

EECS 545: Machine Learning

Lecture 2. Linear Regression (Part 1)

Honglak Lee
1/13/2025



Announcement

- Homework #1 will be out tomorrow (Jan. 14) and will be due 11:55 pm, Jan. 28 (Tue)
 - Note: this is the same date as Add/Drop deadline.
 - Form a study group and start early.
- Honor code
 - Collaboration and discussion is strongly encouraged, but you should write your own solution independently.
 - Write down the names of study group members.
 - **Do not** refer to or copy solutions from any other people or other resources. In addition, please do not let other people copy your solution.

Announcement

- [Project information](#) and [suggested project topics](#) will be released by today (to be updated by Friday).
- The project proposal is due by Feb 4, Tuesday (23:55 PM).

3

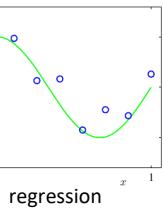
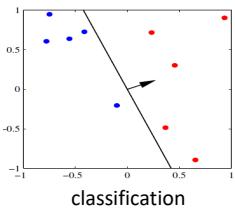
Announcement

- Schedule of review sessions (with Zoom links) announced on Canvas. Recordings will be made available after the session
 - Linear Algebra (by Yiwei) – 1/13, 4pm
 - Probability (by Ishan) – 1/17, 11am
 - Python / NumPy (by Violet) – 1/14, 4pm
- A quiz will be due 24 hours after every lecture
 - E.g. the lecture 2 quiz will be due tomorrow (Tuesday 1/14) at 10:30am
- Questions?

4

Supervised Learning

- Goal:
 - Given data X in feature space and the labels Y
 - Learn to predict Y from X
- Labels could be discrete or continuous
 - Discrete-valued labels: classification
 - Continuous-valued labels: regression (today's topic)



5

Overview of Topics

- Linear Regression
 - Objective function
 - Vectorization
 - Computing gradient
 - Batch gradient vs. Stochastic Gradient
 - Closed form solution

6

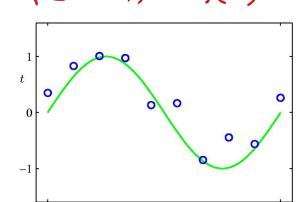
Notation

In this lecture, we will use the following notation:

- $\mathbf{x} \in \mathbb{R}^D$: data (scalar or vector)
- $\phi(\mathbf{x}) \in \mathbb{R}^M$: features for \mathbf{x} (vector)
- $\phi_j(\mathbf{x}) \in \mathbb{R}$: j-th feature for \mathbf{x} (scalar)
- $y \in \mathbb{R}$: continuous-valued label (i.e., target value)
- $\mathbf{x}^{(n)}$: denotes the n-th training example.
- $y^{(n)}$: denotes the n-th training label.

Linear regression (with 1D inputs)

- Consider the 1D case (e.g. $D=1$)
- Given a set of observation $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
- and corresponding target values $\{y^{(1)}, \dots, y^{(N)}\}$
- We want to learn a function $h(\mathbf{x}, \mathbf{w}) \approx y$ to predict future values



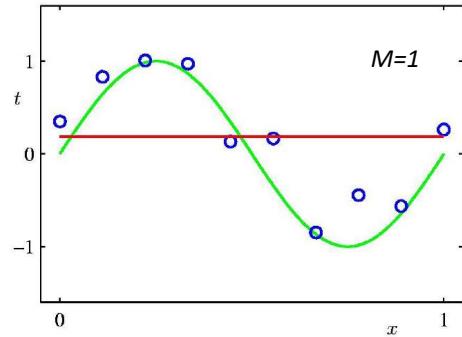
M factors $(1, x_1, x_2, \dots, x^{M-1})$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_{M-1} x^{M-1} = \sum_{j=0}^{M-1} w_j x^j$$

