

AI 6102: Machine Learning Methodologies & Applications

L3: Linear Models: Regression

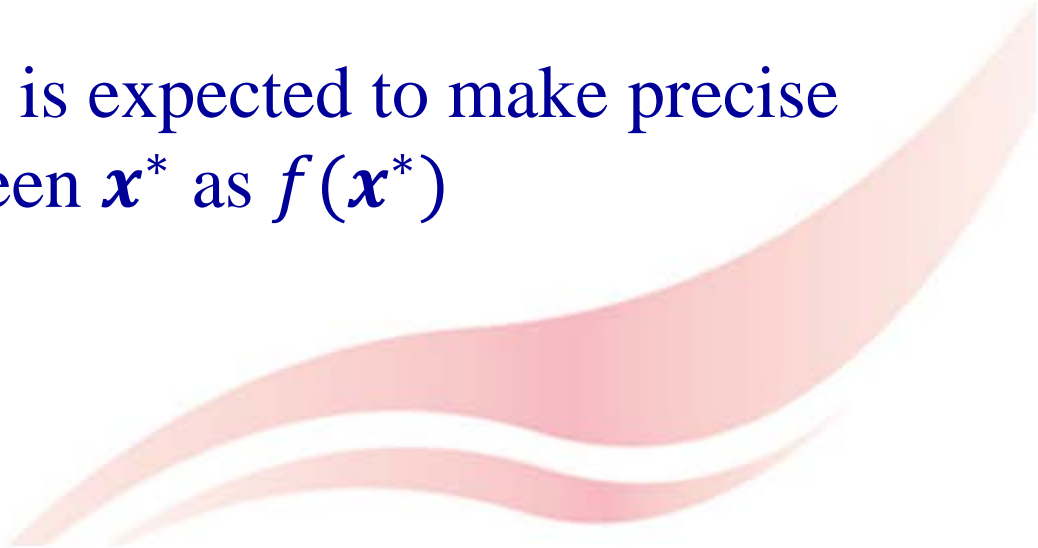
Sinno Jialin Pan

Nanyang Technological University, Singapore


Homepage: <http://www.ntu.edu.sg/home/sinnopan>

Recall: Supervised Learning

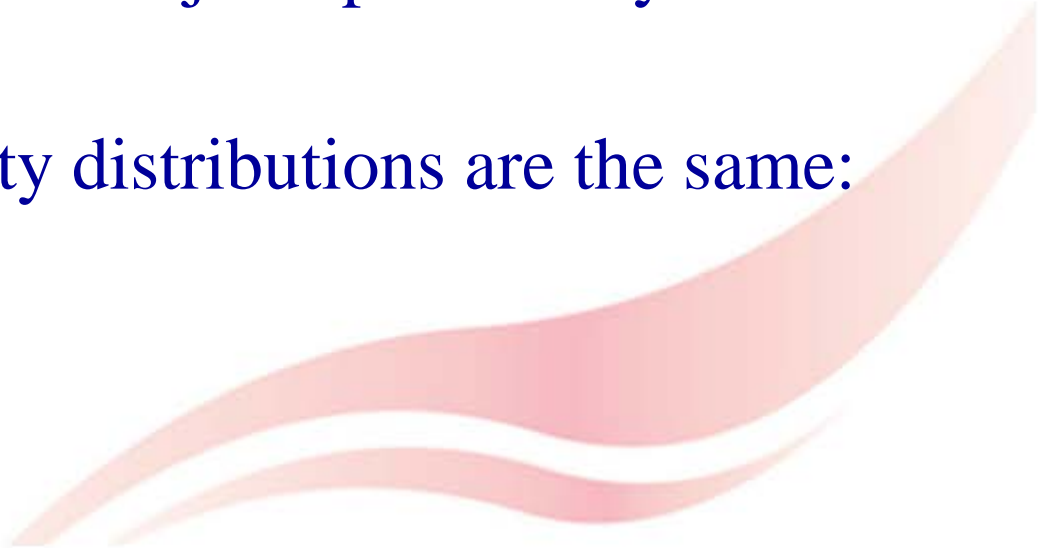
In mathematics

- Given: a set of N labeled data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, where \mathbf{x}_i is m -dimensional vector of numerical values, and y_i is a scalar
 - We aim to learn a mapping $f: \mathbf{x} \rightarrow y$ by requiring $f(\mathbf{x}_i) = y_i$
 - The learned mapping f is expected to make precise predictions on any unseen \mathbf{x}^* as $f(\mathbf{x}^*)$
- 

Hypothesis

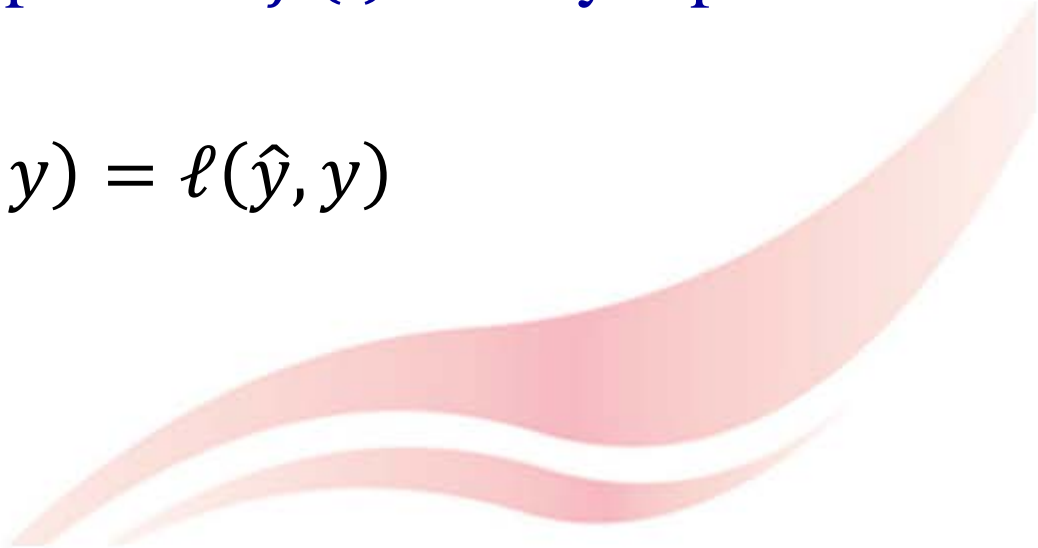
- A mapping or function $f: \mathbf{x} \rightarrow y$ can be considered as an element of some space of possible functions $\mathcal{H}: \mathbb{R}^m \rightarrow \mathbb{R}$, often called hypothesis space
 - Supervised learning aims to find a hypothesis $f \in \mathcal{H}$ from training data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, s.t. for any test \mathbf{x}^* , $f(\mathbf{x}^*) = y^*$
- 

