# 50.007 Machine Learning

2026 Spring

# 4. Regression

Na Zhao

Assistant Professor, ISTD

# Recap - Linear Classifier (separable case)

1. Training Set (Linearly Separable)

$$\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \ldots, \left(x^{(n)}, y^{(n)}\right)$$

2. Model (Linear Classifier)

$$h(x; \theta) = \text{sign}(\theta_1 x_1 + \cdots + \theta_d x_d)$$

3. Training Error (Zero-one Loss)

$$\varepsilon_n(\theta) = \frac{1}{n}\sum_{(x,y)\in\mathcal{S}_n} [\![\, y(\theta^\top x) \leq 0 \,]\!]$$

4. Algorithm (The Perceptron Algorithm)

$$\text{if } y^{(t)} \neq h(x^{(t)}; \theta^{(k)}) \text{ then}$$
$$\theta^{(k+1)} = \theta^{(k)} + y^{(t)} x^{(t)}$$

# Recap - Linear Classifier (Non-separable case)

1. **Training Set**  (<span style="color:darkred">Not Necessarily Linearly Separable</span>)

$$\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \ldots, \left(x^{(n)}, y^{(n)}\right)$$

2. **Model** (<span style="color:blue">Linear Classifier</span>)

$$h(x; \theta) = \mathrm{sign}(\theta_1 x_1 + \cdots + \theta_d x_d)$$

3. **Training Error**  (<span style="color:darkred">Hinge Loss</span>)

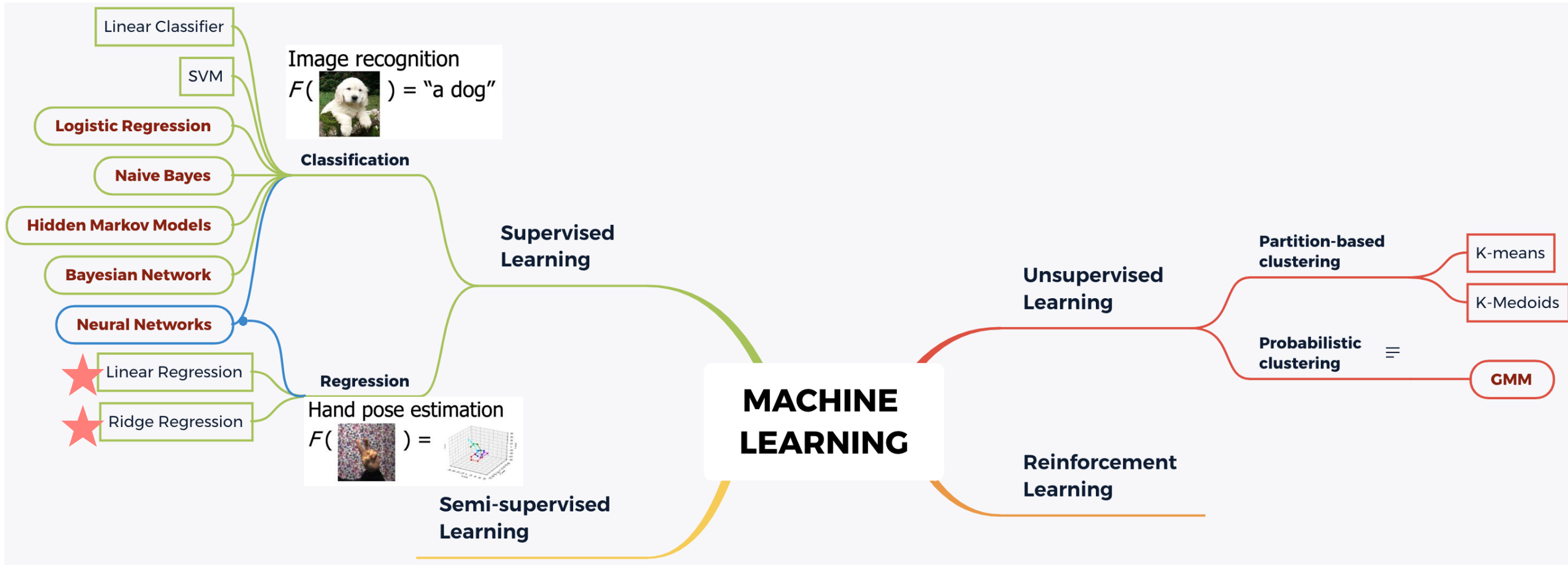$$R_n(\theta) = \frac{1}{n} \sum_{(x,y) \in \mathcal{S}_n} \max\{1 - y(\theta^\top x), 0\}$$

