# AI6102: Machine Learning Methodologies & Applications

## L3: Linear Models: Regression

**ZHANG Hanwang**

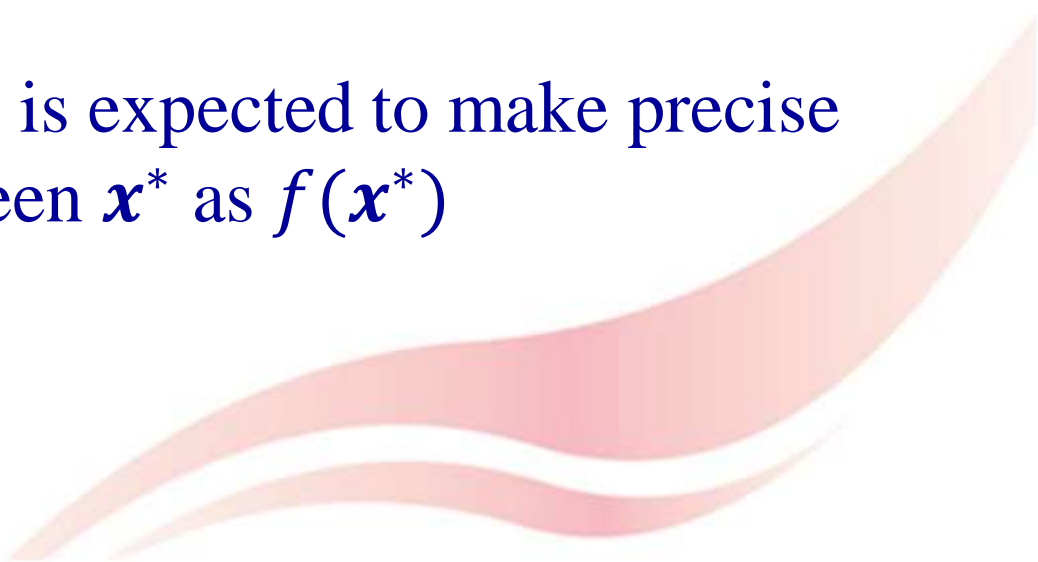**hanwangzhang@ntu.edu.sg**

Nanyang Technological University, Singapore

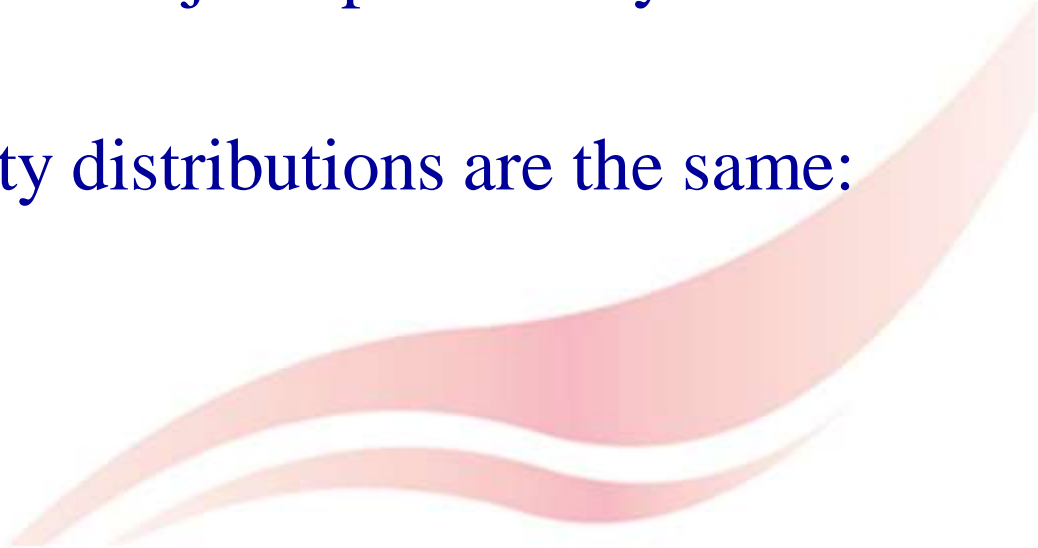Homepage: https://mreallab.github.io/

# Recall: Supervised Learning

In mathematics

- Given: a set of $N$ labeled data $\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)\}$, where $\boldsymbol{x}_i$ is $m$-dimensional vector of numerical values, and $y_i$ is a scalar

- We aim to learn a mapping $f : \boldsymbol{x} \rightarrow y$ by requiring $f(\boldsymbol{x}_i) = y_i$

- The learned mapping $f$ is expected to make precise predictions on any unseen $\boldsymbol{x}^*$ as $f(\boldsymbol{x}^*)$
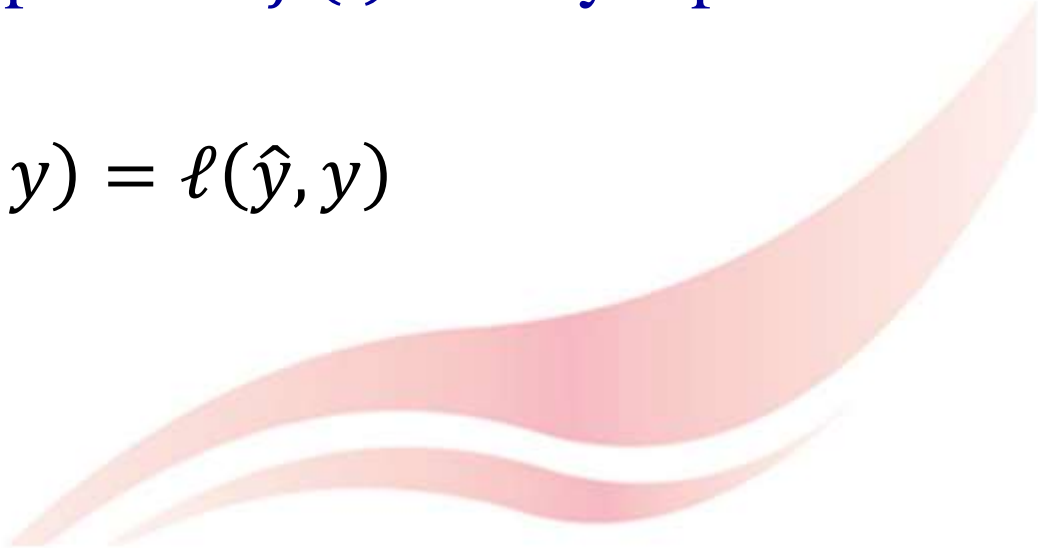
# Hypothesis

- A mapping or function $f: x \rightarrow y$ can be considered as an element of some space of possible functions $\mathcal{H}: \mathbb{R}^m \rightarrow \mathbb{R}$, often called hypothesis space

- Supervised learning aims to find a hypothesis $f \in \mathcal{H}$ from training data $\{(x_1, y_1), \dots, (x_N, y_N)\}$, s.t. for any test $x^*, f(x^*) = y^*$

# Assumption

- The training data instances $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)$, are independent and identically distributed (i.i.d.), and drawn from an unknown joint probability distribution $P_{tr}(\boldsymbol{x}, y)$

  Independent -> probability definition
  Identical ->

- Unseen test data instances $\{(\boldsymbol{x}^*, y^*)\}$ are also i.i.d, and drawn from an unknown joint probability distribution $P_{ts}(\boldsymbol{x}, y)$

- The two joint probability distributions are the same: $P_{tr}(\boldsymbol{x}, y) = P_{ts}(\boldsymbol{x}, y)$

# Loss Function

- Denote by $\hat{y} = f(\boldsymbol{x})$ the prediction of the function $f(\cdot)$ on a data instance $\boldsymbol{x}$, and $y$ is the ground-truth output of $\boldsymbol{x}$

- Let $\ell: \mathbb{R} \times \mathbb{R} \to \mathbb{R}_+ \geq 0$ be a loss function to measure the difference between the ground-truth output $y$ and the prediction $\hat{y}$ of a hypothesis $f(\cdot)$ on any input data instance $\boldsymbol{x}$

$$\ell(f(\boldsymbol{x}), y) = \ell(\hat{y}, y)$$

# Risk Minimization

- The risk associated with a hypothesis $f(\cdot)$ is defined as the expectation of the loss function over all possible input-output pairs drawn from a joint probabilistic distribution $P(\boldsymbol{x}, y)$:

$$R(f) = \mathbb{E}_{(\boldsymbol{x}, y) \sim P}[\ell(f(\boldsymbol{x}), y)]$$

- Recall: in supervised learning, the learned hypothesis $f(\cdot)$ is expected to make precise predictions on any test data instance $\boldsymbol{x}^*$, i.e., $f(\boldsymbol{x}^*) = y^*$

$$f^* = \arg\min_f R_{ts}(f) = \arg\min_f \mathbb{E}_{(\boldsymbol{x}, y) \sim P_{ts}}[\ell(f(\boldsymbol{x}), y)]$$

Test data is unseen in training! And even in the test phase, $y$ is not observed!