Assumption

- The training data instances $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, are independent and identically distributed (i.i.d.), and drawn from an unknown joint probability distribution $P_{tr}(\mathbf{x}, y)$
 - Unseen test data instances $\{(\mathbf{x}^*, y^*)\}$ are also i.i.d, and drawn from an unknown joint probability distribution $P_{ts}(\mathbf{x}, y)$
 - The two joint probability distributions are the same:
 $P_{tr}(\mathbf{x}, y) = P_{ts}(\mathbf{x}, y)$
- 

Loss Function

- Denote by $\hat{y} = f(\mathbf{x})$ the prediction of the function $f(\cdot)$ on a data instance \mathbf{x} , and y is the ground-truth output of \mathbf{x}
- Let $\ell: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+ \geq 0$ 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 \mathbf{x}

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


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 $P(\mathbf{x}, y)$:

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

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

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



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

Risk Minimization (cont.)

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

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

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


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

Definition of expectation

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

Risk Minimization (cont.)

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

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


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

$= 1$ Assumption
 $P_{tr}(\mathbf{x}, y) = P_{ts}(\mathbf{x}, y)$

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

Definition of expectation

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

Empirical Risk Minimization

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

- The distribution $P_{tr}(\mathbf{x}, y)$ is unknown, thus we are not able to sample (infinite) input-output pairs $\{(\mathbf{x}, y)\}$ to learn a hypothesis $f(\cdot)$!
- In practice, only a finite number of training pairs are available, $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- Approximate the expected risk by empirical risk:

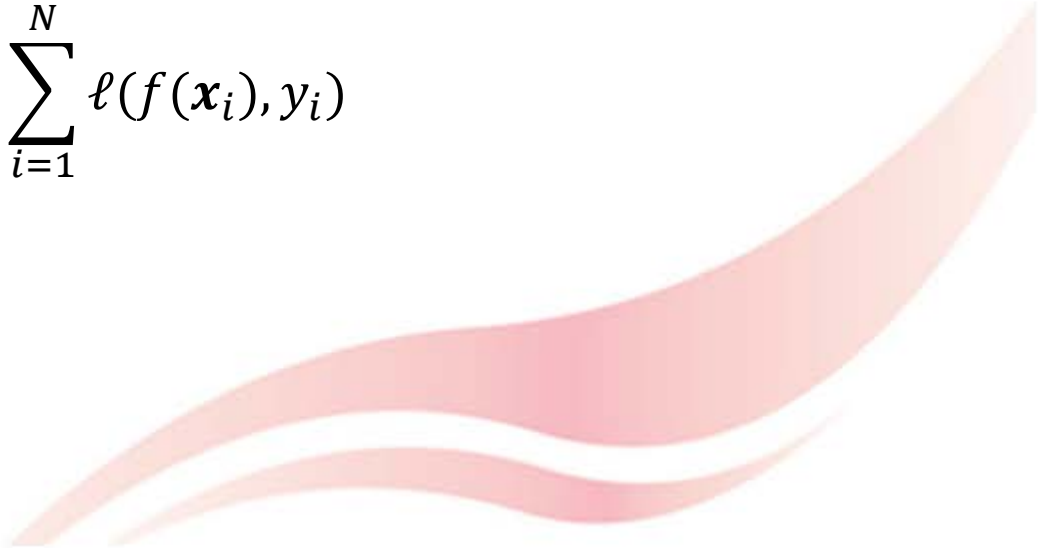
$$\mathbb{E}_{(\mathbf{x}, y) \sim P_{tr}} [\ell(f(\mathbf{x}), y)] \approx \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i) = \hat{R}_{tr}(f)$$

Empirical Risk Minimization (cont.)

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

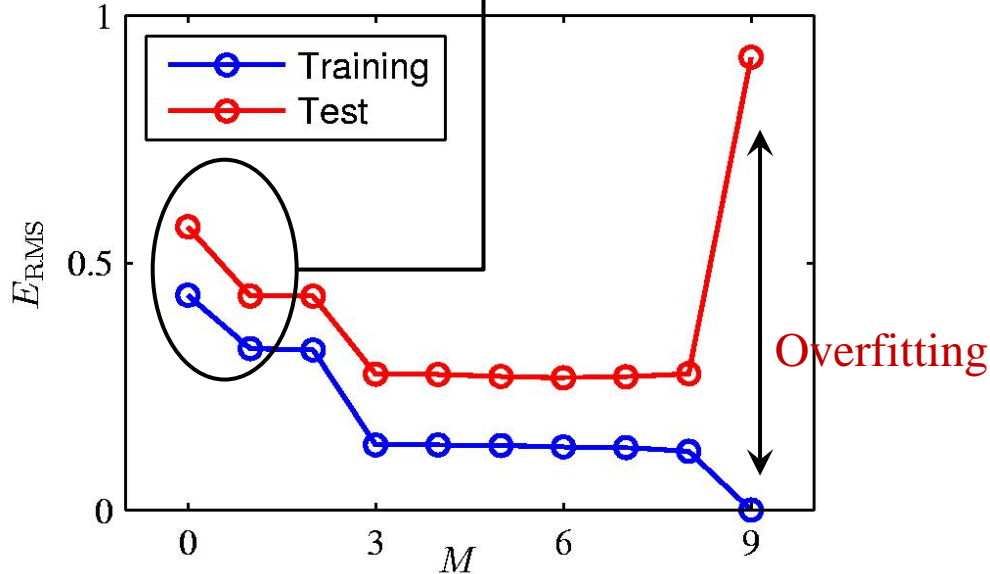
$$\hat{f} = \arg \min_f \hat{R}_{tr}(f) = \arg \min_f \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i)$$