4. **Algorithm** (<span style="color:darkred">Stochastic Sub-Gradient Descent</span>)

$$\text{select } t \in \{1, \ldots, n\} \text{ at random,}$$

$$\text{if } y^{(t)}(\theta^{(k)} \cdot x^{(t)}) \leq 1, \text{ then}$$

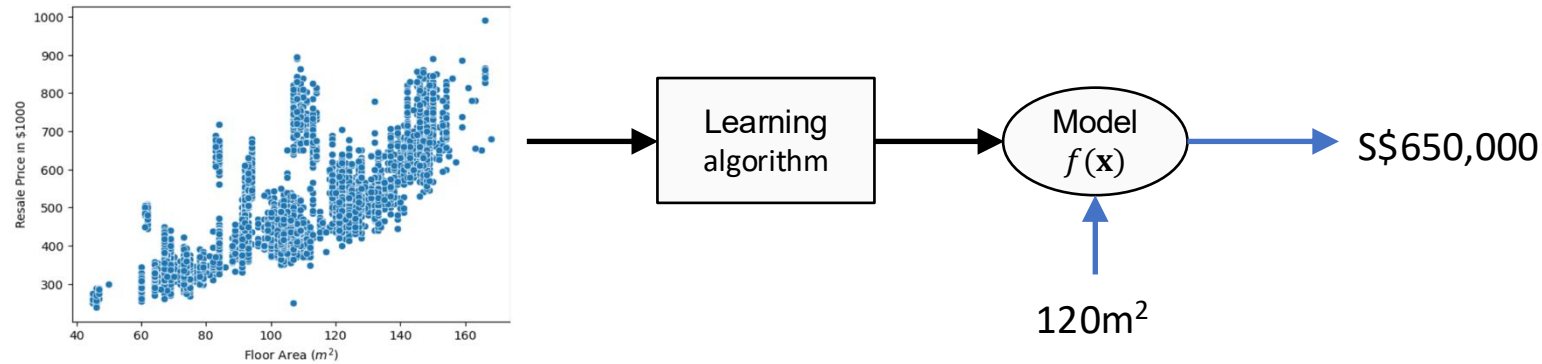$$\theta^{(k+1)} = \theta^{(k)} + \eta_k y^{(t)} x^{(t)}$$