# Risk Minimization (cont.)

$$f^* = \arg\min_f R_{ts}(f)$$

$$= \arg\min_f \mathbb{E}_{(\boldsymbol{x},y) \sim P_{ts}} \left[ \ell(f(\boldsymbol{x}), y) \right]$$

$$= \arg\min_f \mathbb{E}_{(\boldsymbol{x},y) \sim P_{ts}} \left[ \frac{P_{tr}(\boldsymbol{x}, y)}{P_{tr}(\boldsymbol{x}, y)} \ell(f(\boldsymbol{x}), y) \right]$$

$$= \arg\min_f \int_y \int_{\boldsymbol{x}} P_{ts}(\boldsymbol{x}, y) \left( \frac{P_{tr}(\boldsymbol{x}, y)}{P_{tr}(\boldsymbol{x}, y)} \ell(f(\boldsymbol{x}), y) \right) d\boldsymbol{x}\, dy$$

Definition of expectation

$$\mathbb{E}_{\boldsymbol{x} \sim P}[g(\boldsymbol{x})] = \int_{\boldsymbol{x}} P(\boldsymbol{x}) g(\boldsymbol{x})\, d\boldsymbol{x} = \int_{x_1} \ldots \int_{x_m} P(\boldsymbol{x}) g(\boldsymbol{x})\, dx_1 \ldots dx_m$$

# Risk Minimization (cont.)

$$f^* = \arg\min_f R_{ts}(f)$$

$$= \arg\min_f \int_y \int_x \boxed{P_{ts}(\boldsymbol{x},y)}\left(\frac{\boxed{P_{tr}(\boldsymbol{x},y)}}{P_{tr}(\boldsymbol{x},y)}\ell(f(\boldsymbol{x}),y)\right)d\boldsymbol{x}\,dy$$

$$= 1 \qquad \text{Assumption } P_{tr}(\boldsymbol{x},y) = P_{ts}(\boldsymbol{x},y)$$

$$= \arg\min_f \int_y \int_x P_{tr}(\boldsymbol{x},y)\left(\boxed{\frac{P_{ts}(\boldsymbol{x},y)}{P_{tr}(\boldsymbol{x},y)}}\ell(f(\boldsymbol{x}),y)\right)d\boldsymbol{x}\,dy$$

$$= \arg\min_f \int_y \int_x P_{tr}(\boldsymbol{x},y)\ell(f(\boldsymbol{x}),y)\,d\boldsymbol{x}\,dy$$

Definition of expectation

$$= \arg\min_f \mathbb{E}_{(\boldsymbol{x},y)\sim P_{tr}}[\ell(f(\boldsymbol{x}),y)]$$

# Empirical Risk Minimization

$$f^* = \arg\min_f \mathbb{E}_{(\boldsymbol{x},y)\sim P_{tr}}[\ell(f(\boldsymbol{x}),y)] = \arg\min_f R_{tr}(f)$$

- The distribution $P_{tr}(\boldsymbol{x},y)$ is unknown, thus we are not able to sample (infinite) input-output pairs $\{(\boldsymbol{x},y)\}$ to learn a hypothesis $f(\cdot)$

- In practice, only a finite number of training pairs are available, $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_N, y_N)$

- Approximate the expected risk by empirical risk:

$$\mathbb{E}_{(\boldsymbol{x},y)\sim P_{tr}}[\ell(f(\boldsymbol{x}),y)] \approx \frac{1}{N}\sum_{i=1}^{N}\ell(f(\boldsymbol{x}_i),y_i) = \hat{R}_{tr}(f)$$

# Empirical Risk Minimization (cont.)

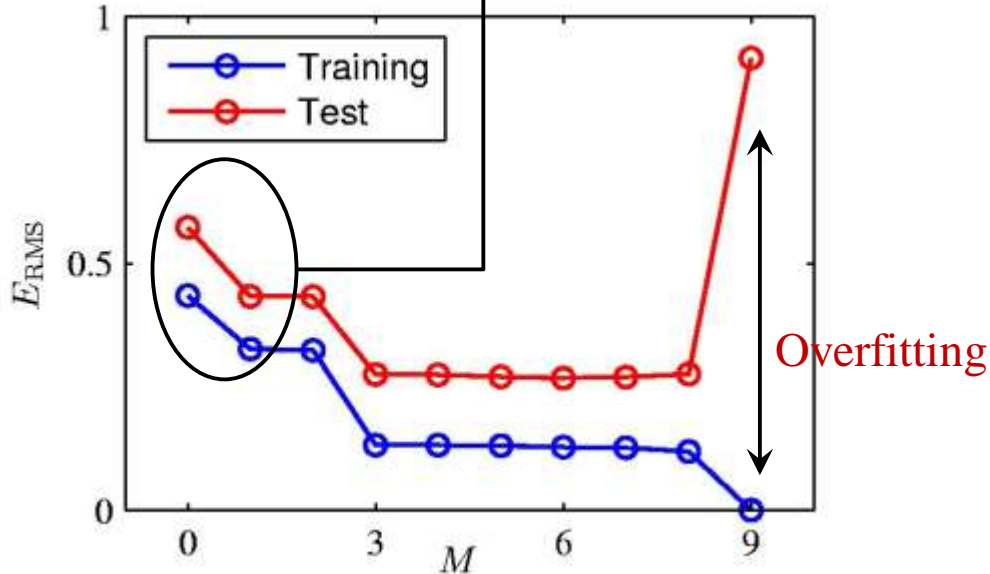- In practice, the hypothesis $f$ can be learned by minimizing the empirical risk

$$\hat{f} = \arg\min_f \hat{R}_{tr}(f) = \arg\min_f \frac{1}{N} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i)$$

- Given a training data set, $N$ is a constant, thus for convenience in presentation, $\frac{1}{N}$ is dropped
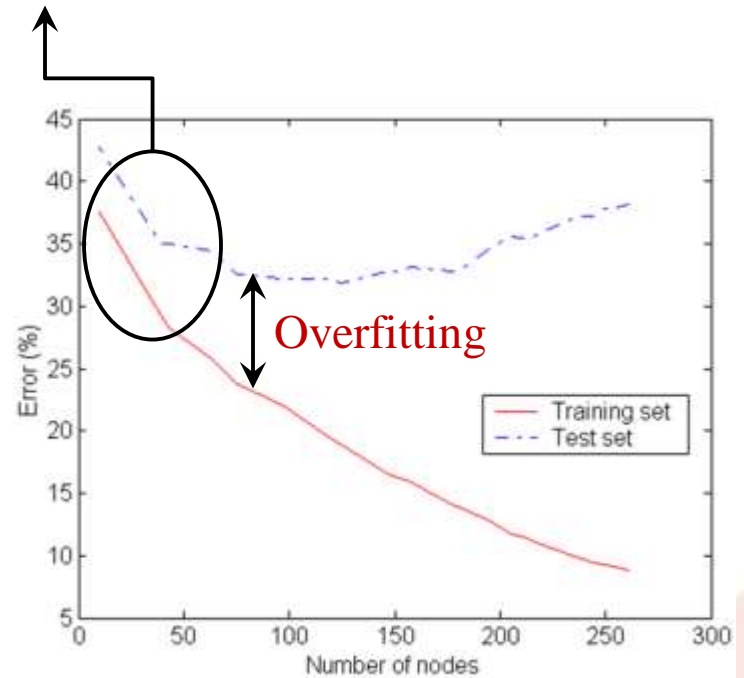
$$\hat{f} = \arg\min_f \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i)$$