- Given a training data set, N is a constant, thus for convenience in presentation, $\frac{1}{N}$ is dropped

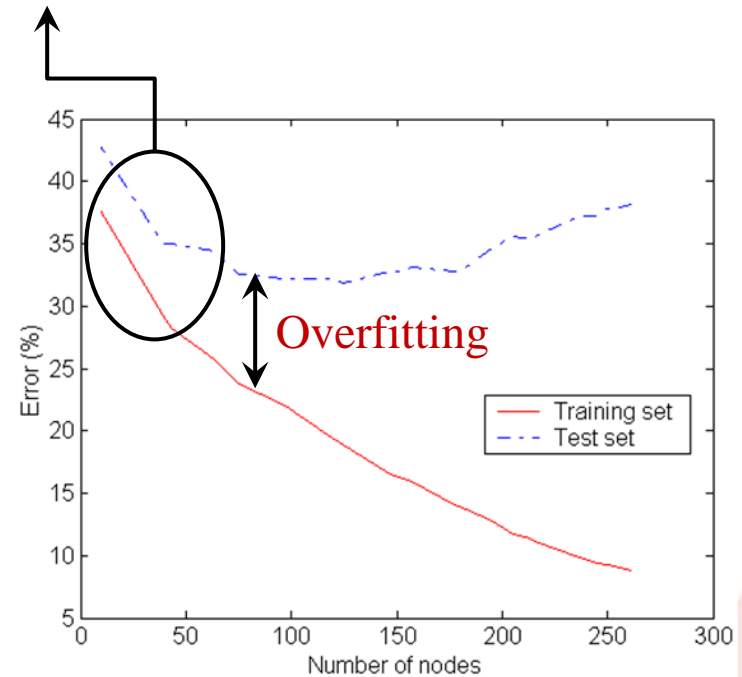
$$\hat{f} = \arg \min_f \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i)$$


Overfitting Revisit

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



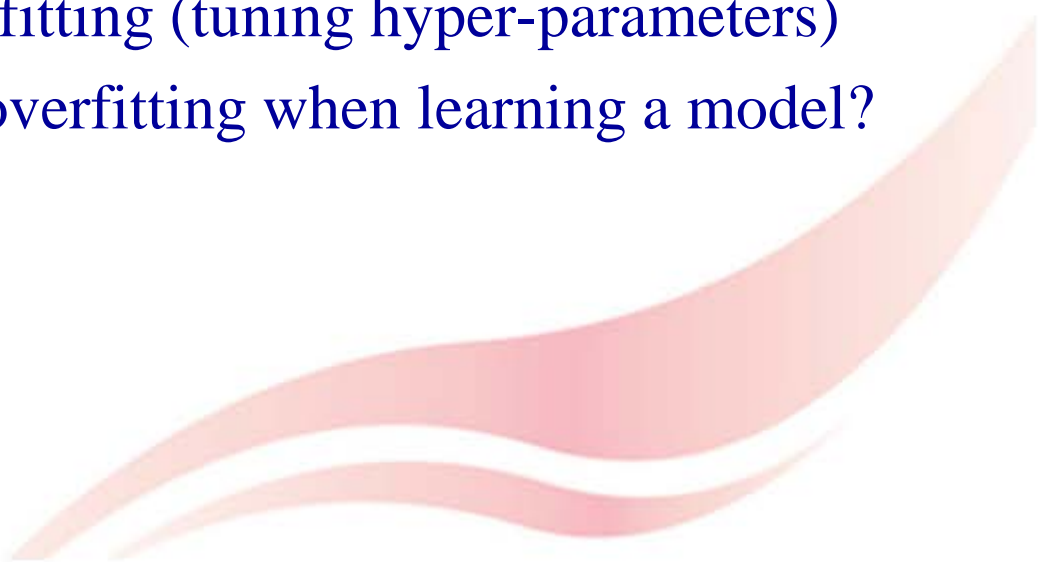
Polynomial curve fitting




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

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

- Structural Risk Minimization

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


- $\Omega(f)$ 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.)

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

- 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., $f(\mathbf{x}; \boldsymbol{\theta})$
- The parameterized $f(\mathbf{x}; \boldsymbol{\theta})$ defines a family of functions with different values of $\boldsymbol{\theta}$
- Learning f is equivalent to learning the values of $\boldsymbol{\theta}$

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

Structural Risk Minimization (cont.)

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

- 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 θ_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 θ_i , and fewer parameters with non-zero values
 - Induce sparsity, i.e., some θ_i 's tend to be zeros

Linear Models: Regression

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

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

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

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

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

$$\mathbf{x}_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \dots \\ x_{mi} \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_m \end{bmatrix}$$

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

Linear Models: Classification

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

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

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

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

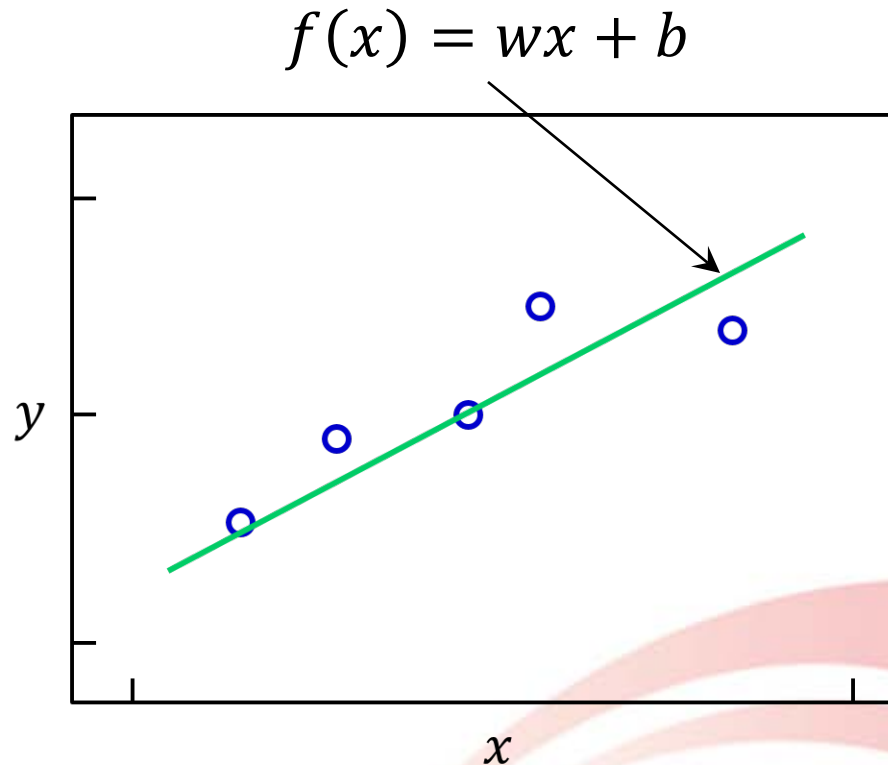
- For example,

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

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 $\{x_i, y_i\}, i = 1, \dots, N$



1D Linear Regression (cont.)