# Roadmap

# Learning Objectives
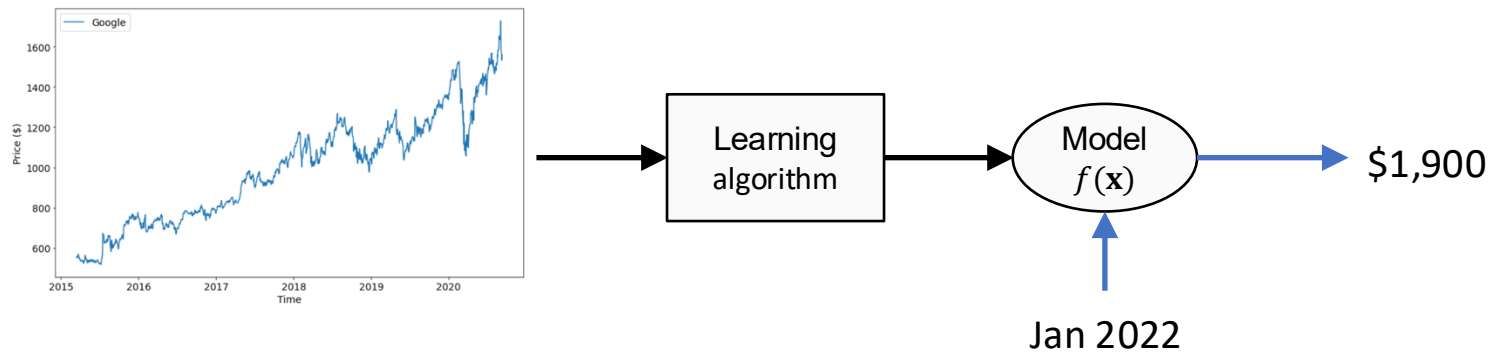
You should be able to know:

1. What is linear regression?

2. How to learn a linear regression model?

3. What is ridge regression and how is ridge regression different from linear regression?

4. What is regularization and why do we need to do regularization?

# Regression Examples

- **House Price Prediction**



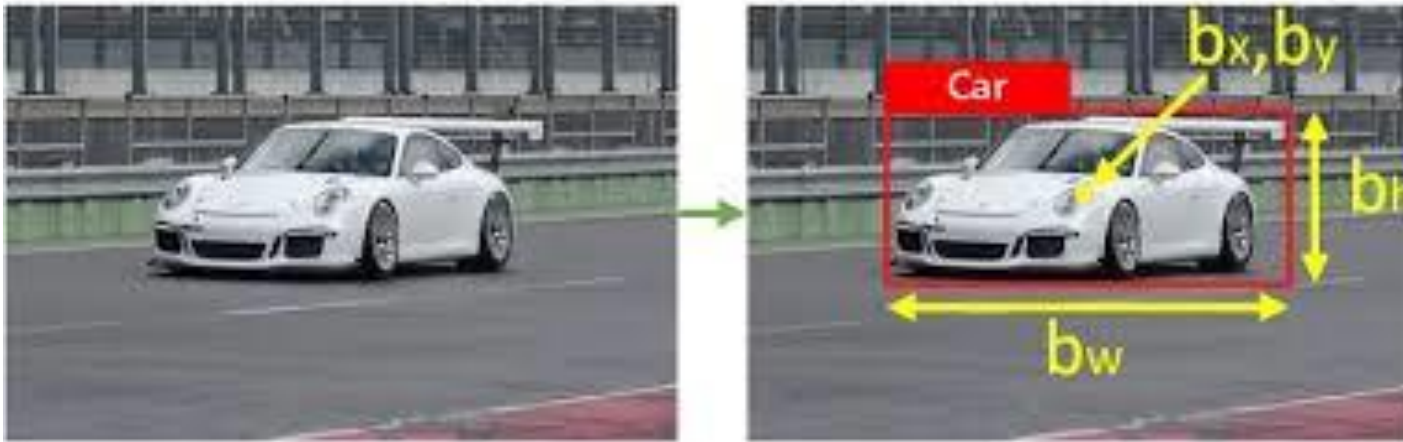- **Stock Price Prediction**



Learning a function
$$y = f(x)$$
$$x \in \mathbb{R}^d$$
$$y \in \mathbb{R}^m$$

$m$=1 in these two examples

# Regression Examples

- **Object Detection (Localization)**
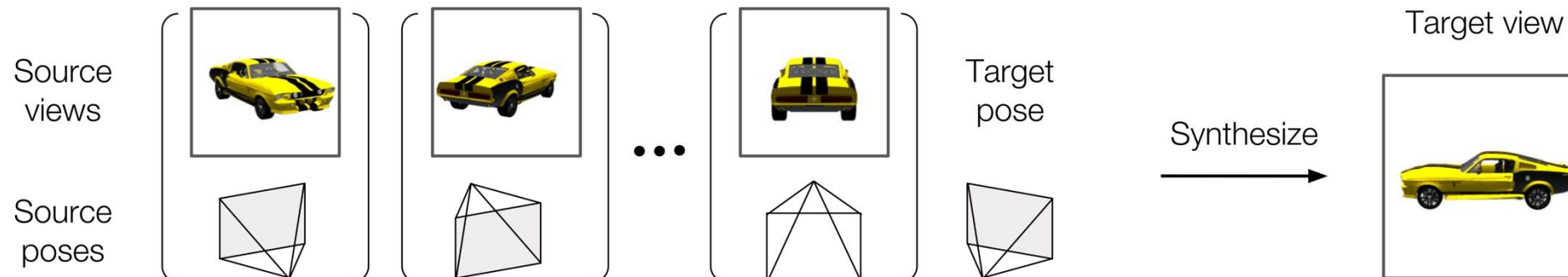


Learning a function
$$y = f(x)$$
$$x \in \mathbb{R}^d$$
$$y \in \mathbb{R}^m$$

$m>1$ in these two examples

- **Novel View Synthesis**

# Linear Regression

- A **linear regression** function is simply a linear function of the feature vectors:

$$f(x; \theta, \theta_0) = \theta \cdot x + \theta_0 = \sum_{i=1}^{d} \theta_i x_i + \theta_0$$

*↖ training sample (fixed)*

- Similar as in linear classification, <span style="color:red">different parameter choices</span> $\theta \in \mathbb{R}^d$, $\theta_0 \in \mathbb{R}$, give rise to the **hypothesis set** $\mathcal{F}$. *determined by $x$ & $\theta$*

- **Training Objective**: find $f(x; \hat{\theta}, \hat{\theta}_0)$ that would yield accurate predictions on yet unseen examples based on $\mathcal{S}_n$.

**Training data**

$$\mathcal{S}_n = \left\{ \left( x^{(t)}, y^{(t)} \right) \mid t = 1, \ldots, n \right\}$$

- Features/Inputs $x^{(t)} = \left( x_1^{(t)}, \ldots, x_d^{(t)} \right)^\top \in \mathbb{R}^d$
- Response/Output $y^{(t)} \in \mathbb{R}$

How does it differ from the linear classification function?

# Linear Regression – Cost Function

*shows difference btw true label & predicted label*

- Similar to the classification setting, we will measure training error in terms of **empirical risk**:

*n training sample*

*true label*  *predicted label*

$$R_n(\theta) = \frac{1}{n}\sum_{t=1}^{n}\text{Loss}(y^{(t)} - \theta \cdot x^{(t)})$$

$\theta' \cdot x^{(t)}$,

$$\theta' = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, \quad x' = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Note $\theta_0$ is dropped here for simplicity, remember you can always add it back.