# Overfitting Revisit



Underfitting: when model is too simple, both training and test error are large
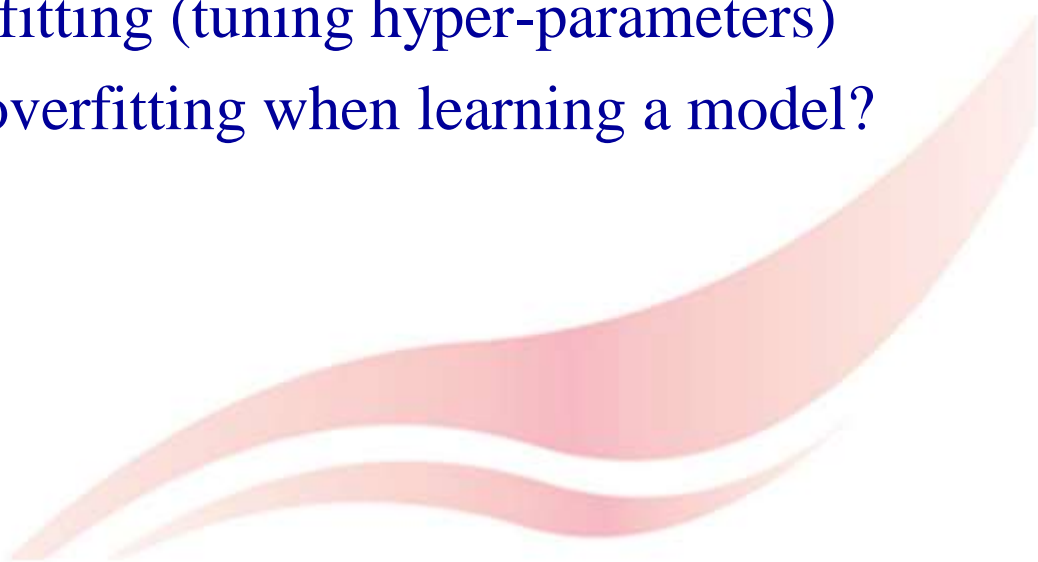
Polynomial curve fitting

Decision tree classification

Overfitting: when test error begins to increase even though training error continues to decrease

# Overfitting v.s. Model Complexity

- Observations of the two examples:
  - Increasing model complexity could make training error or training loss to keep being decreased
  - When model complexity keeps increasing, test error or test loss will increase after some point
- After a model is learned, we can use validation set to evaluate it to reduce the risk of overfitting (tuning hyper-parameters)
- Can we reduce the risk of overfitting when learning a model?

# Occam's Razor Principle

- Given two models of similar performance, we should prefer the simpler model over the more complex model

- For complex models, there is a greater chance that it is fitted accidentally by noise in data
  - Overfitting results in models that are more complex than necessary

- Therefore, we should include model complexity when learning a model

# Structural Risk Minimization (cont.)

- Empirical Risk Minimization

$$\hat{f} = \arg\min_f \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i)$$

- Structural Risk Minimization

$$\hat{f} = \arg\min_f \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i) + \boxed{\lambda\Omega(f)}$$

- $\Omega(f)$ is known as a penalty or regularization term to control the model complexity of $f$
- $\lambda > 0$ is a trade-off hyper-parameter

# Structural Risk Minimization (cont.)

$$\hat{f} = \arg \min_{f} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i) + \lambda \Omega(f)$$

- How to learn $f$?

  - Design a specific form of $f$ in terms of some parameters, denoted by a vector $\boldsymbol{\theta} \in \mathbb{R}^{t \times 1}$, i.e., $f(x; \boldsymbol{\theta})$

  - The parameterized $f(x; \boldsymbol{\theta})$ defines a family of functions with different values of $\boldsymbol{\theta}$

  - Learning $f$ is equivalent to learning the values of $\boldsymbol{\theta}$

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i; \boldsymbol{\theta}), y_i) + \lambda \Omega(\boldsymbol{\theta})$$

# Structural Risk Minimization (cont.)

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i; \boldsymbol{\theta}), y_i) + \lambda \Omega(\boldsymbol{\theta})$$

- Popular regularization terms include
  - the squared L2 norm: $\|\boldsymbol{\theta}\|_2^2$
    - $\|\boldsymbol{\theta}\|_2^2 = \sum_{i=1}^{t} \theta_i^2$
    - Tends to prefer a model with a smaller value for each parameter $\theta_i$
  - the L1 norm: $\|\boldsymbol{\theta}\|_1$
    - $\|\boldsymbol{\theta}\|_1 = \sum_{i=1}^{t} |\theta_i|$
    - Tends to prefer a model with a smaller value for each parameter $\theta_i$, and fewer parameters with non-zero values
    - Induce sparsity, i.e., some $\theta_i$'s tend to be zeros

# Linear Models: Regression

$$\widehat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i; \boldsymbol{\theta}), y_i) + \lambda \Omega(\boldsymbol{\theta})$$

- In general, for regression, $f(\boldsymbol{x}_i; \boldsymbol{\theta})$ is defined as

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \boldsymbol{w} \cdot \boldsymbol{x} + b$$

$\boldsymbol{\theta}$ is a concatenation of $\boldsymbol{w}$ and $b$

- Given $\boldsymbol{x}_i$, the prediction of $f(\boldsymbol{x}; \boldsymbol{\theta})$ is the linear combination of its $m$ feature values with weights $\boldsymbol{w}$ plus a bias term $b$

$\boldsymbol{x}_i$

| $x_{1i}$ |
| $x_{2i}$ |
| ... |
| $x_{mi}$ |

$\boldsymbol{w}$

| $w_1$ |
| $w_2$ |
| ... |
| $w_m$ |

$$\hat{y} = f(\boldsymbol{x}; \boldsymbol{\theta}) = \boldsymbol{w} \cdot \boldsymbol{x} + b = \sum_{i=1}^{m} x_i w_i + b$$

# Linear Models: Classification

$$\widehat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i; \boldsymbol{\theta}), y_i) + \lambda \Omega(\boldsymbol{\theta})$$

- In general, for classification, $f(\boldsymbol{x}_i; \boldsymbol{w})$ is defined as

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = h(\boldsymbol{w} \cdot \boldsymbol{x} + b)$$

  where $h(z)$ is function to map continuous values to discrete values (denoting different categories)
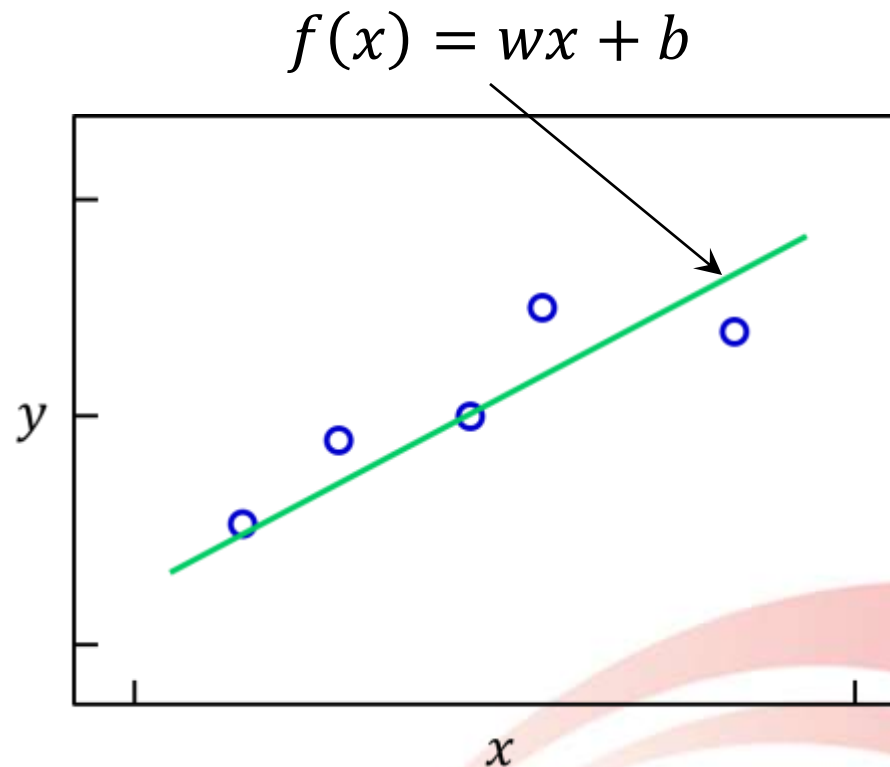
- For example,

$$h(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{if } z < 0 \end{cases}$$