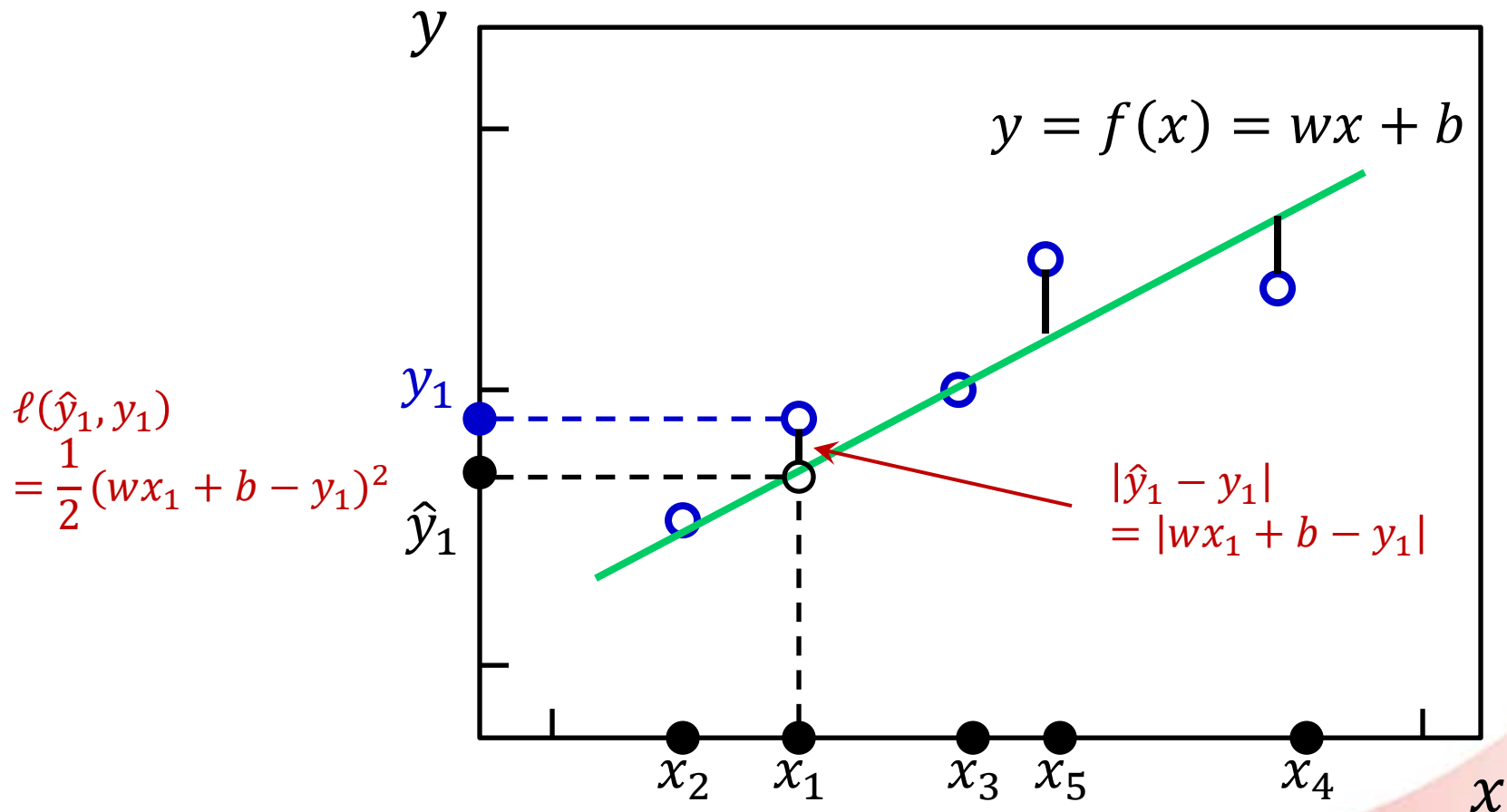
$$[\hat{w}, \hat{b}] = \arg \min_{[w, b]} \sum_{i=1}^N \ell(\hat{y}_i, y_i)$$

Drop the regularization term for simplicity at first

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



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

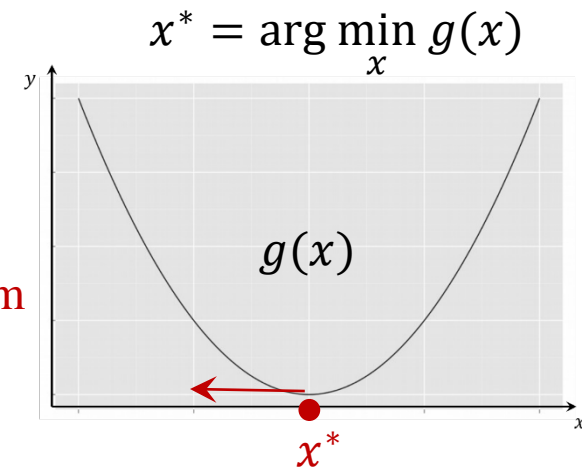
Optimization

- Learn w and b by minimizing the square loss

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

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

- 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

$$\frac{\partial \left(\frac{1}{2} \sum_{i=1}^N (wx_i + b - y_i)^2 \right)}{\partial w} = 0$$

$$\frac{\partial \left(\frac{1}{2} \sum_{i=1}^N (wx_i + b - y_i)^2 \right)}{\partial b} = 0$$

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

Chain rule of calculus

$$y = g(x)$$

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

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

$$z_i = wx_i + b - y_i$$

$$\begin{aligned} \frac{\partial z_i^2}{\partial w} &= \frac{\partial z_i^2}{\partial z_i} \frac{\partial z_i}{\partial w} \\ &= 2z_i \frac{\partial (wx_i + b - y_i)}{\partial w} \\ &= 2z_i x_i \end{aligned}$$

Closed-form Solution (cont.)