- Cost/loss function: "**Least Squares**" criterion

why putting square here?

$$R_n(\theta) = \frac{1}{n}\sum_{t=1}^{n}\text{Loss}(y^{(t)} - \theta \cdot x^{(t)}) = \frac{1}{n}\sum_{t=1}^{n}(y^{(t)} - \theta \cdot x^{(t)})^2/2$$
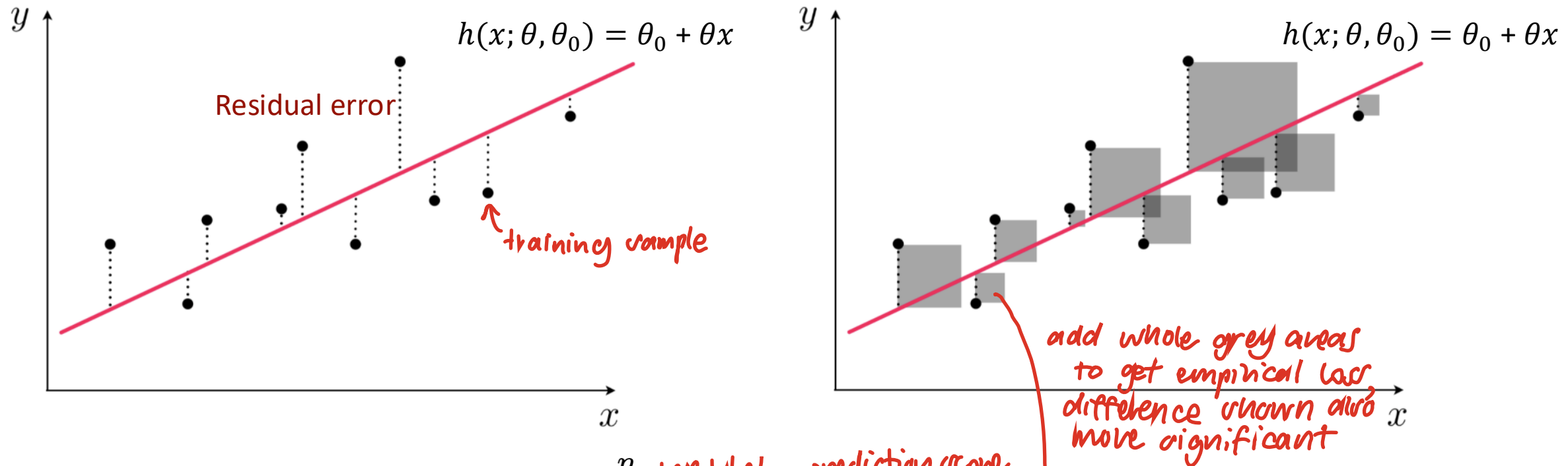
why is the result divided by 2?

- ▪ The square is to avoid cancellation due to positive and negative differences
- ▪ And it also can permit small discrepancies but heavily penalize large deviations

*# difference × be >1*

# Cost Function - Graphical Interpretation



$h(x; \theta, \theta_0) = \theta_0 + \theta x$

Residual error

training sample

$h(x; \theta, \theta_0) = \theta_0 + \theta x$

add whole grey areas to get empirical loss, difference shown also more significant

true label

prediction score

$$R_n(\theta) = \frac{1}{n} \sum_{t=1}^{n} (y^{(t)} - h(x^{(t)}; \theta, \theta_0))^2 / 2$$

# Linear Regression – Optimization

*↗iterative way to optimize linear regression*

- **Stochastic Gradient Descent**

$\theta^{(0)} = 0$ (vector)    *# Random initialization*
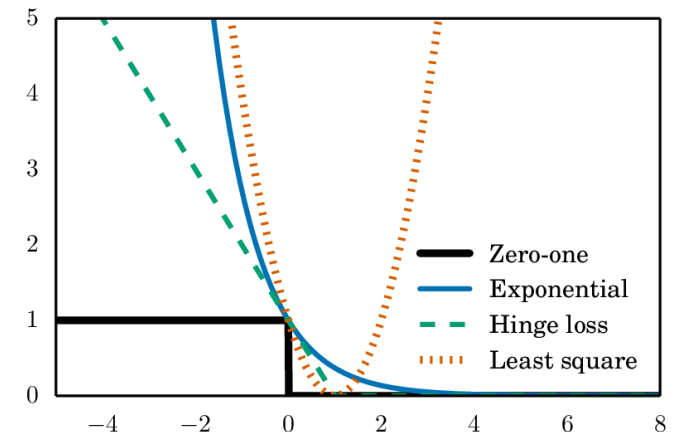
select $t \in \{1, \dots, n\}$ at random,

*←difference*

$\theta^{(k+1)} = \theta^{(k)} - \eta_k \nabla_\theta (y^{(t)} - \theta \cdot x^{(t)})^2 / 2 \big|_{\theta=\theta^{(k)}}$    *# Update*

Repeat the update step until stopping criterion is met.