**Next Lecture**

# Linear Regression: One-Dimension

- Each instance is represented by only one input feature
- To learn a linear function $f(x)$ in terms of $w$ and $b$ (both are scalars) from $\{x_i, y_i\}, i = 1, \dots, N$

$$f(x) = wx + b$$

# 1D Linear Regression (cont.)
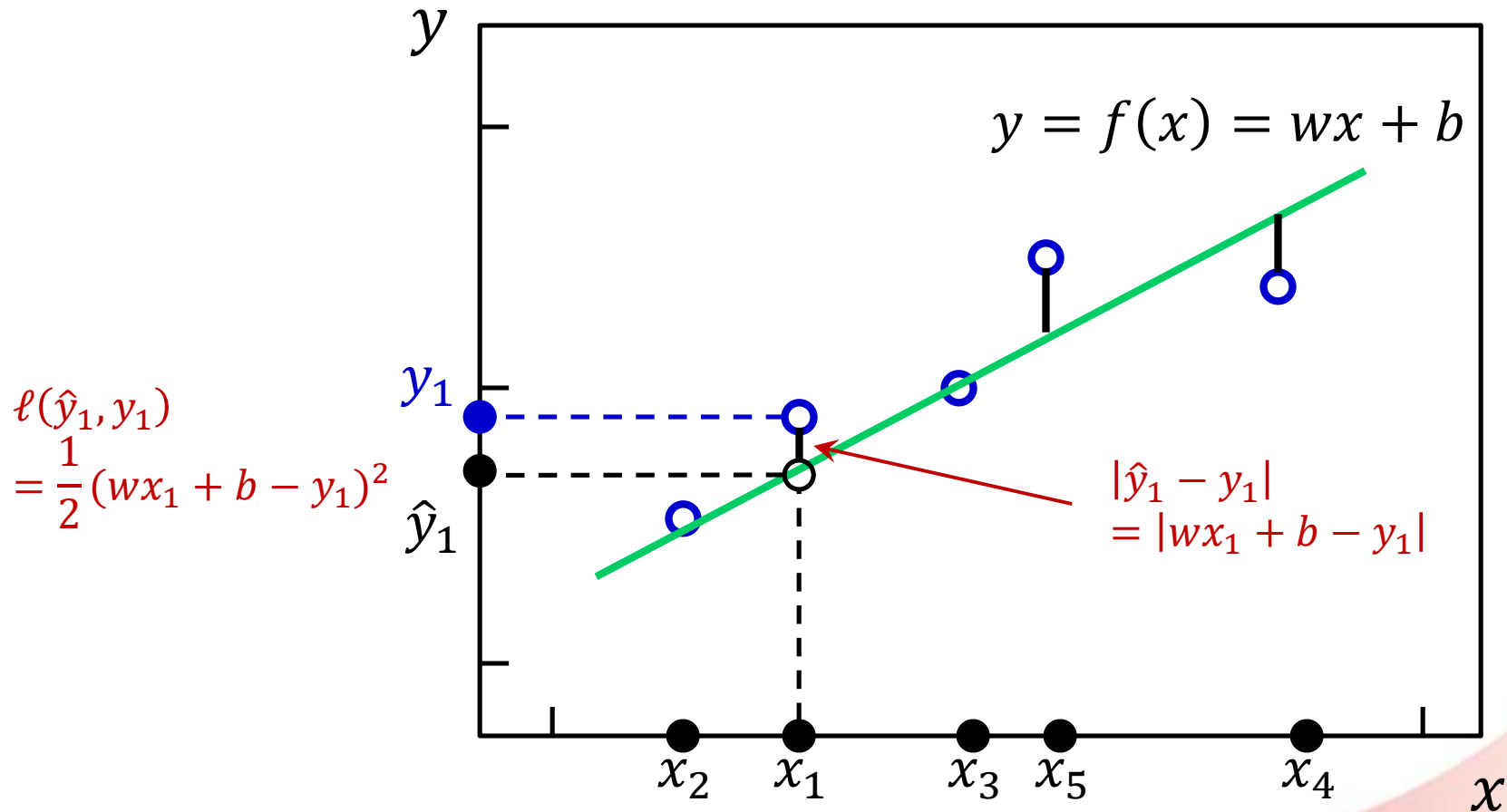
$$[\widehat{w}, \widehat{b}] = \arg \min_{[w,b]} \sum_{i=1}^{N} \ell(\hat{y}_i, y_i)$$

<span style="color:red">Drop the regularization term for simplicity at first</span>

where $\hat{y}_i = wx_i + b$

- The loss function $\ell(\hat{y}_i, y_i)$ is to measure the difference between $\hat{y}_i$ and $y_i$
  - For regression, the magnitude of the difference, i.e., $|\hat{y}_i - y_i|$
- To make the resultant optimization problem easier to solve
  - We expect the loss function has some good properties, e.g., differentiable everywhere
  - The square of magnitude, $|\hat{y}_i - y_i|^2 = (\hat{y}_i - y_i)^2$ or $\frac{1}{2}(\hat{y}_i - y_i)^2$

# Regression Loss Function

$$y = f(x) = wx + b$$

$$\ell(\hat{y}_1, y_1) = \frac{1}{2}(wx_1 + b - y_1)^2$$

$$|\hat{y}_1 - y_1| = |wx_1 + b - y_1|$$

$$\frac{1}{2}\sum_{i=1}^{5} \ell(\hat{y}_i, y_i) = \frac{1}{2}\sum_{i=1}^{5}(wx_i + b - y_i)^2$$
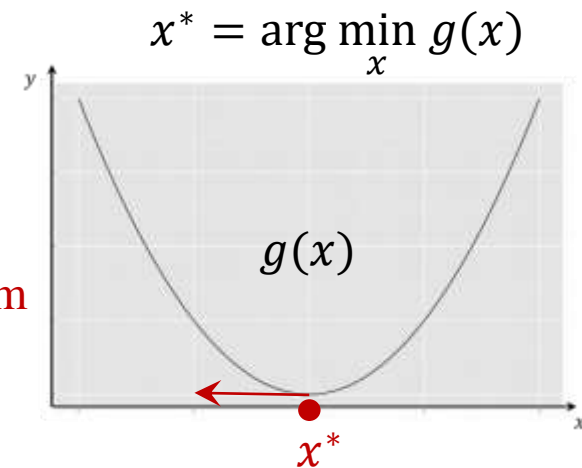
# Optimization

- Learn $w$ and $b$ by minimizing the square loss

$$[\hat{w}, \hat{b}] = \arg\min_{[w,b]} \boxed{\frac{1}{2}\sum_{i=1}^{N}(wx_i + b - y_i)^2}$$

- The objective of the optimization problem
- The objective is convex

$$x^* = \arg\min_x g(x)$$

$g(x)$

$x^*$

- Unconstrained optimization problem
- Set the derivatives of the objective w.r.t. $w$ and $b$ to zero, respectively
- $\hat{w}$ and $\hat{b}$ can be obtained by solving the equations

# Closed-form Solution

$$\begin{cases} \dfrac{\partial \left( \frac{1}{2} \sum_{i=1}^{N} (wx_i + b - y_i)^2 \right)}{\partial w} = 0 \\[3em] \dfrac{\partial \left( \frac{1}{2} \sum_{i=1}^{N} (wx_i + b - y_i)^2 \right)}{\partial b} = 0 \end{cases}$$

$$\begin{cases} \sum_{i=1}^{N} (wx_i + b - y_i)x_i = 0 \\[2em] \sum_{i=1}^{N} (wx_i + b - y_i) = 0 \end{cases}$$

Chain rule of calculus

$$y = g(x)$$

$$z = f(y) = f\big(g(x)\big)$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$$z_i = wx_i + b - y_i$$

$$\frac{\partial z_i^2}{\partial w} = \frac{\partial z_i^2}{\partial z_i} \frac{\partial z_i}{\partial w}$$

$$= 2z_i \frac{\partial(wx_i + b - y_i)}{\partial w}$$

$$= 2z_i x_i$$

# Closed-form Solution (cont.)

$$\sum_{i=1}^{N}(wx_i + b - y_i)x_i = 0$$

$$\sum_{i=1}^{N}(wx_i + b - y_i) = 0$$

$$\sum_{i=1}^{N}(wx_i + b - y_i)x_i = 0$$

$$\sum_{i=1}^{N}(wx_i - y_i) + Nb = 0$$

$$b = \frac{1}{N}\sum_{i=1}^{N}(y_i - wx_i)$$

# Closed-form Solution (cont.)

$$\sum_{i=1}^{N}(wx_i + b - y_i)x_i = 0$$

$$b = \frac{1}{N}\sum_{i=1}^{N}(y_i - wx_i)$$

$$\sum_{i=1}^{N}(wx_i + b - y_i)x_i = 0$$

$$b = \frac{1}{N}\sum_{i=1}^{N}y_i - w\frac{1}{N}\sum_{i=1}^{N}x_i = \bar{y} - w\bar{x}$$