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

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

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

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

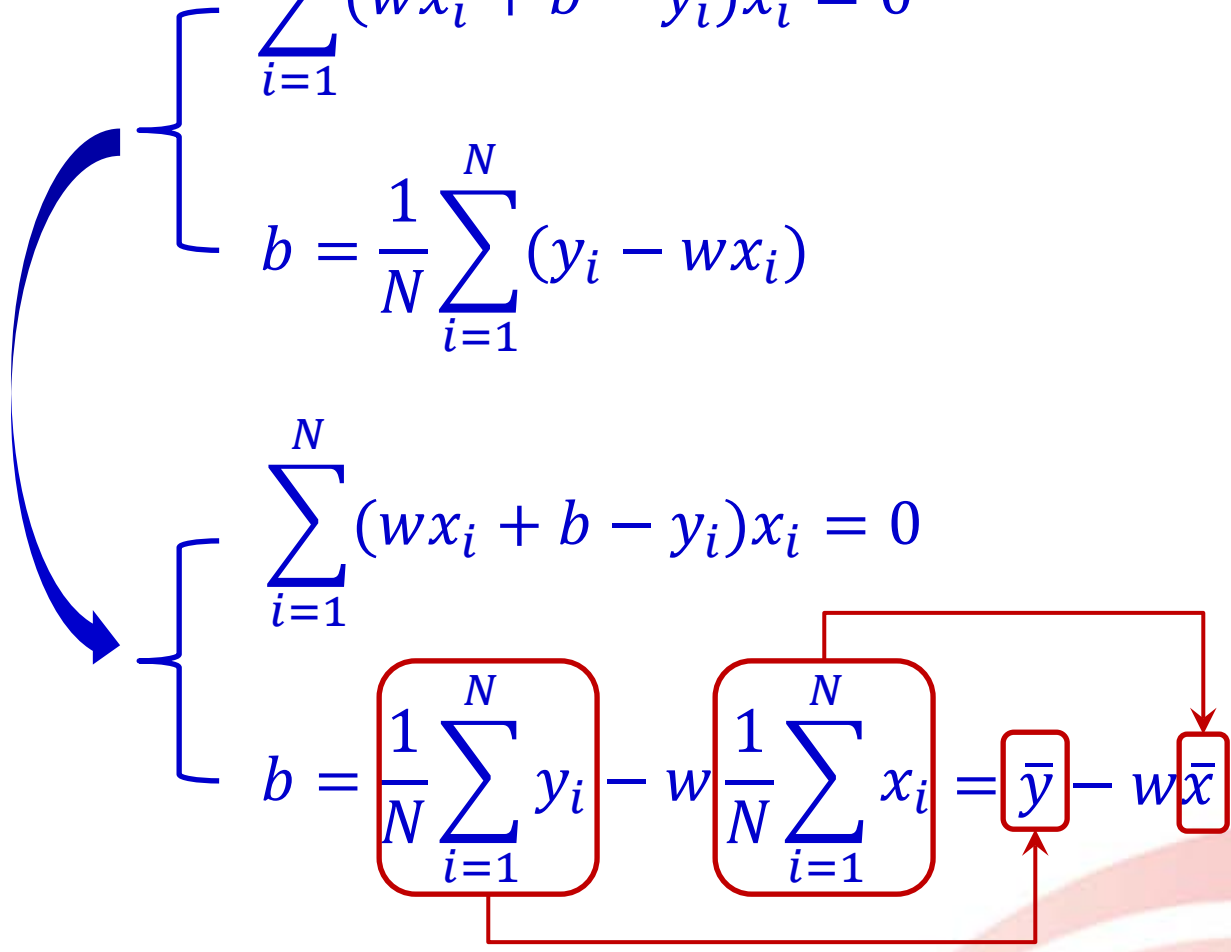
$$\sum_{i=1}^N (wx_i - y_i) + Nb = 0$$

Closed-form Solution (cont.)

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

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

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

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


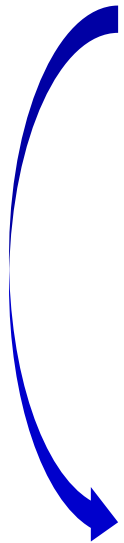
Closed-form Solution (cont.)

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

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

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

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



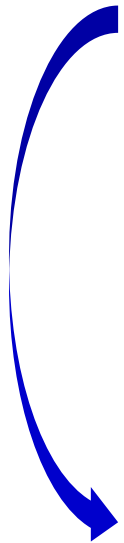
Closed-form Solution (cont.)

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

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

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

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



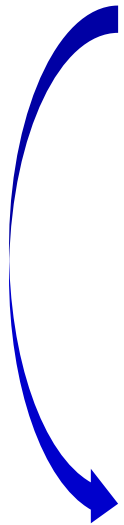
Closed-form Solution (cont.)

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

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

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

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



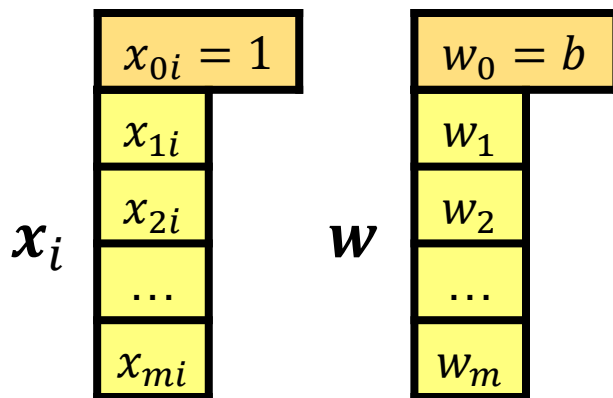
Multi-Dimension Case

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

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

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

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



Both \mathbf{w} and \mathbf{x} have
 $m + 1$ dimensions

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

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

Optimization

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

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

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

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

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

$$\mathbf{x} (m \times 1)$$

x_1
x_2
\dots
x_m

$$\mathbf{x}^T (1 \times m)$$

x_1	x_2	\dots	x_m
-------	-------	---------	-------

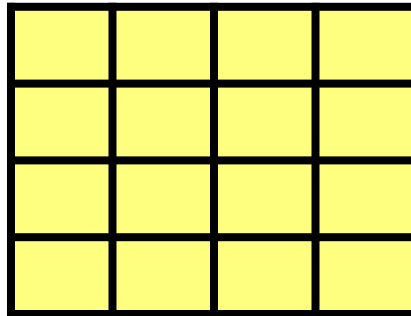
$$\mathbf{A} (m \times N)$$

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

Matrix/Vector Concepts

- Square matrix
 - If a matrix \mathbf{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 \mathbf{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{Ax}