$\nabla_\theta (y^{(t)} - \theta \cdot x^{(t)})^2 / 2 = (y^{(t)} - \theta \cdot x^{(t)}) \nabla_\theta (y^{(t)} - \theta \cdot x^{(t)})$

$= -(y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}$

$\theta^{(k+1)} = \theta^{(k)} + \eta_k (y^{(t)} - \theta \cdot x^{(t)}) x^{(t)}$



Legend:
- Zero-one
- Exponential
- Hinge loss
- Least square

$R_n(\theta) = \dfrac{1}{n} \displaystyle\sum_{t=1}^{n} (y^{(t)} - \theta \cdot x^{(t)})^2 / 2$

- *Convex!*
- *Everywhere differentiable!*

# Linear Regression – Optimization

- **Closed Form Solution**
  - **Find** $\hat{\theta}$: minimizing empirical risk directly by **setting gradient to zero**.

$$\nabla R_n(\theta)_{\theta=\hat{\theta}} = \frac{1}{n} \sum_{t=1}^{n} \nabla_\theta \left\{ (y^{(t)} - \theta \cdot x^{(t)})^2 / 2 \right\}_{|\theta=\hat{\theta}}$$

*gradient*

$$= \frac{1}{n} \sum_{t=1}^{n} \{ -(y^{(t)} - \hat{\theta} \cdot x^{(t)}) x^{(t)} \}$$

$$= -\frac{1}{n} \sum_{t=1}^{n} y^{(t)} x^{(t)} + \frac{1}{n} \sum_{t=1}^{n} (\hat{\theta} \cdot x^{(t)}) x^{(t)} \qquad \boxed{\hat{\theta} \cdot x^{(t)} = (x^{(t)})^T \hat{\theta}}$$

$$= -\underbrace{\frac{1}{n} \sum_{t=1}^{n} y^{(t)} x^{(t)}}_{=b} + \underbrace{\frac{1}{n} \sum_{t=1}^{n} x^{(t)} (x^{(t)})^T}_{=A} \overset{\text{← optimized}}{\hat{\theta}} = -b + A\hat{\theta}$$

# Linear Regression – Optimization

- **Closed Form Solution**
  - **Find** $\hat{\theta}$: minimizing empirical risk directly by **setting gradient to zero**.

$$\nabla R_n(\theta)_{\theta=\hat{\theta}} = -\underbrace{\frac{1}{n}\sum_{t=1}^{n} y^{(t)} x^{(t)}}_{=b} + \underbrace{\frac{1}{n}\sum_{t=1}^{n} x^{(t)}(x^{(t)})^T}_{=A}\hat{\theta} = -b + A\hat{\theta} = 0$$

$$X = \begin{bmatrix} x^1 \\ x_2 \\ \vdots \\ x^n \end{bmatrix}$$

  - If $A$ is **invertible** $(n \geq d)$, we can find the *optimal* parameter $\hat{\theta}$ as:

$$\hat{\theta} = A^{-1}b. \qquad {\in \mathbb{R}^d}$$

where, $b = \dfrac{1}{n}X^T\vec{y}, \quad A = \dfrac{1}{n}X^T X$

$b \in \mathbb{R}^d$, $A \in \mathbb{R}^{d\times d}$

$X = [x^{(1)}, \ldots, x^{(n)}]^T \in \mathbb{R}^{n\times d}$
$\vec{y} = [y^{(1)}, \ldots, y^{(n)}]^T \in \mathbb{R}^{n}$

  - **Note**: The time complexity of inverting $A$ is $\mathcal{O}(d^3)$. It would be costly if $d$ is large.

# What can go wrong in Linear Regression?

- $n < d$ (A is *not* invertible)
  - It is not unusual to see the number of input variables greatly exceed the number of observations, *e.g.* micro-array data analysis, environmental pollution studies.

- Overfitting
  - The training objective is just about fitting the data as well as possible, but we might also want to keep the hypothesis from getting *too* attached to the training data, as this may hinder its generalization to unseen data.

☞ Regularizing parameters*!*

# Ridge Regression

- The cost/loss function of **ridge regression**:

$$J_{n,\lambda}(\theta) = \frac{\lambda}{2}\|\theta\|^2 + R_n(\theta) = \boxed{\frac{\lambda}{2}\|\theta\|^2} + \boxed{\frac{1}{n}\sum_{t=1}^{n}(y^{(t)} - \theta \cdot x^{(t)})^2/2}$$

<span style="color:red">Pressure to simplify model</span>   <span style="color:blue">Pressure to fit data</span>

- $\lambda > 0$ is the regularization parameter
- It quantifies the *trade-off* between <u>keeping the parameters small</u> and <u>fitting to the training data</u>.

- $\|\theta\|^2$ is also known as **$L_2$ regularization**, it can be used in many contexts aside from linear regression, *e.g.,* logistic regression and support vector machines.