7

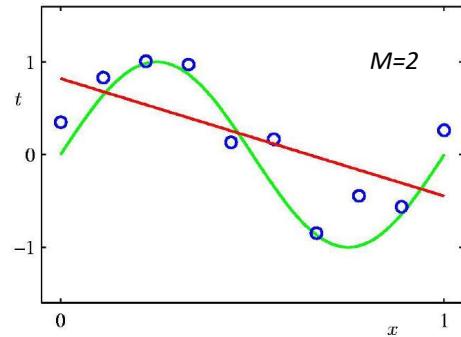
8

0th Order Polynomial



9

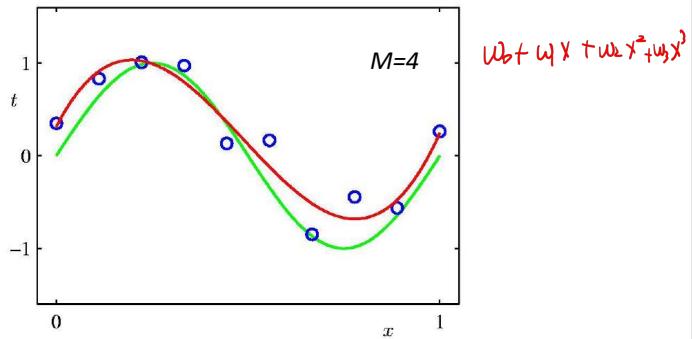
1st Order Polynomial



$$w_0 + w_1 \cdot x$$

10

3rd Order Polynomial



$$w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

11

Linear Regression (general case)

$$h(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$

$$\phi_0(\mathbf{x}) = 1$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{bmatrix}$$

$$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^\top$$

(\mathbf{w} and $\phi(\mathbf{x})$ are column vectors)

- The function $h(\mathbf{x}, \mathbf{w})$ is linear in parameters \mathbf{w} .

– Goal: Find the best value for the weights \mathbf{w} .

- For simplicity, add a *bias term (constant function)*:

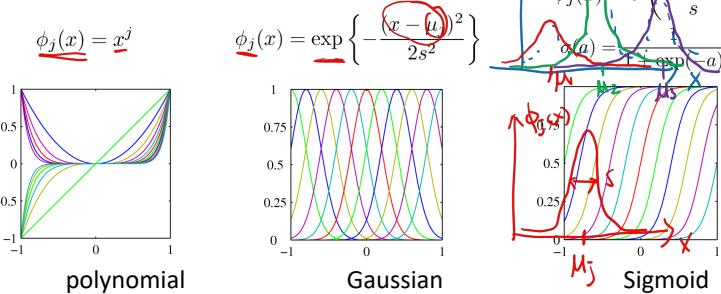
$$h(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

$$\text{where } \mathbf{w} = (w_0, \dots, w_{M-1})^\top$$

$$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^\top$$

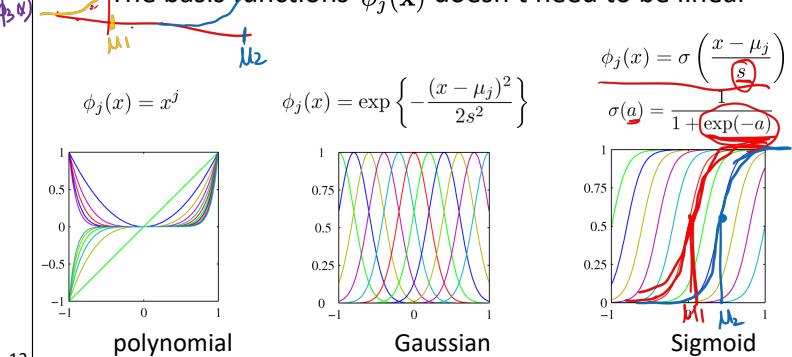
Basis Functions

- The basis functions $\phi_j(\mathbf{x})$ doesn't need to be linear

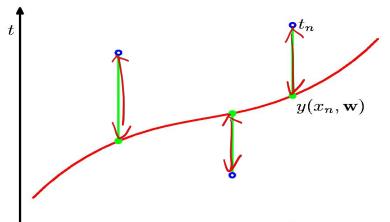


Basis Functions

- The basis functions $\phi_j(\mathbf{x})$ doesn't need to be linear



Objective: Sum-of-Squares Error Function

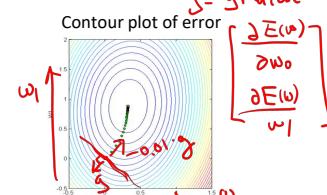
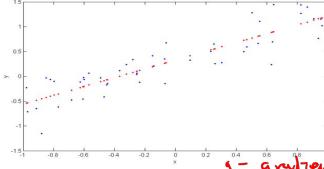