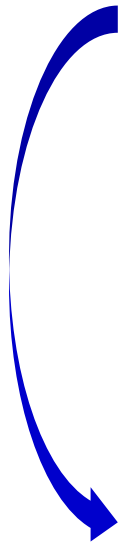# Closed-form Solution (cont.)

$$\sum_{i=1}^{N}(wx_i + b - y_i)x_i = 0$$

$$b = \bar{y} - w\bar{x}$$

$$\sum_{i=1}^{N}(wx_i + \bar{y} - w\bar{x} - y_i)x_i = 0$$

$$b = \bar{y} - w\bar{x}$$

# Closed-form Solution (cont.)

$$\sum_{i=1}^{N}(wx_i + \bar{y} - w\bar{x} - y_i)x_i = 0$$

$$b = \bar{y} - w\bar{x}$$

$$w\sum_{i=1}^{N}(x_i - \bar{x})x_i = \sum_{i=1}^{N}(y_i - \bar{y})x_i$$
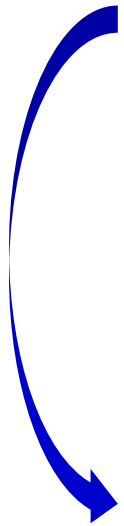
$$b = \bar{y} - w\bar{x}$$

# Closed-form Solution (cont.)

$$w \sum_{i=1}^{N} (x_i - \bar{x}) x_i = \sum_{i=1}^{N} (y_i - \bar{y}) x_i$$

$$b = \bar{y} - w\bar{x}$$

$$w = \frac{\sum_{i=1}^{N} (y_i - \bar{y}) x_i}{\sum_{i=1}^{N} (x_i - \bar{x}) x_i}$$

$$b = \bar{y} - w\bar{x}$$

# Multi-Dimension Case

- Each instance has $m$ dimensions, a linear function $f(x)$ is defined as

$$f(x) = w \cdot x + b$$

- By defining $w_0 = b$, and $x_0 = 1$, $f(x)$ can be rewritten as

$$f(x) = \boxed{w \cdot x}$$

Both $w$ and $x$ have $m + 1$ dimensions

$x_i$

| $x_{0i} = 1$ |
|---|
| $x_{1i}$ |
| $x_{2i}$ |
| ... |
| $x_{mi}$ |

$w$

| $w_0 = b$ |
|---|
| $w_1$ |
| $w_2$ |
| ... |
| $w_m$ |

$$= \sum_{k=0}^{m} x_{ki} w_k$$

$$= \sum_{k=1}^{m} x_{ki} w_k + x_{0i} w_0$$

# Optimization

- Learn $\boldsymbol{w}$ by minimizing the total square loss

$$\widehat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} \frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2$$

- A closed-form solution can be obtained by setting the derivative of the objective w.r.t. $\boldsymbol{w}$ to zero, and solving the resultant equations

$$\frac{\partial\left(\frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2\right)}{\partial\boldsymbol{w}} = \boldsymbol{0}$$

# Brief Linear Algebra Review

- Linear algebra plays a crucial role in deriving solutions for various machine learning methods

- You are highly recommended to refer to Part I of the Deep Learning book at https://www.deeplearningbook.org/

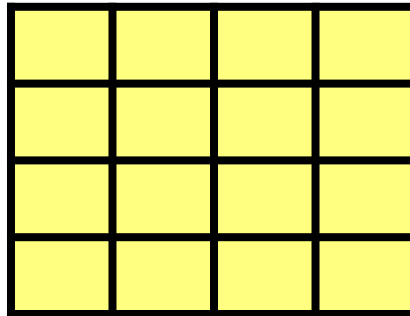- Transpose of a vector or matrix

$x$ ($m \times 1$)

$x^T$ ($1 \times m$)

$\mathbf{A}$ ($m \times N$)

$\mathbf{A}^T$ ($N \times m$)

# Matrix/Vector Concepts

- Square matrix
  - If a matrix **A** has the same number of rows and columns, then it is said to be square matrix

$$\mathbf{A}\ (m \times m)$$

- Symmetric matrix
  - If a **square** matrix **A** satisfies $\mathbf{A} = \mathbf{A}^T$

# Matrix Multiplication

- Matrix multiplication is associative
  - $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$
- Matrix multiplication is distributive
  - $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$
- Matrix multiplication is NOT commutative in general
  - $\mathbf{AB} \neq \mathbf{BA}$
- The identity matrix, $\mathbf{I}$ $(m \times m)$, is a symmetric matrix with ones on the diagonal and zeros everywhere else
  - If $\mathbf{A}$ is a square matrix $(m \times m)$: $\mathbf{AI} = \mathbf{IA} = \mathbf{A}$
  - If $\mathbf{A}$ is $(N \times m)$: $\mathbf{AI} = \mathbf{A}$
  - If $\mathbf{A}$ is $(m \times N)$: $\mathbf{IA} = \mathbf{A}$

$\mathbf{I}$ $(m \times m)$

| 1 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

# Matrix Operations

- The transpose of $\mathbf{A}^T$:
$$(\mathbf{A}^T)^T = \mathbf{A}$$

- The transpose of $\mathbf{AB}$:
$$(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T$$

- The transpose of $\mathbf{A}\boldsymbol{x}$
$$(\mathbf{A}\boldsymbol{x})^T = \boldsymbol{x}^T\mathbf{A}^T$$

- The transpose of $\boldsymbol{x}^T\boldsymbol{y}$
$$(\boldsymbol{x}^T\boldsymbol{y})^T = \boldsymbol{y}^T\boldsymbol{x}$$

- The transpose of a scalar is the scalar itself
$$a^T = a$$

# Matrix Operations (cont.)

- For a square matrix $\mathbf{A}$ $(m \times m)$, if it is invertible, then there exists a unique matrix, denoted by $\mathbf{A}^{-1}$, such that
$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

- If it is not invertible, then such a matrix $\mathbf{A}^{-1}$ does not exist

- Non-square matrices do not have inverses by definition

- Properties of the inverse ($\mathbf{A}$ and $\mathbf{B}$ are invertible)
  - $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
  - $(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$
  - $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$

# Linear System

- Given the following system of linear equations

$$2x_1 + x_3 = 5$$
$$3x_1 - 4x_2 + 2x_3 = 4$$
$$2x_2 - 3x_3 = -3$$
$$-x_1 + 2x_2 - 5x_3 = 1$$

- They can be written in a more compact form as

$$\mathbf{A}x = b$$

$$\mathbf{A}\ (4 \times 3)$$

| 2 | 0 | 1 |
|----|----|----|
| 3 | −4 | 2 |
| 0 | 2 | −3 |
| −1 | 2 | −5 |

$$x\ (3 \times 1)$$

| $x_1$ |
|----|
| $x_2$ |
| $x_3$ |

$$b\ (4 \times 1)$$

| 5 |
|----|
| 4 |
| −3 |
| 1 |

# Linear System (cont.)

$$\mathbf{A}x = b$$

$$\mathbf{A} \; (k \times d) \qquad x \; (d \times 1) \qquad b \; (k \times 1)$$

- If **A** is square, and invertible, then we multiply both sides of the equation by $\mathbf{A}^{-1}$ to obtain a **unique** solution
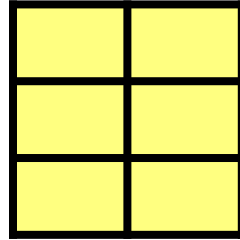
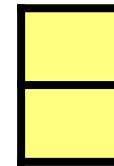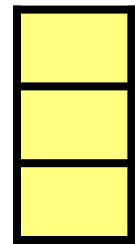$$\mathbf{A}^{-1}\mathbf{A}x = \mathbf{A}^{-1}b \implies x = \mathbf{A}^{-1}b$$