$$(\mathbf{Ax})^T = \mathbf{x}^T \mathbf{A}^T$$

- The transpose of $\mathbf{x}^T \mathbf{y}$

$$(\mathbf{x}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{x}$$


- 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)
 - $(\mathbf{A}^{-1})^{-1} = \mathbf{A}$
 - $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$
 - $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$
- 

Linear System

- Given the following system of linear equations

$$\begin{aligned}2x_1 + x_3 &= 5 \\3x_1 - 4x_2 + 2x_3 &= 4 \\2x_2 - 3x_3 &= -3 \\-x_1 + 2x_2 - 5x_3 &= 1\end{aligned}$$

- They can be written in a more compact form as

$$\mathbf{Ax} = \mathbf{b}$$

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

2	0	1
3	-4	2
0	2	-3
-1	2	-5

$\mathbf{x} (3 \times 1)$

x_1
x_2
x_3

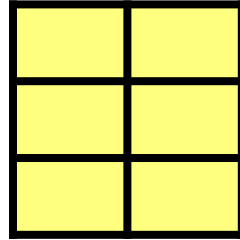
$\mathbf{b} (4 \times 1)$

5
4
-3
1

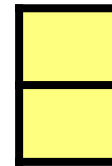
Linear System (cont.)

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

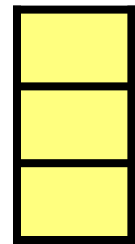
$$\mathbf{A} \ (k \times d)$$



$$\mathbf{x} \ (d \times 1)$$



$$\mathbf{b} \ (k \times 1)$$



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

$$\mathbf{A}^{-1}\mathbf{A}\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} \implies \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

Linear System (cont.)

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

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

$\mathbf{x} (d \times 1)$

$\mathbf{b} (k \times 1)$

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

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

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

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

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

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

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

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

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

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

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

The Matrix Cookbook

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Closed-form Solution (cont.)

$$a\mathbf{x} = \mathbf{x}a$$

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

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

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

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

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

Closed-form Solution (cont.)

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

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

$(m+1) \times 1$

 \mathbf{x}_i

x_{0i}
x_{1i}
...
x_{mi}

\mathbf{x}_i^T

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

$1 \times (m+1)$

$\mathbf{x}_i^T \mathbf{x}_i$ Inner product, scalar

$\mathbf{x}_i \mathbf{x}_i^T$ $(m+1)$ 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.)

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

$$x_i x_i^T$$

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

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

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

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

Closed-form Solution (cont.)

$$\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T)$$

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

$$\sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) = \mathbf{X} \mathbf{X}^T$$

$(m + 1)$ by N

\mathbf{X}

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

N by $(m + 1)$

\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.)

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

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

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

$(m+1) \text{ by } N$

\mathbf{X}

x_{01}	x_{02}	\dots	x_{0N}
x_{11}	x_{12}	\dots	x_{1N}
\dots	\dots	\dots	\dots
x_{m1}	x_{m2}	\dots	x_{mN}


$N \text{ by } 1$

\mathbf{y}

y_1
y_2
\dots
y_N

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


Closed-form Solution (cont.)


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

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



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


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


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

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



When \mathbf{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

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


Regularized Linear Regression

- Recall structural risk minimization

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

$\lambda > 0$ tradeoff
hyper-parameter

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

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

A regularization term to control
the complexity of the model

Also known as Ridge regression

Optimization

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

- Still an unconstrained optimization problem, and the objective is still convex
- A closed-form solution can be obtained by setting the derivative of the objective w.r.t. \mathbf{w} , and solving the resultant equation

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

Closed-form Solution

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

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

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

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

The matrix cookbook

Closed-form Solution (cont.)

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

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

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

Always invertible as
long as λ is positive

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

Why $\mathbf{XX}^T + \lambda \mathbf{I}$ Invertible?

- A square matrix is invertible if and only if it does not have a zero eigenvalue
- If a symmetric matrix \mathbf{A} is positive semidefinite, then all of its eigenvalues are non-negative (≥ 0)
 - When a symmetric matrix \mathbf{A} ($d \times d$) is said to be positive semidefinite iif for any non-zero column vector \mathbf{x} ($d \times 1$), $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$
- If a symmetric matrix \mathbf{A} is positive definite, then all of its eigenvalues are positive (> 0)
 - When a symmetric matrix \mathbf{A} ($d \times d$) is said to be positive definite iif for any non-zero column vector \mathbf{x} ($d \times 1$), $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$
- 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 $\mathbf{z}^T \mathbf{XX}^T \mathbf{z} = \boxed{\mathbf{y}^T \mathbf{y}} \quad \|\mathbf{y}\|_2^2 \geq 0$
- $\mathbf{XX}^T + \lambda \mathbf{I}$ 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: $\mathbf{z}^T (\mathbf{XX}^T + \lambda \mathbf{I}) \mathbf{z} = \boxed{\mathbf{z}^T \mathbf{XX}^T \mathbf{z}} + \boxed{\lambda \mathbf{z}^T \mathbf{z}} > 0$

≥ 0

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

Large-scale Issue

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

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

- The computation complexity of computing an inverse of a $(m + 1) \times (m + 1)$ matrix is $O((m + 1)^3)$
- When m is large, it is time consuming
- Rather than computing the inverse to obtain a closed-form solution, consider the following linear system

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

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

Large-scale Issue (cont.)