# Ridge Regression – Optimization

$$\nabla_\theta \left\{ \frac{\lambda}{2}\|\theta\|^2 + (y^{(t)} - \theta \cdot x^{(t)})^2/2 \right\}_{|\theta = \theta^{(k)}} = \lambda\theta^{(k)} - (y^{(t)} - \theta^{(k)} \cdot x^{(t)})x^{(t)}$$

$$\nabla_\theta (y^{(t)} - \theta \cdot x^{(t)})^2/2 = -(y^{(t)} - \theta \cdot x^{(t)})x^{(t)}$$

- **Stochastic Gradient Descent**

$\theta^{(0)} = 0$ (vector)   *# Random initialization*

select $t \in \{1, \dots, n\}$ at random,

$$\theta^{(k+1)} = \boxed{(1 - \lambda\eta_k)}\theta^{(k)} + \eta_k(y^{(t)} - \theta \cdot x^{(t)})x^{(t)}$$   *# Update*

Repeat the update step until stopping criterion is met.

$$\theta^{k+1} = \theta^k - \eta_k \nabla_\theta L_r$$

$$\eta_k \left( \lambda\theta^k + \nabla_\theta L_i \right)$$
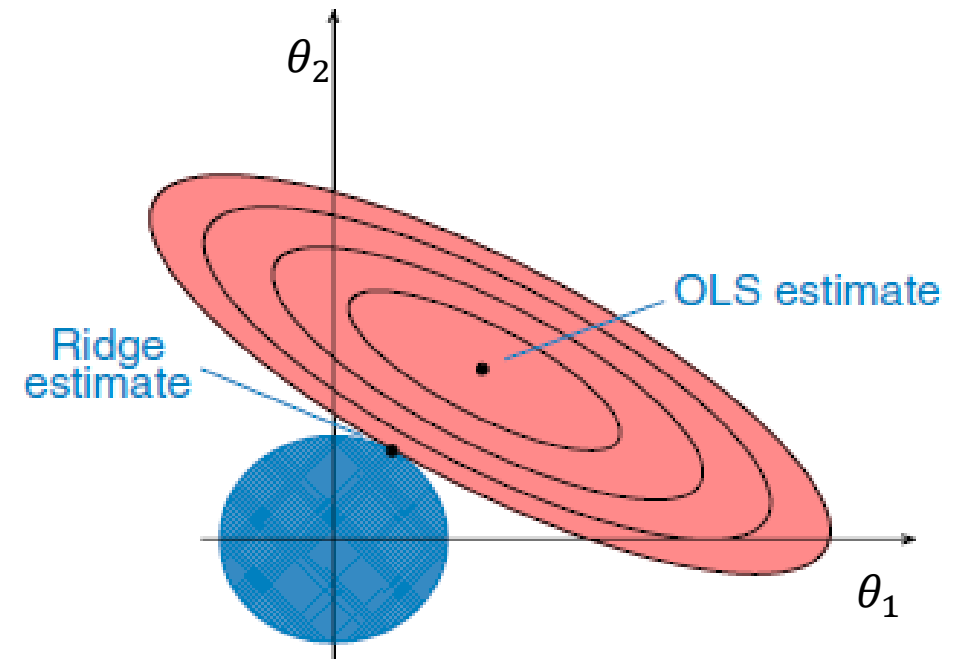
- **Closed Form Solution**

$$\nabla \frac{\lambda}{2}\|\theta\|^2 + R_n(\theta) = -\underbrace{\frac{1}{n}\sum_{t=1}^{n} y^{(t)}x^{(t)}}_{=b} + \underbrace{\frac{1}{n}\sum_{t=1}^{n} (\lambda + x^{(t)}(x^{(t)})^T)\,\hat{\theta}}_{=A} = -b + A\hat{\theta} = 0$$

regularisation term

$$\boxed{A = \lambda I + (1/n)X^T X}$$   In ridge regression ($\underline{\lambda > 0}$), A is *always* invertible.

# Ridge Regression - Geometric Interpretation

- The **ellipses** correspond to the contours of residual sum of squares (RSS)
  - The inner ellipse has smaller RSS
  - RSS is minimized at ordinal least square (OLS) estimates.

- The **circle** corresponds to the L2 regularization.

- In **ridge regression**, we are trying to minimize the ellipse size and circle *simultaneously*.
  - The ridge estimate ($\hat{\theta}$) is given by the point at which <u>the ellipse and the circle touch</u>.