# Linear System (cont.)

$$\mathbf{A}x = b$$

$\mathbf{A}\ (k \times d)$

$x\ (d \times 1)$

$b\ (k \times 1)$

- If **A** is not invertible, solutions are **not** unique, we can find a solution by using the pseudo inverse (also known as generalized inverse) of **A** instead, denoted by $\mathbf{A}^{\dagger}$

$$x = \mathbf{A}^{\dagger} b$$

- Special case: when **A** is square $(k \times k)$ but not invertible

$$\mathbf{A}^{\dagger}\mathbf{A} = \mathbf{A}\mathbf{A}^{\dagger} = \mathrm{diag}(\lambda_1, \dots, \lambda_k), \text{where } \lambda_i \in \{0,1\}, \text{at least one } \lambda_i = 0$$

$$\mathbf{A}\mathbf{A}^{\dagger} = \mathbf{A}^{\dagger}\mathbf{A}$$

| 1 | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

# Closed-form Solution for Linear Regression

$$\frac{\partial \left(\frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2\right)}{\partial \boldsymbol{w}} = 0$$

$$\frac{1}{2}\sum_{i=1}^{N}\frac{\partial(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2}{\partial \boldsymbol{w}} = 0$$

$$\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)\boldsymbol{x}_i = 0$$

$$\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i)\boldsymbol{x}_i - \sum_{i=1}^{N}y_i\boldsymbol{x}_i = 0$$

$$z_i = \boldsymbol{w}\cdot\boldsymbol{x}_i - y_i$$

$$\frac{\partial z_i^2}{\partial \boldsymbol{w}} = \frac{\partial z_i^2}{\partial z_i}\frac{\partial z_i}{\partial \boldsymbol{w}} = 2z_i\frac{\partial(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)}{\partial \boldsymbol{w}}$$

$$\frac{\partial(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)}{\partial \boldsymbol{w}} = \frac{\partial(\boldsymbol{w}\cdot\boldsymbol{x}_i)}{\partial \boldsymbol{w}} - 0 = \boldsymbol{x}_i$$

The Matrix Cookbook
https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf

# Closed-form Solution (cont.)

$$a\boldsymbol{x} = \boldsymbol{x}a$$

$$\sum_{i=1}^{N} \boxed{(\boldsymbol{w} \cdot \boldsymbol{x}_i)} \boldsymbol{x}_i - \sum_{i=1}^{N} y_i \boldsymbol{x}_i = \mathbf{0}$$

scalar

$$\sum_{i=1}^{N} \boxed{\boldsymbol{x}_i(\boldsymbol{w} \cdot \boldsymbol{x}_i)} - \sum_{i=1}^{N} y_i \boldsymbol{x}_i = \mathbf{0}$$

$$\boldsymbol{x}_i(\boldsymbol{w} \cdot \boldsymbol{x}_i) = \boldsymbol{x}_i(\boldsymbol{w}^T \boldsymbol{x}_i) = \boldsymbol{x}_i(\boldsymbol{x}_i^T \boldsymbol{w}) = (\boldsymbol{x}_i \boldsymbol{x}_i^T)\boldsymbol{w}$$

$$\left( \sum_{i=1}^{N} \boxed{(\boldsymbol{x}_i \boldsymbol{x}_i^T)} \right) \boldsymbol{w} - \sum_{i=1}^{N} y_i \boldsymbol{x}_i = \mathbf{0}$$

$\boldsymbol{x}_i \in \mathbb{R}^{(m+1) \times 1}$ and $\boldsymbol{x}_i^T \in \mathbb{R}^{1 \times (m+1)}$, thus
$\boldsymbol{x}_i \boldsymbol{x}_i^T$ is a $(m+1)$ by $(m+1)$ matrix

# Closed-form Solution (cont.)

$$\left(\sum_{i=1}^{N} \boxed{(x_i x_i^T)}\right) w - \sum_{i=1}^{N} y_i x_i = 0$$

$x_i \in \mathbb{R}^{(m+1) \times 1}$ and $x_i^T \in \mathbb{R}^{1 \times (m+1)}$, thus
$x_i x_i^T$ is a $(m+1)$ by $(m+1)$ matrix

$(m+1) \times 1$

$x_i^T x_i$ — Inner product, scalar

$x_i$

| $x_{0i}$ |
|----------|
| $x_{1i}$ |
| ... |
| $x_{mi}$ |

$x_i^T$

| $x_{0i}$ | $x_{1i}$ | ... | $x_{mi}$ |
|----------|----------|-----|----------|

$1 \times (m+1)$

$x_i x_i^T$   $(m+1)$ by $(m+1)$ matrix

| $x_{0i}x_{0i}$ | $x_{0i}x_{1i}$ | ... | $x_{0i}x_{mi}$ |
|----------------|----------------|-----|----------------|
| $x_{1i}x_{0i}$ | $x_{1i}x_{1i}$ | ... | $x_{1i}x_{mi}$ |
| ... | ... | ... | ... |
| $x_{mi}x_{0i}$ | $x_{mi}x_{1i}$ | ... | $x_{mi}x_{mi}$ |

# Closed-form Solution (cont.)

$$\left(\sum_{i=1}^{N}\left(x_i x_i^T\right)\right)w - \sum_{i=1}^{N} y_i x_i = 0$$

$$x_i x_i^T$$

$(m+1)$ by $(m+1)$ matrix

| $x_{0i}x_{0i}$ | $x_{0i}x_{1i}$ | ... | $x_{0i}x_{mi}$ |
|---|---|---|---|
| $x_{1i}x_{0i}$ | $x_{1i}x_{1i}$ | ... | $x_{1i}x_{mi}$ |
| ... | ... | ... | ... |
| $x_{mi}x_{0i}$ | $x_{mi}x_{1i}$ | ... | $x_{mi}x_{mi}$ |

$$\sum_{i=1}^{N}\left(x_i x_i^T\right)$$

| $\sum_{i=1}^{N} x_{0i}x_{0i}$ | $\sum_{i=1}^{N} x_{0i}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{0i}x_{mi}$ |
|---|---|---|---|
| $\sum_{i=1}^{N} x_{1i}x_{0i}$ | $\sum_{i=1}^{N} x_{1i}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{1i}x_{mi}$ |
| ... | ... | ... | ... |
| $\sum_{i=1}^{N} x_{mi}x_{0i}$ | $\sum_{i=1}^{N} x_{mi}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{mi}x_{mi}$ |

# Closed-form Solution (cont.)

$$\sum_{i=1}^{N}(x_i x_i^T)$$

| $\sum_{i=1}^{N} x_{0i}x_{0i}$ | $\sum_{i=1}^{N} x_{0i}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{0i}x_{mi}$ |
|---|---|---|---|
| $\sum_{i=1}^{N} x_{1i}x_{0i}$ | $\sum_{i=1}^{N} x_{1i}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{1i}x_{mi}$ |
| ... | ... | ... | ... |
| $\sum_{i=1}^{N} x_{mi}x_{0i}$ | $\sum_{i=1}^{N} x_{mi}x_{1i}$ | ... | $\sum_{i=1}^{N} x_{mi}x_{mi}$ |

$$\sum_{i=1}^{N}(x_i x_i^T) = \mathbf{XX}^T$$

$(m+1)$ by $N$           $N$ by $(m+1)$

**X**

| $x_{01}$ | $x_{02}$ | ... | $x_{0N}$ |
|---|---|---|---|
| $x_{11}$ | $x_{12}$ | ... | $x_{1N}$ |
| ... | ... | ... | ... |
| $x_{m1}$ | $x_{m2}$ | ... | $x_{mN}$ |

$\mathbf{X}^T$

| $x_{01}$ | $x_{11}$ | ... | $x_{m1}$ |
|---|---|---|---|
| $x_{02}$ | $x_{12}$ | ... | $x_{m2}$ |
| ... | ... | ... | ... |
| $x_{0N}$ | $x_{1N}$ | ... | $x_{mN}$ |

# Closed-form Solution (cont.)

$$\left( \sum_{i=1}^{N} \left( \boldsymbol{x}_i \boldsymbol{x}_i^T \right) \right) \boldsymbol{w} - \sum_{i=1}^{N} y_i \boldsymbol{x}_i = 0$$

$$\mathbf{X}\mathbf{X}^T \boldsymbol{w} - \sum_{i=1}^{N} y_i \boldsymbol{x}_i = 0$$

$$\mathbf{X}\mathbf{X}^T \boldsymbol{w} - \mathbf{X}\boldsymbol{y} = 0$$

$(m + 1)$ by $N$

**X**

| $x_{01}$ | $x_{02}$ | ... | $x_{0N}$ |
|---|---|---|---|
| $x_{11}$ | $x_{12}$ | ... | $x_{1N}$ |
| ... | ... | ... | ... |
| $x_{m1}$ | $x_{m2}$ | ... | $x_{mN}$ |

$N$ by 1

**$\boldsymbol{y}$**

| $y_1$ |
|---|
| $y_2$ |
| ... |
| $y_N$ |

$$\sum_{i=1}^{N} y_i \boldsymbol{x}_i = \mathbf{X}\boldsymbol{y}$$

# Closed-form Solution (cont.)

$$\mathbf{X}\mathbf{X}^T w - \mathbf{X}y = 0$$

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

When $\mathbf{X}\mathbf{X}^T$ is invertible

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

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

What if $\mathbf{X}\mathbf{X}^T$ is NOT invertible?

# When $XX^T$ is NOT Invertible

- We can use the pseudo inverse (also known as generalized inverse) of $\mathbf{XX}^T$ instead, i.e., $(\mathbf{XX}^T)^\dagger$

- In this case

$$w = (\mathbf{XX}^T)^\dagger \mathbf{X}y$$

# Regularized Linear Regression

- Recall structural risk minimization

$$\hat{f} = \arg \min_f \sum_{i=1}^{N} \ell(f(\boldsymbol{x}_i), y_i) + \lambda \Omega(f)$$

$\lambda > 0$ tradeoff
hyper-parameter

$$\|\boldsymbol{w}\|_2^2 = \boldsymbol{w} \cdot \boldsymbol{w}$$

$$\widehat{\boldsymbol{w}} = \arg \min_{\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{N} (\boldsymbol{w} \cdot \boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

A regularization term to control
the complexity of the model

Also known as Ridge regression

# Optimization

$$\widehat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} \frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2$$

- Still an unconstrained optimization problem, and the objective is still convex

- A closed-form solution can be obtain by setting the derivative of the objective w.r.t. $\boldsymbol{w}$, and solving the resultant equation

$$\frac{\partial\left(\frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2\right)}{\partial\boldsymbol{w}} = \mathbf{0}$$

# Closed-form Solution

$$\frac{\partial \left(\frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2\right)}{\partial \boldsymbol{w}} = \boldsymbol{0}$$

$$\frac{\partial \frac{1}{2}\sum_{i=1}^{N}(\boldsymbol{w}\cdot\boldsymbol{x}_i - y_i)^2}{\partial \boldsymbol{w}} + \frac{\partial \frac{\lambda}{2}\|\boldsymbol{w}\|_2^2}{\partial \boldsymbol{w}} = \boldsymbol{0}$$

$$\frac{\partial \|\boldsymbol{w}\|_2^2}{\partial \boldsymbol{w}} = \frac{\partial (\boldsymbol{w}\cdot\boldsymbol{w})}{\partial \boldsymbol{w}} = 2\boldsymbol{w}$$

The matrix cookbook

$$\left(\sum_{i=1}^{N}(\boldsymbol{x}_i\boldsymbol{x}_i^T)\right)\boldsymbol{w} - \sum_{i=1}^{N} y_i\boldsymbol{x}_i + \lambda\boldsymbol{w} = \boldsymbol{0}$$

# Closed-form Solution (cont.)

$$\left( \sum_{i=1}^{N} \left( \boldsymbol{x}_i \boldsymbol{x}_i^T \right) \right) \boldsymbol{w} - \sum_{i=1}^{N} y_i \boldsymbol{x}_i + \lambda \boldsymbol{w} = \mathbf{0}$$

$$\mathbf{X}\mathbf{X}^T \boldsymbol{w} - \mathbf{X}\boldsymbol{y} + \lambda \mathbf{I}\boldsymbol{w} = \mathbf{0}$$

$$(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})\boldsymbol{w} - \mathbf{X}\boldsymbol{y} = \mathbf{0}$$

Always invertible as
long as $\lambda$ is positive

$$(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})\boldsymbol{w} = \mathbf{X}\boldsymbol{y} \qquad \Longrightarrow \qquad \boldsymbol{w} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1}\mathbf{X}\boldsymbol{y}$$

# Why $XX^T + \lambda I$ Invertible?

- A square matrix is invertible if and only if it does not have a zero eigenvalue

- If a symmetric matrix **A** is <u>positive semidefinite</u>, then all of its eigenvalues are non-negative ($\geq 0$)
  - When a symmetric matrix **A** ($d \times d$) is said to be positive semidefinite iif for any non-zero column vector $\boldsymbol{x}$ ($d \times 1$), $\boldsymbol{x}^T \mathbf{A} \boldsymbol{x} \geq 0$

- If a symmetric matrix **A** is <u>positive definite</u>, then all of its eigenvalues are positive ($> 0$)
  - When a symmetric matrix **A** ($d \times d$) is said to be positive definite iif for any non-zero column vector $\boldsymbol{x}$ ($d \times 1$), $\boldsymbol{x}^T \mathbf{A} \boldsymbol{x} > 0$

- A positive definite matrix is invertible (all eigenvalues $> 0$, and of course non-zero)

# Why $\mathbf{XX}^T + \lambda \mathbf{I}$ Invertible? (cont.)

- $\mathbf{XX}^T$ is a positive semidefinite
  - Need to prove for any non-zero $(m+1)$- dimensional column vector $\mathbf{z}$, $\mathbf{z}^T \mathbf{XX}^T \mathbf{z} \geq 0$
  - Proof: Denote $\mathbf{y} = \mathbf{X}^T \mathbf{z}$. Then $\mathbf{z}^T \mathbf{XX}^T \mathbf{z} = \boxed{\mathbf{y}^T \mathbf{y}}$   $\|\mathbf{y}\|_2^2 \geq 0$
- $\mathbf{XX}^T + \lambda \mathbf{I}$ is positive definite if $\lambda > 0$    $\|\mathbf{y}\|_2^2 = 0$ iff $\mathbf{y} = \mathbf{X}^T \mathbf{z} = \mathbf{0}$
  - Need to prove for any non-zero $(m+1)$- dimensional column vector $\mathbf{z}$, $\mathbf{z}^T (\mathbf{XX}^T + \lambda \mathbf{I})\mathbf{z} = \mathbf{z}^T \mathbf{XX}^T \mathbf{z} + \lambda \mathbf{z}^T \mathbf{z} > 0$
  - Proof: $\mathbf{z}^T (\mathbf{XX}^T + \lambda \mathbf{I})\mathbf{z} = \boxed{\mathbf{z}^T \mathbf{XX}^T \mathbf{z}} + \boxed{\lambda \mathbf{z}^T \mathbf{z}} > 0$

  $\qquad\qquad\qquad\qquad\quad \geq 0 \qquad\quad > 0$ since $\lambda > 0$
  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ and $\mathbf{z}$ is non-zero

# Large-scale Issue

$$(m + 1) \times (m + 1)$$

$$w = \boxed{(\mathbf{XX}^T + \lambda \mathbf{I})}^{-1} \mathbf{X} y$$

- The computation complexity of computing an inverse of a $(m + 1) \times (m + 1)$ matrix is $O\big((m + 1)^3\big)$

- When $m$ is large, it is time consuming

- Rather than computing the inverse to obtain a closed-form solution, consider the following linear system

$$(\mathbf{XX}^T + \lambda \mathbf{I})w = \mathbf{X}y$$       Linear system: $\mathbf{A}x = b$

- We can solve it by using various numerical methods, e.g., Gaussian elimination, etc.

# Large-scale Issue (cont.)

$$\min_{\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^{N} (\boldsymbol{w} \cdot \boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

- Alternatively, rather than trying to derive an analytical solution, we can apply numerical methods to iteratively minimize the objective, e.g., gradient descent

- Denote the objective by

$$E(\boldsymbol{w}) = \frac{1}{2} \sum_{i=1}^{N} (\boldsymbol{w} \cdot \boldsymbol{x}_i - y_i)^2 + \frac{\lambda}{2} \|\boldsymbol{w}\|_2^2$$

Objective to be minimized

$$\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \rho \frac{\partial E(\boldsymbol{w})}{\partial \boldsymbol{w}}$$

Learning rate
$\rho \in (0,1]$

# Gradient Descent

$$w^* = \arg \min_w E(w)$$



$$w_{i+1} = w_i - \rho \left. \frac{\partial E(w)}{\partial w} \right|_{w=w_i}$$

$$i = i + 1$$

$E(w)$

$E(w_1)$

$w^*$

$w_3 \ w_2 \ w_1$

$w$

# Implementation using scikit-learn

- API: sklearn.linear_model: Linear Models
  https://scikit-learn.org/stable/modules/classes.html#module-sklearn.linear_model

  – Classical linear regressors

  – linear_model.LinearRegression → linear regression without regularization

  – linear_model.Ridge → regularized linear regression

## Classical linear regressors

| | |
|---|---|
| linear_model.LinearRegression(*[, ...]) | Ordinary least squares Linear Regression. |
| linear_model.Ridge([alpha, fit_intercept, ...]) | Linear least squares with l2 regularization. |
| linear_model.RidgeCV([alphas, ...]) | Ridge regression with built-in cross-validation. |
| linear_model.SGDRegressor([loss, penalty, ...]) | Linear model fitted by minimizing a regularized empirical loss with SGD |

# Example

```
>>> from sklearn.linear_model import LinearRegression
```

```
>>> from sklearn.linear_model import Ridge
>>> import numpy as np


>>> n_samples, n_features = 10, 5
>>> rng = np.random.RandomState(0)
>>> y = rng.randn(n_samples)
>>> X = rng.randn(n_samples, n_features)
```

```
>>> rr = Ridge(alpha=0.1)
>>> rr.fit(X, y)
>>> pred_train_rr= rr.predict(X)
```
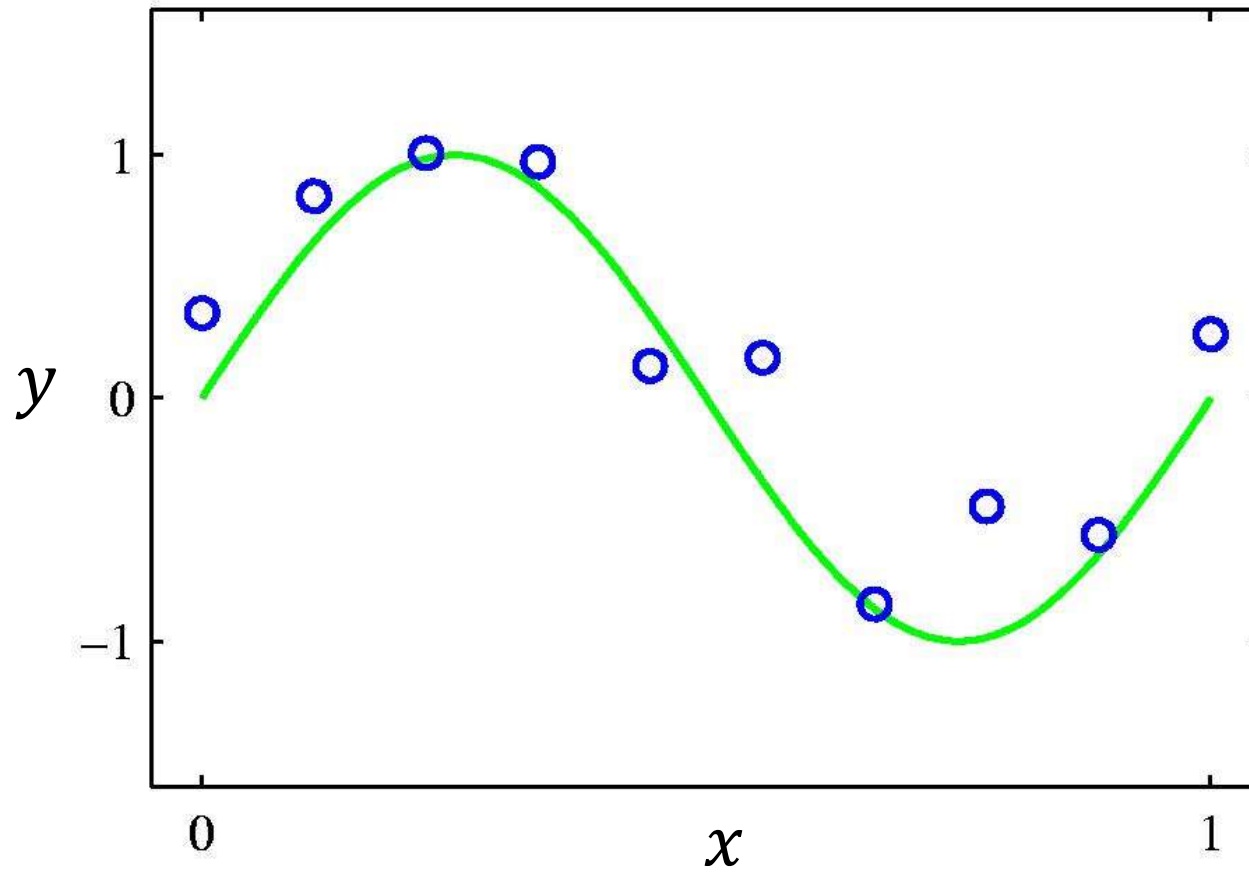
```
>>> lr = LinearRegression()
>>> lr.fit(X, y)
>>> pred_train_lr= lr.predict(X)
```

Model training and testing

# Nonlinear Regression
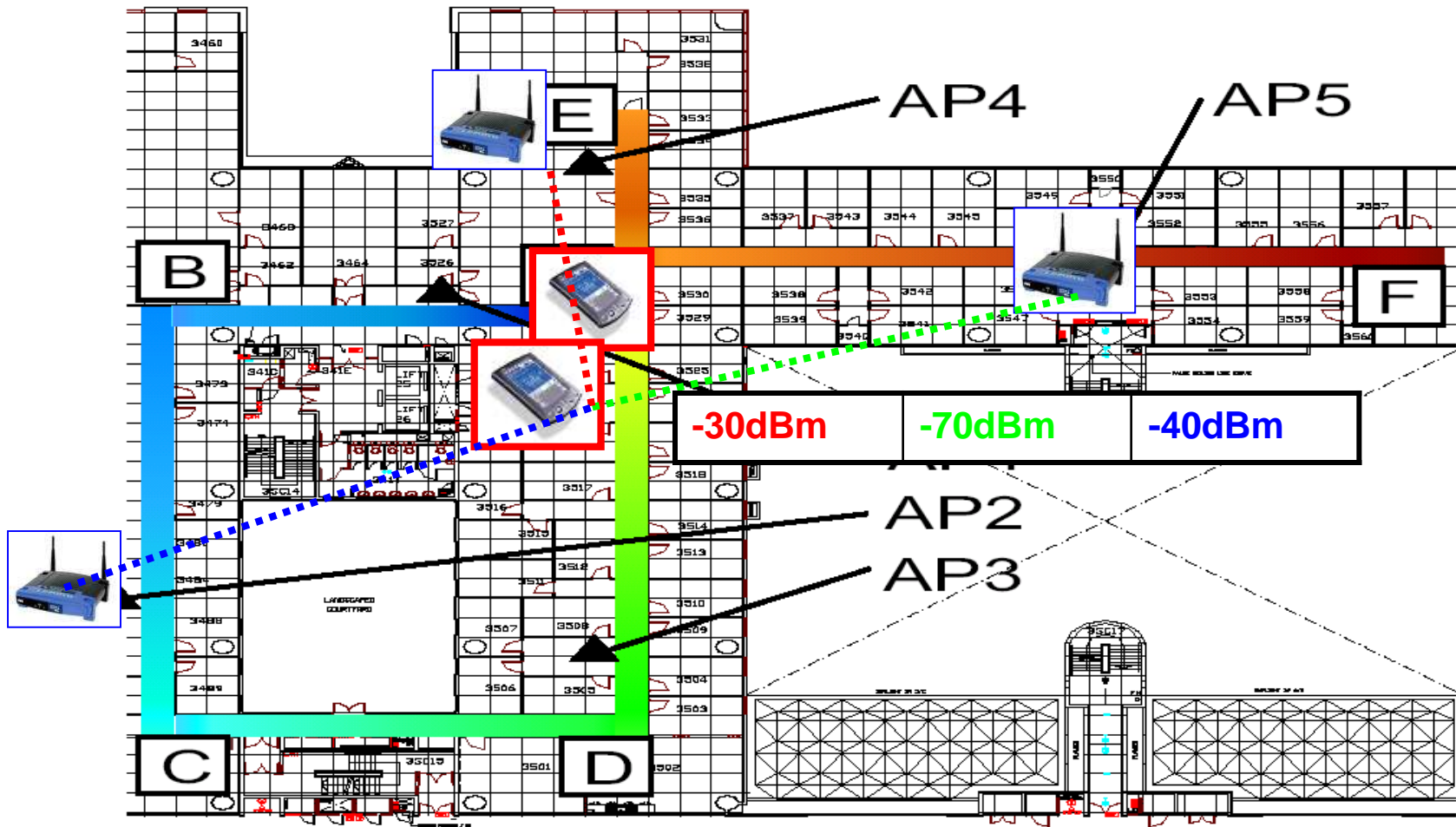
# Regression: Real-world Example

Indoor WiFi localization

# Thank you!