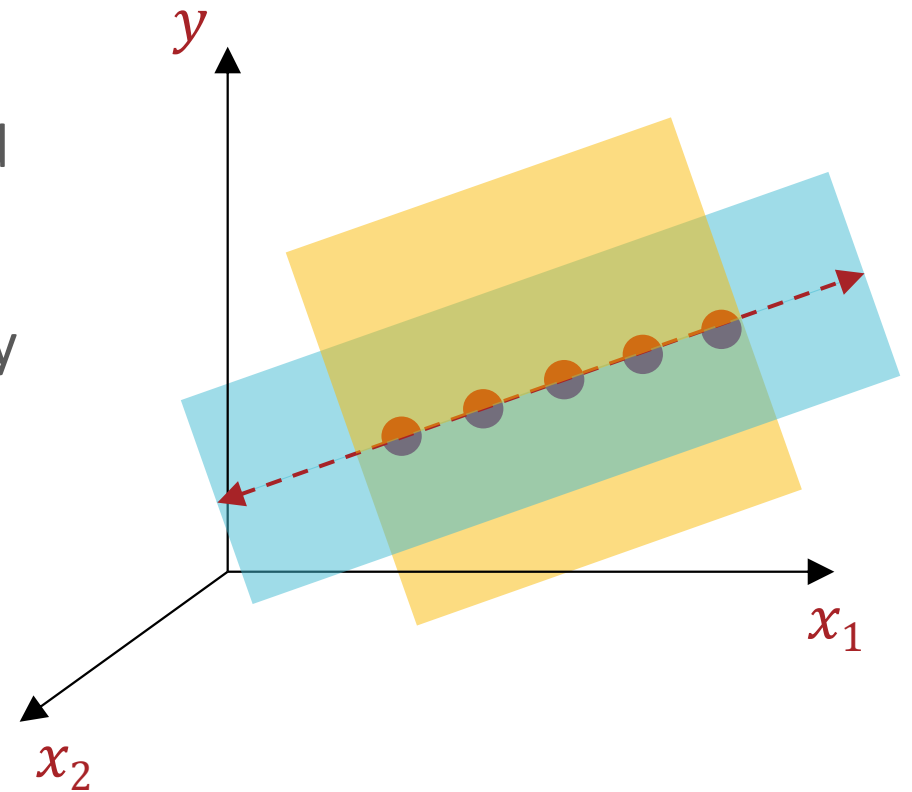
Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights \mathbf{w}) are there for the given dataset?



Linear Regression: Uniqueness

- Consider a 2D linear regression model trained to minimize the mean squared error: how many optimal solutions (i.e., sets of weights \mathbf{w}) are there for the given dataset?



Closed Form Solution

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1. Is $\mathbf{X}^T \mathbf{X}$ invertible?

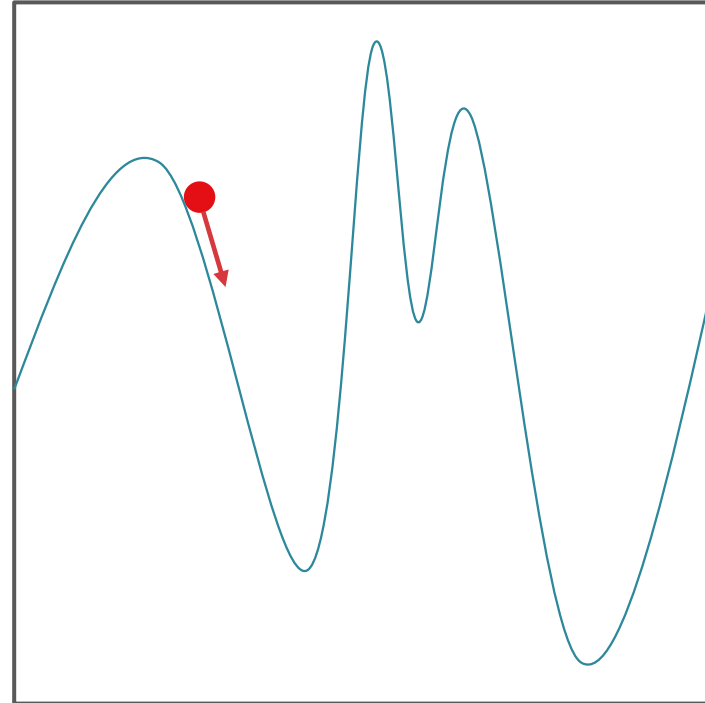
- When $N \gg D + 1$, $\mathbf{X}^T \mathbf{X}$ is (almost always) full rank and therefore, invertible
- If $\mathbf{X}^T \mathbf{X}$ is not invertible (occurs when one of the features is a linear combination of the others) then there are infinitely many solutions.

2. If so, how computationally expensive is inverting $\mathbf{X}^T \mathbf{X}$?

$O(D^3)$ because $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(D+1) \times (D+1)}$
($\exists O(D^{2.4 \dots})$)

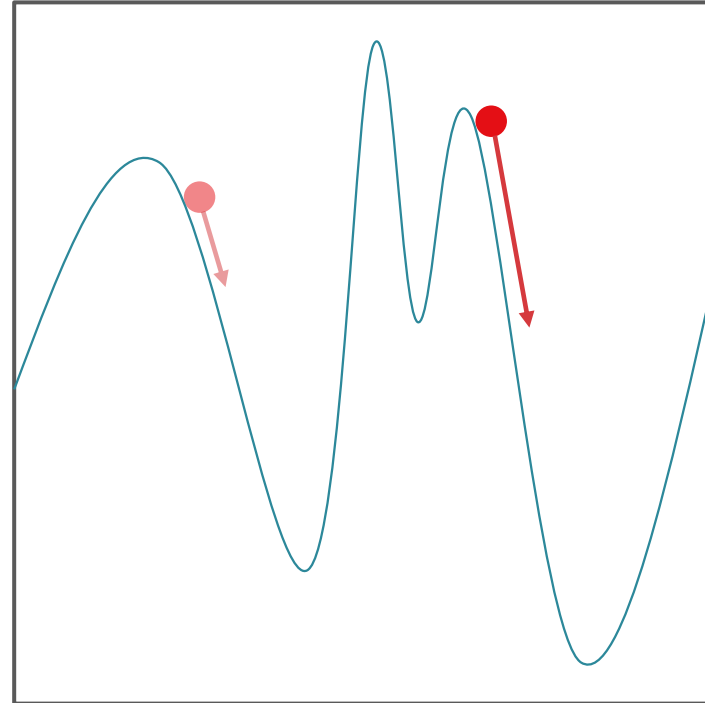
Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



Gradient Descent: Intuition

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere