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(h(\mathbf{x}^{(n)}, \mathbf{w}) - y^{(n)} \right)^2$$

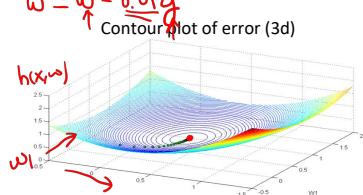
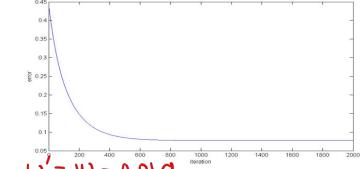
We want to find \mathbf{w} that minimizes $E(\mathbf{w})$ over the training data.

Linear regression via gradient descent (illustration)

Training data (blue) vs. prediction (red)



Error curve vs. training epoch



15

Least squares problem

- Objective function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2$$

- Gradient

$$\begin{aligned} \frac{\partial E(w)}{\partial w_k} &= \frac{\partial}{\partial w_k} \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ &= \frac{1}{2} \cdot 2 \cdot \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \frac{\partial}{\partial w_k} \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \\ &\quad \text{---} \\ &\quad \text{---} \\ &\quad \text{---} \end{aligned}$$

17

Least squares problem

- Objective function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2$$

- Gradient

$$\begin{aligned} \frac{\partial E(w)}{\partial w_k} &= \frac{\partial}{\partial w_k} \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ &= \sum_{n=1}^N \left[\left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \frac{\partial}{\partial w_k} \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \right] \end{aligned}$$

18

Least squares problem

- Objective function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2$$

- Gradient

$$\begin{aligned} \frac{\partial E(w)}{\partial w_k} &= \frac{\partial}{\partial w_k} \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ &= \sum_{n=1}^N \left[\left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \frac{\partial}{\partial w_k} \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \right] \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_k(\mathbf{x}^{(n)}) \end{aligned}$$

Concatenate each component of the gradient:

$$\frac{\partial E(w)}{\partial w_k} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_k(\mathbf{x}^{(n)})$$

We get a vectorized form of the gradient:

$$\begin{aligned} \nabla_w E(\mathbf{w}) &= \left[\begin{array}{c} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{array} \right] \\ g &= \left[\begin{array}{c} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_0(\mathbf{x}^{(n)}) \\ \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_{M-1}(\mathbf{x}^{(n)}) \end{array} \right] \end{aligned}$$

Same

19

20

Concatenate each component of the gradient:

$$\frac{\partial E(w)}{\partial w_k} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_k(\mathbf{x}^{(n)})$$

We get a vectorized form of the gradient:

$$\begin{aligned} \nabla_w E(\mathbf{w}) &= \left[\begin{array}{c} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{array} \right] \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix} \in \mathbb{R}^M \end{aligned}$$

Concatenate each component of the gradient:

$$\frac{\partial E(w)}{\partial w_k} = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi_k(\mathbf{x}^{(n)})$$

We get a vectorized form of the gradient:

$$\begin{aligned} \nabla_w E(\mathbf{w}) &= \left[\begin{array}{c} \frac{\partial}{\partial w_0} E(\mathbf{w}) \\ \frac{\partial}{\partial w_1} E(\mathbf{w}) \\ \vdots \\ \frac{\partial}{\partial w_{M-1}} E(\mathbf{w}) \end{array} \right] \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \begin{bmatrix} \phi_0(\mathbf{x}^{(n)}) \\ \phi_1(\mathbf{x}^{(n)}) \\ \vdots \\ \phi_{M-1}(\mathbf{x}^{(n)}) \end{bmatrix} \\ &= \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)}) \\ &= \sum_{n=1}^N (\mathbf{w}^\top \phi(\mathbf{x}^{(n)}) - y^{(n)}) \phi(\mathbf{x}^{(n)}) \end{aligned}$$

21

23

Batch Gradient Descent

- Given data (\mathbf{x}, \mathbf{y}) and an initial \mathbf{w}

- Repeat until convergence:

$$\eta = 0.01$$

$$\mathbf{w} := \mathbf{w} - \eta \nabla_w E(\mathbf{w})$$

where

$$\begin{aligned} g &= \nabla_w E(\mathbf{w}) = \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)}) \\ &= \sum_{n=1}^N (\mathbf{w}^\top \phi(\mathbf{x}^{(n)}) - y^{(n)}) \phi(\mathbf{x}^{(n)}) \end{aligned}$$

Stochastic Gradient Descent

- Main idea: instead of computing batch gradient (over entire training data), just compute gradient for individual example and update

- Repeat until convergence

- for $n=1, \dots, N$

$$\mathbf{w} := \mathbf{w} - \eta \nabla_w E(\mathbf{w} | \mathbf{x}^{(n)})$$

where

$$\begin{aligned} \nabla_w E(\mathbf{w} | \mathbf{x}^{(n)}) &= \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)}) \\ &= (\mathbf{w}^\top \phi(\mathbf{x}^{(n)}) - y^{(n)}) \phi(\mathbf{x}^{(n)}) \end{aligned}$$

$t=1 \rightarrow \eta_1$
 $t \approx T \rightarrow \eta_T$

24

25

Stochastic Gradient Descent

- Repeat until convergence:
 - for $n=1, \dots, N$ *random shuffle*

$$\mathbf{w} := \mathbf{w} - \eta \nabla_{\mathbf{w}} E(\mathbf{w} | \mathbf{x}^{(n)})$$

Note: Typically the learning rate is gradually decreased as training time (t) goes on:
e.g., $\eta_t \propto \frac{1}{t}$ or $\eta_t = \eta_1 \frac{1}{(1 + (t-1)/\tau)}$

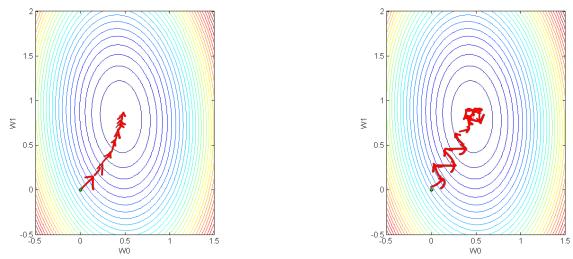
$$\text{where } \nabla_{\mathbf{w}} E(\mathbf{w} | \mathbf{x}^{(n)}) = \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right) \phi(\mathbf{x}^{(n)}) \\ = (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) - y^{(n)}) \phi(\mathbf{x}^{(n)})$$

- Implementation tips in practice:

— For each step of gradient computation in SGD, a small number of samples ("minibatch") may be used for computing the gradient instead of just one sample. Then we iterate this over the entire dataset with multiple epochs until convergence.

26

Batch gradient vs. Stochastic gradient



Closed form solution

- Main idea:
 - Compute gradient and set gradient to 0. (condition for optimal solution)
 - Solve the equation in a closed form,

- The objective function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2$$

- We will derive the gradient from matrix calculus

28

Closed form solution

- Objective function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) - y^{(n)})^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}))^2 - \sum_{n=1}^N y^{(n)} \mathbf{w}^T \phi(\mathbf{x}^{(n)}) + \frac{1}{2} \sum_{n=1}^N (y^{(n)})^2$$

$(a-b)^2 = a^2 - 2ab + b^2$

Closed form solution

- Objective function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) - y^{(n)})^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}))^2 - \sum_{n=1}^N y^{(n)} \mathbf{w}^T \phi(\mathbf{x}^{(n)}) + \frac{1}{2} \sum_{n=1}^N (y^{(n)})^2 \\ = \frac{1}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} - \mathbf{w}^T \Phi^T \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{y}$$

- Trick: vectorization (by defining data matrix)

$\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) = \mathbf{w}^T \Phi_{n,1}$ Φ is an $N \times M$ matrix, applying $\Phi \mathbf{w} = \mathbf{y}$

The data matrix $\begin{bmatrix} h(x^{(1)}, u) \\ \vdots \\ h(x^{(N)}, u) \end{bmatrix} = \begin{bmatrix} (\Phi \mathbf{w})_1 \\ \vdots \\ (\Phi \mathbf{w})_N \end{bmatrix}$

The design matrix is an $N \times M$ matrix, applying $\Phi \mathbf{w} \approx \mathbf{y}$

– the M basis functions (columns)
– to N data points (rows)

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}_{N \times M}$$

$$\Phi \mathbf{w} \approx \mathbf{y}$$

32

31

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix} \quad (\Phi \mathbf{w})^T = \mathbf{w}^T \Phi^T$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ = (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y}) \\ = (\mathbf{w}^T \Phi^T - \mathbf{y}^T)(\Phi \mathbf{w} - \mathbf{y}) \\ = \mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2 \mathbf{w}^T \Phi^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}$$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}) - y^{(n)})^2 \\ = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi(\mathbf{x}^{(n)}))^2 - \sum_{n=1}^N y^{(n)} \mathbf{w}^T \phi(\mathbf{x}^{(n)}) + \frac{1}{2} \sum_{n=1}^N (y^{(n)})^2$$

$(a-b)^2 = a^2 - 2ab + b^2$

34

36

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}^{(1)}) & \phi_1(\mathbf{x}^{(1)}) & \dots & \phi_{M-1}(\mathbf{x}^{(1)}) \\ \phi_0(\mathbf{x}^{(2)}) & \phi_1(\mathbf{x}^{(2)}) & \dots & \phi_{M-1}(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}^{(N)}) & \phi_1(\mathbf{x}^{(N)}) & \dots & \phi_{M-1}(\mathbf{x}^{(N)}) \end{pmatrix}$$

$\stackrel{\text{p.s.d.}}{=} \Phi^T \Phi$ $\Phi \in \mathbb{R}^{N \times M}$

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(\sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}^{(n)}) - y^{(n)} \right)^2$$

$$= \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^\top \Phi(\mathbf{x}^{(n)}) - y^{(n)})^2$$

$$\frac{1}{2} \sum_{n=1}^N (\Phi \mathbf{w})_n^2 = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^\top \Phi(\mathbf{x}^{(n)}))^2 - \sum_{n=1}^N y^{(n)} \mathbf{w}^\top \Phi(\mathbf{x}^{(n)}) + \frac{1}{2} \sum_{n=1}^N (y^{(n)})^2$$

$$(\Phi \mathbf{w})^\top \Phi \mathbf{w} = \frac{1}{2} \mathbf{w}^\top \Phi^\top \Phi \mathbf{w} - \mathbf{w}^\top \Phi^\top y + \frac{1}{2} y^\top y$$

$$= \mathbf{w}^\top \Phi^\top \Phi \mathbf{w}$$

37

Useful trick: Matrix Calculus

- Idea so far:
 - Compute gradient and set gradient to $\mathbf{0}$ (condition for optimal solution)
 - Solve the equation in a closed form using matrix calculus
- Need to compute the first derivative in matrix form

38

Matrix calculus: The Gradient

- Suppose that $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ is a function that takes as an input matrix A of size $[m \times n]$ and returns a real value (scalar).
- Then the gradient of f with respect to $A \in \mathbb{R}^{m \times n}$ is the matrix of partial derivatives, defined as:

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \dots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \dots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \dots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

$$(\nabla_A f(A))_{ij} = \frac{\partial f(A)}{\partial A_{ij}}$$

39

40

Matrix calculus: The Gradient

Note that the size of $\nabla_A f(A)$ is always the same as the size of A . So if, in particular, A is just a vector $x \in \mathbb{R}^n$, then

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

- $\nabla_x(f(x) + g(x)) = \nabla_x f(x) + \nabla_x g(x)$.
- For $t \in \mathbb{R}$, $\nabla_x(t f(x)) = t \nabla_x f(x)$.

Gradient of Linear Functions

- Linear function: $f(\mathbf{x}) = \sum_{i=1}^n b_i x_i = \mathbf{b}^\top \mathbf{x}$
- Gradient: $\frac{\partial f(\mathbf{x})}{\partial x_k} = \frac{\partial}{\partial x_k} \left(\sum_{i=1}^n b_i x_i \right) = b_k$
- Compact form: $\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{b}$

Gradient of Quadratic Functions

* Assumption: A is a symmetric matrix: i.e., $A_{ij} = A_{ji}$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial}{\partial \mathbf{x}} \left(\sum_{i,j=1}^n x_i A_{ij} x_j \right) = \mathbf{x}^\top A \mathbf{x}$$

$\sum_{i,j=1}^n x_i A_{ij} x_j = \sum_{i=1}^n x_i \left(\sum_{j=1}^n A_{ij} x_j \right) = \sum_{i=1}^n x_i (A \mathbf{x})_i$

$\frac{\partial f(\mathbf{x})}{\partial x_k} = \frac{\partial}{\partial x_k} \left(\sum_{i,j=1}^n x_i A_{ij} x_j \right) = \sum_{j=1}^n A_{kj} x_j = (A \mathbf{x})_k$

$\nabla_{\mathbf{x}} f(\mathbf{x}) = 2A\mathbf{x}$

41

42

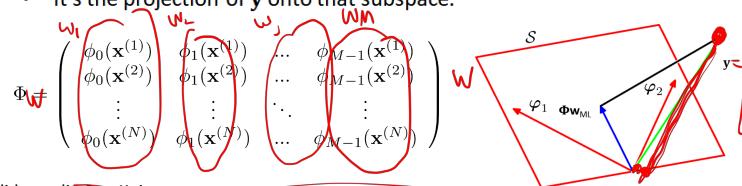
Putting together: Solution via matrix calculus

- Compute gradient and set to zero
- $\nabla_{\mathbf{w}} E(\mathbf{w}) = \nabla_{\mathbf{w}} \left(\frac{1}{2} \mathbf{w}^\top \Phi^\top \Phi \mathbf{w} - \mathbf{w}^\top \Phi^\top \mathbf{y} + \frac{1}{2} \mathbf{y}^\top \mathbf{y} \right)$
- $= \Phi^\top \Phi \mathbf{w} - \Phi^\top \mathbf{y}$
- $= \mathbf{0}$
- Solve the resulting equation (normal equation) $\Phi^\top \Phi \mathbf{w} = \Phi^\top \mathbf{y}$
- $\mathbf{w}_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y}$

This is the *Moore-Penrose pseudo-inverse*: $\Phi^\dagger = (\Phi^\top \Phi)^{-1} \Phi^\top$ applied to: $\Phi \mathbf{w} \approx \mathbf{y}$

Geometric Interpretation $\Phi \mathbf{w} \approx \mathbf{y}$

- Assuming many more observations (N) than the M basis functions $\phi_j(x)$ ($j=0, \dots, M-1$)
- View the observed target values $\mathbf{y} = \{y^{(1)}, \dots, y^{(N)}\}$ as a vector in an N -dim. space.
- The M basis functions $\phi_j(x)$ span the N -dimensional subspace.
 - Where the N -dim vector for ϕ_j is $\{\phi_j(\mathbf{x}^{(1)}), \dots, \phi_j(\mathbf{x}^{(N)})\}$
- $\Phi \mathbf{w}_{ML}$ is the point in the subspace with minimal squared error from \mathbf{y} .
- It's the projection of \mathbf{y} onto that subspace.



43

44

Slide credit: Ben Kuipers