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

- 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

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

The diagram illustrates the gradient descent update rule. The equation $\mathbf{w}_{t+1} = \mathbf{w}_t - \rho \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}}$ is shown. A red arrow points from the text "Objective to be minimized" to the term $E(\mathbf{w})$ in the numerator of the gradient. Another red arrow points from the text "Learning rate $\rho \in (0,1]$ " to the variable ρ in the denominator. The terms $E(\mathbf{w})$ and ρ are enclosed in red boxes.

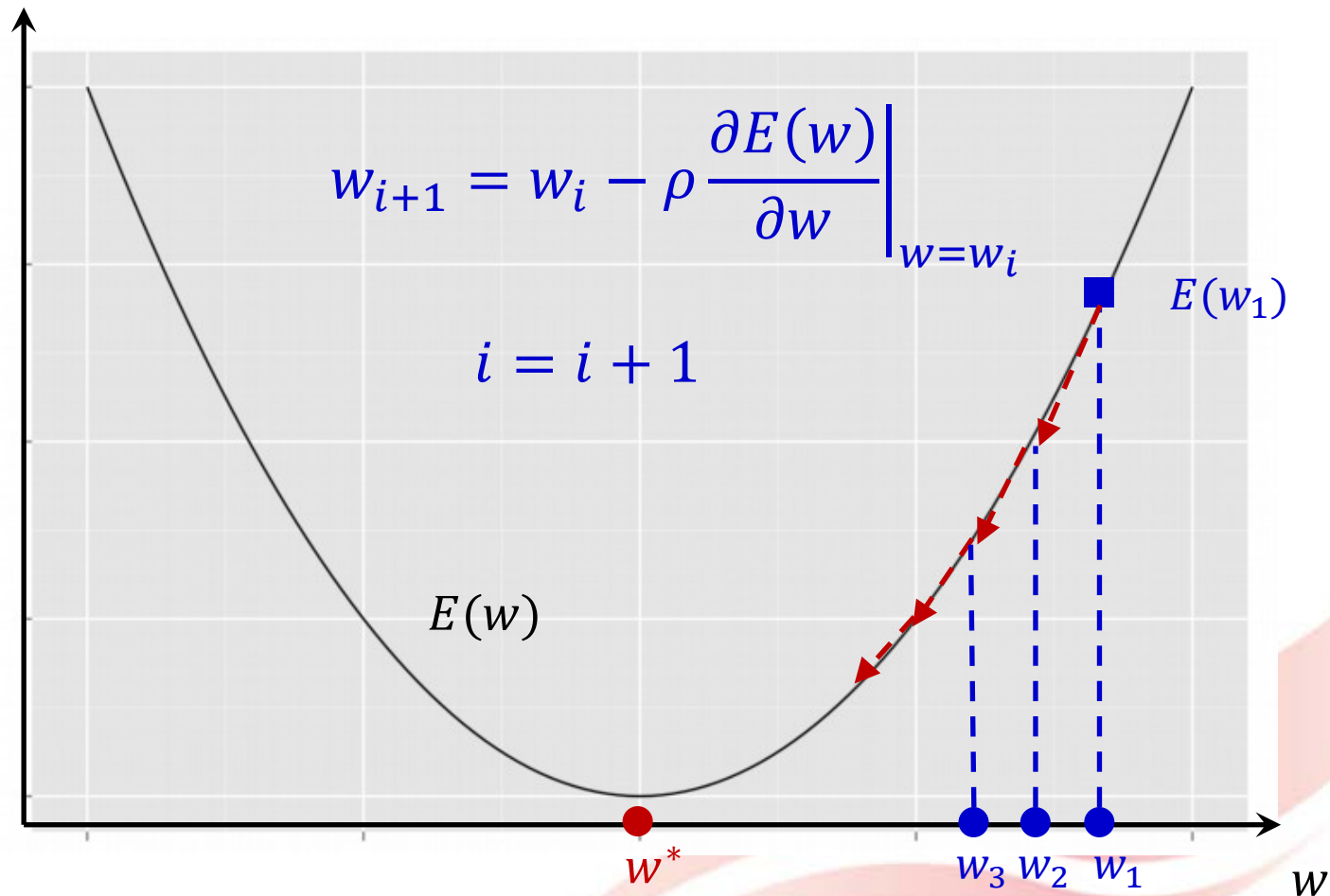
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \rho \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}}$$