- The **trade-off** between the penalty (L2) term and RSS:
  - A large $\theta$ would give you a better RSS but then it will push the penalty term higher.
  - While a smaller $\theta$ would a worse RSS.

$$J_{n,\lambda}(\theta) = \frac{\lambda}{2}\|\theta\|^2 + \frac{1}{n}\sum_{t=1}^{n}(y^{(t)} - \theta \cdot x^{(t)})^2/2$$



$\theta_2$
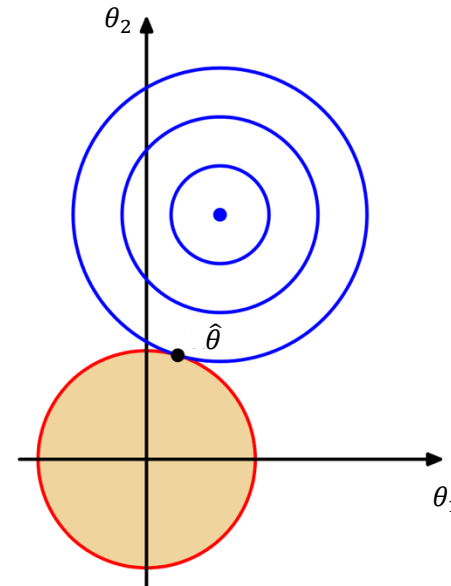
OLS estimate

Ridge estimate
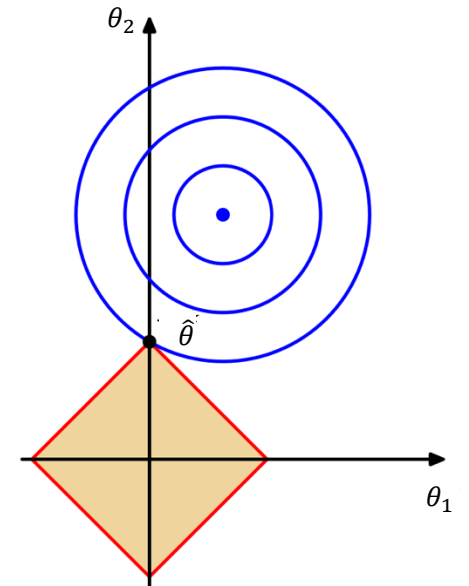
$\theta_1$

# More General Regularizer

- Allow the regularization term to be raised to different powers $p$:

$$J_{n,\lambda}(\theta) = \frac{\lambda}{2} |\theta|^p + R_n(\theta)$$

- $p = 2 \;\longrightarrow\;$ **Ridge Regression**
- $p = 1 \;\longrightarrow\;$ **Lasso Regression**
  - *a.k.a.* **L1 Regularization**
  - if $\lambda$ is sufficiently large, some of the weights $\theta_i$ are driven to zero, leading to a *sparse* model.
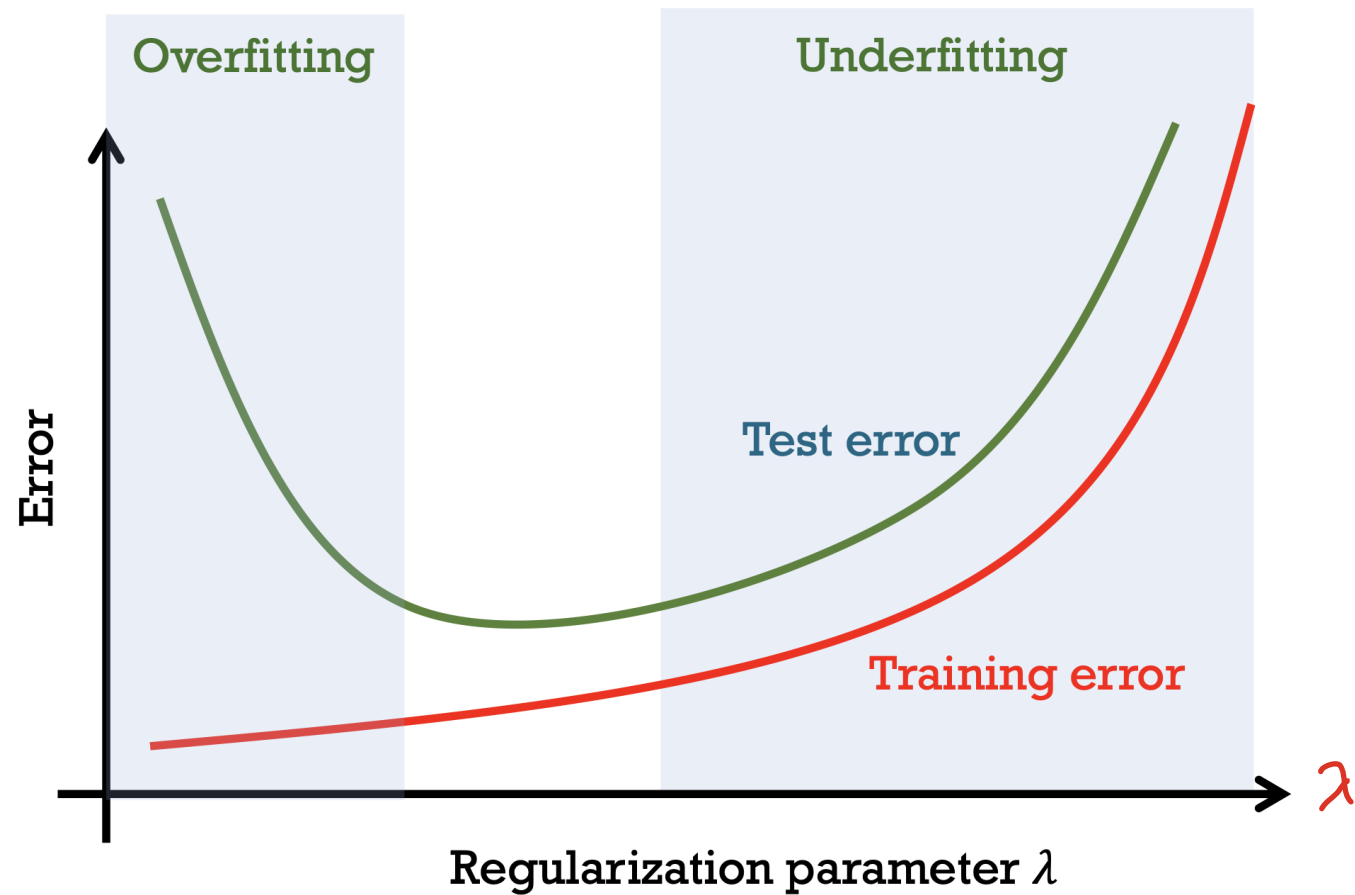  - can be used for **feature selection**.

Ridge Regression

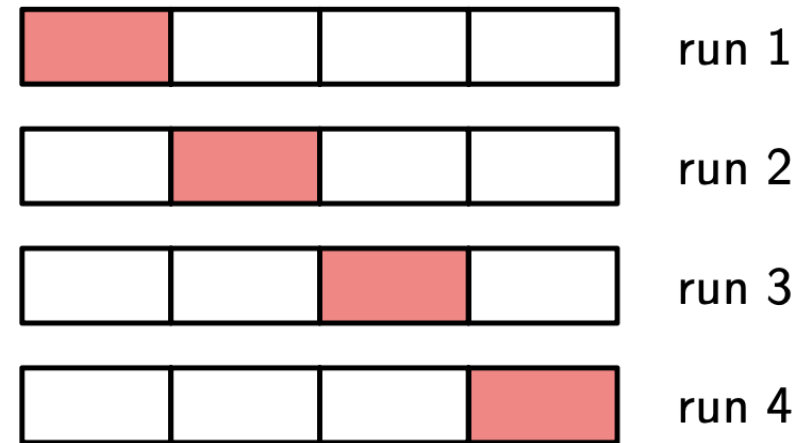Lasso Regression

# Effect of Regularization

# How to pick $\lambda$?

- The regularization parameter $\lambda$ is an example of a *hyperparameter*, which affects the model complexity.

- We don't usually have access to the test data. How to pick optimal $\lambda$ that minimizes the test error?

- The solution is to create a **validation** data set, as a proxy to the test data, and to compute the *validation error*.

# Cross-Validation

- Dividing the dataset into a fixed training and validation sets can be *problematic* if the amount of training data is **limited**:
    - We need as much of the available data as possible for training to build good models.
    - The small validation set will give a relatively noisy estimate of predictive performance.
- **Cross-Validation** is one solution for this dilemma.
    - $k$**-fold cross-validation**
    - Use a proportion $(k-1)/k$ of the available data for training
    - Leave-one-out for validation (assess performance)
    - Repeated for all $k$ possible choices for the held-out group, then average the performance scores from the $k$ runs.

$k$ = no. of folds

run 1

run 2

run 3

run 4

4-fold cross-validation

# Summary

1. What is linear regression?

$$f(x; \theta, \theta_0) = \theta \cdot x + \theta_0 = \sum_{i=1}^{d} \theta_i x_i + \theta_0$$

2. How to learn a linear regression model?

$$R_n(\theta) = \frac{1}{n} \sum_{t=1}^{n} \text{Loss}(y^{(t)} - \theta \cdot x^{(t)}) = \frac{1}{n} \sum_{t=1}^{n} (y^{(t)} - \theta \cdot x^{(t)})^2 / 2$$

Can optimize the cost function using either SGD or closed form solution.

3. What is ridge regression and how is ridge regression different from linear regression?

Add a L2 regularization term to the linear regression cost function.

4. What is regularization and why do we need to do regularization?

Regularization adds constraint on the parameter values to prevent the model overfitting to the training data.

The main aim of regularization is to reduce the over-complexity of the models and help the model learn a simpler function to promote generalization.

# Acknowledgements

- Some slides and content of this lecture are adopted from:

  - MIT 6.036 Introduction to Machine Learning

  - SUTD 50.007 Machine Learning, Spring 2023 (Asst Prof. Malika Meghjani)

  - PennState STAT 897D Applied Data Mining and Statistical Learning

  - Bishop, C. M. (2006). Pattern recognition and machine learning. Springer. (Chapter 3)