Objective to be minimized

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

Gradient Descent

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



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

<code>linear_model.LinearRegression(*[, ...])</code>	Ordinary least squares Linear Regression.
<code>linear_model.Ridge([alpha, fit_intercept, ...])</code>	Linear least squares with l2 regularization.
<code>linear_model.RidgeCV([alphas, ...])</code>	Ridge regression with built-in cross-validation.
<code>linear_model.SGDRegressor([loss, penalty, ...])</code>	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()
```

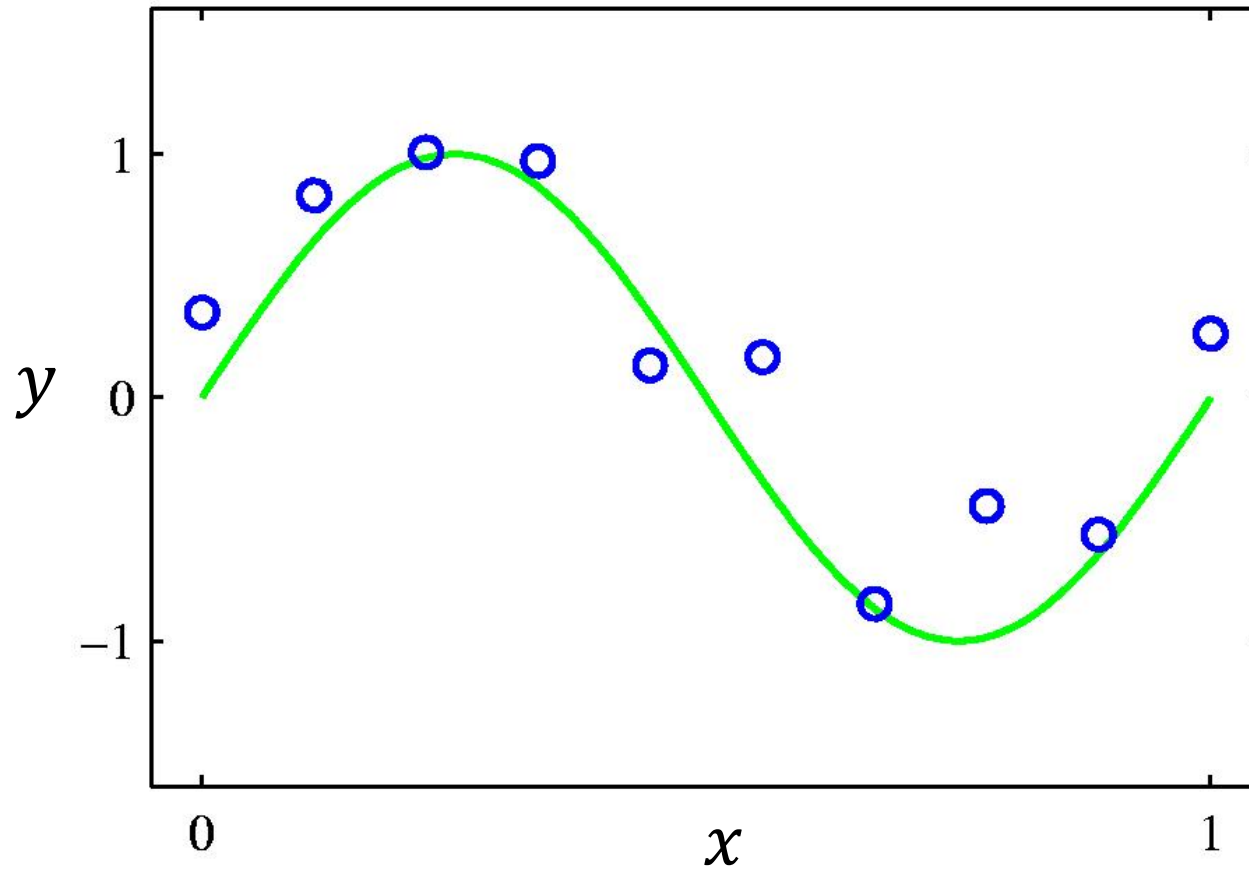
```
>>> lr.fit(X, y)
```

```
>>> pred_train_lr= lr.predict(X)
```

Model training and testing

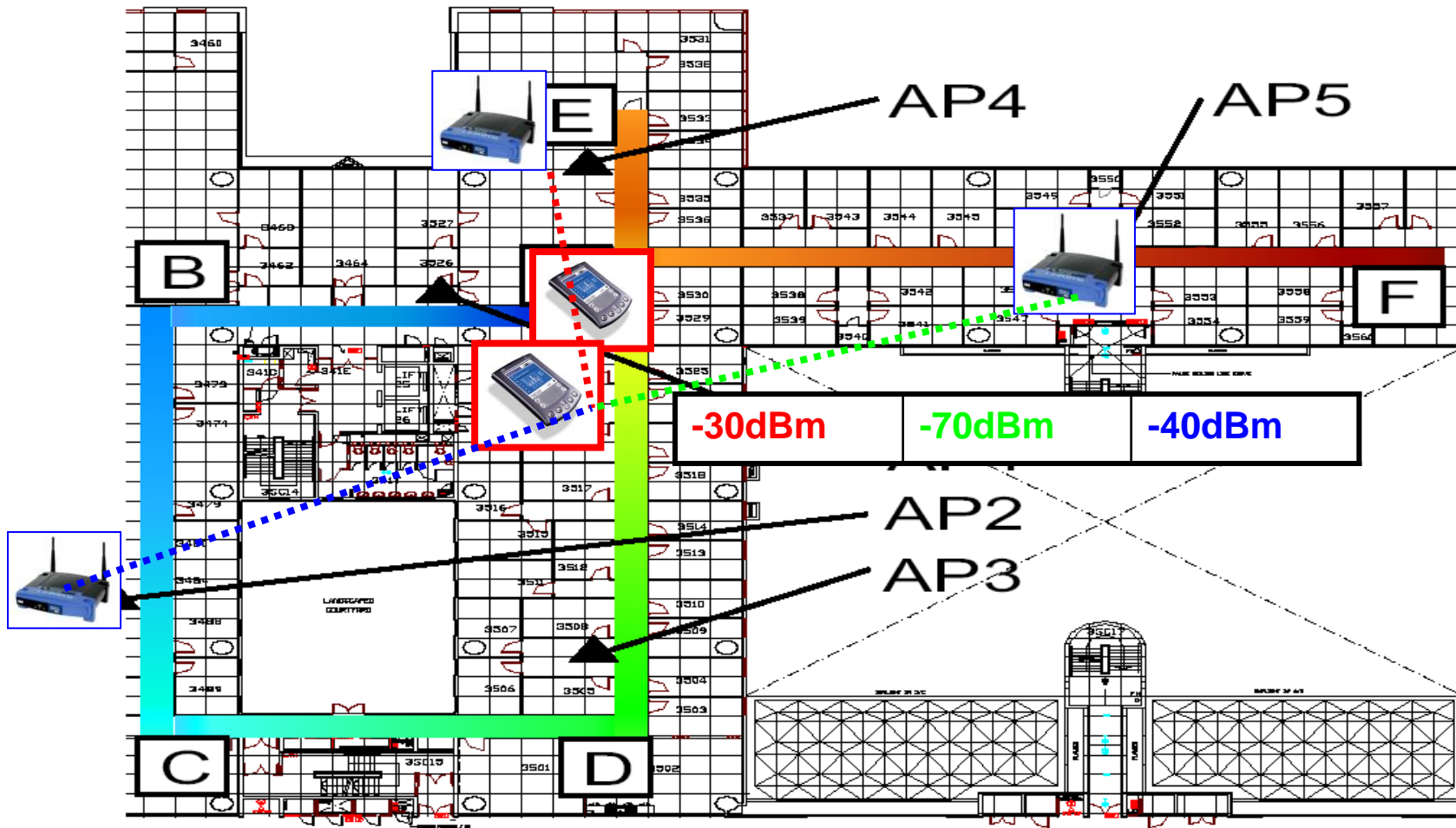
Nonlinear Regression

Kernel methods (L5)



Regression: Real-world Example

Indoor WiFi localization



Thank you!

