

Introduction to Machine Learning - CS725  
Instructor: Prof. Ganesh Ramakrishnan  
Lecture 3 - Linear Regression - Probabilistic  
Interpretation and Regularization

# Recap: Linear Regression is **not** Naively Linear

Linear in  $w$ 's

Need to determine  $\mathbf{w}$  for the linear function  $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}_j) = \phi^T \mathbf{w}$  which minimizes our error function  $E(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$

Extremely

non-linear

in  $\mathbf{x}$ 's

Owing to basis function  $\phi$ , "Linear Regression" is *linear* in  $\mathbf{w}$  but NOT in  $\mathbf{x}$  (which could be arbitrarily non-linear)!

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_n(\mathbf{x}_1) \\ \vdots & \vdots & & \vdots \\ \phi_1(\mathbf{x}_m) & \phi_2(\mathbf{x}_m) & \dots & \phi_n(\mathbf{x}_m) \end{bmatrix} \quad (1)$$

# Recap: Linear Regression is **not** Naively Linear

- Need to determine  $\mathbf{w}$  for the linear function

$f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}_j) = \Phi^T \mathbf{w}$  which minimizes our error function  $E(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$

- Least Squares error and corresponding estimates:

$$E^* = \min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \min_{\mathbf{w}} \left\{ \sum_{j=1}^m \left( \sum_{i=1}^n w_i \phi_i(\mathbf{x}_j) - y_j \right)^2 \right\} \quad (2)$$

*Sum across rows of  $\Phi$*  (pointing to  $m$ )  
*sum within row of  $\Phi$*  (pointing to  $n$ )

$E^*$  is attained at

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \arg \min_{\mathbf{w}} (\Phi \mathbf{w} - \mathbf{y})^T (\Phi \mathbf{w} - \mathbf{y})$$

*in matrix form?*

On test point  $\mathbf{x}$ , use  $f(\mathbf{x}, \mathbf{w}^*) = \mathbf{w}^* \cdot \Phi(\mathbf{x})$  & choice of  $\phi \dots$

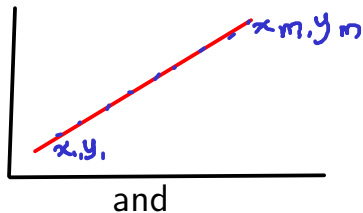
# Recap: Linear Regression is **not** Naively Linear

- Need to determine  $\mathbf{w}$  for the linear function  $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^n w_i \phi_i(\mathbf{x}_j) = \Phi \mathbf{w}$  which minimizes our error function  $E(f(\mathbf{x}, \mathbf{w}), \mathcal{D})$
- Least Squares error and corresponding estimates:

$$E^* = \min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \min_{\mathbf{w}} \left\{ \sum_{j=1}^m \left( \sum_{i=1}^n w_i \phi_i(\mathbf{x}_j) - y_j \right)^2 \right\} \quad (2)$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w}, \mathcal{D}) = \arg \min_{\mathbf{w}} \left( \underbrace{\mathbf{w}^T \Phi^T \Phi \mathbf{w} - 2 \mathbf{y}^T \Phi \mathbf{w} + \mathbf{y}^T \mathbf{y}}_{\|\Phi \mathbf{w} - \mathbf{y}\|_2^2} \right) \quad (3)$$

- The minimum value of the squared loss is zero
- If zero were attained at  $\mathbf{w}^*$ , we would have  $\forall j, \phi^T(x_j)\mathbf{w}^* = y_j$ , or equivalently  $\Phi\mathbf{w}^* = \mathbf{y}$ , where



$$\Phi = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_p(x_1) \\ \dots & \dots & \dots \\ \phi_1(x_m) & \dots & \phi_p(x_m) \end{bmatrix}$$

$\phi^T(x_1)$   
 $\phi^T(x_m)$

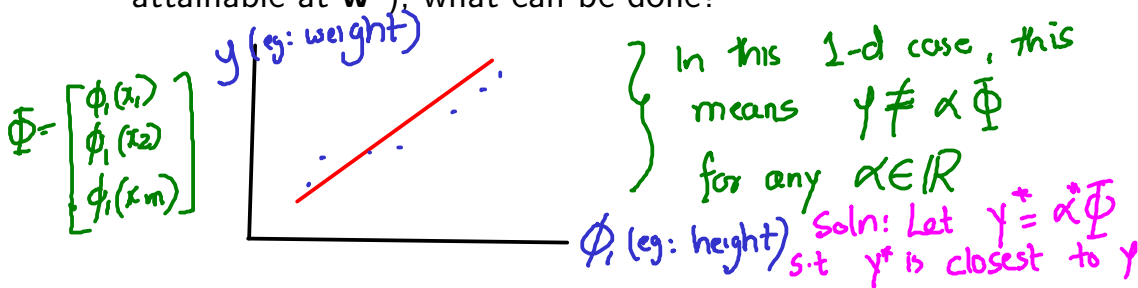
$$E = \|\Phi\mathbf{w}^* - \mathbf{y}\|^2 = 0$$

What if there is no perfect line fit?

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \dots \\ y_m \end{bmatrix}$$

- It has a solution if  $\mathbf{y}$  is in the column space (the subspace of  $\mathbb{R}^m$  formed by the column vectors) of  $\Phi$

- The minimum value of the squared loss is zero
- If  $\mathbf{y}$  is NOT in the column space of  $\Phi$  (that is, if zero is NOT attainable at  $\mathbf{w}^*$ ), what can be done?



# Geometric Interpretation of Least Square Solution

- Let  $\mathbf{y}^*$  be a solution in the column space of  $\Phi$
- The least squares solution is such that the distance between  $\mathbf{y}^*$  and  $\mathbf{y}$  is minimized
- Therefore.....

If # of pts increases beyond 2, the number of coordinates in the figure will increase

If # of features ( $\phi_i$ 's) increases beyond 1, the dimensionality of  $c(\Phi)$  will grow to a plane

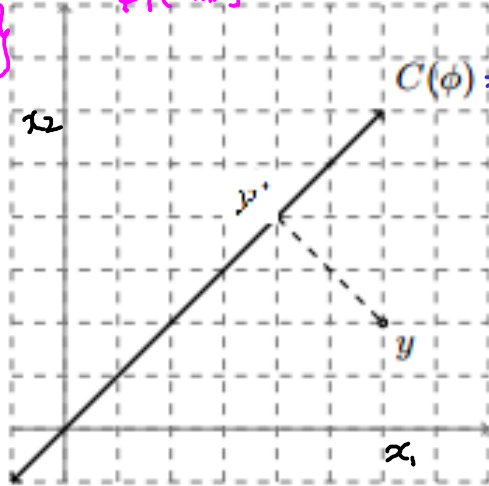
For 1-dim case,  $\Phi = \begin{bmatrix} \phi_1(x_1) \\ \phi_1(x_m) \end{bmatrix} \in \mathbb{R}^m$

$$C(\Phi) = \{ \Phi \alpha \mid \alpha \in \mathbb{R} \}$$

= one dim  
vector  
in  $\mathbb{R}^m$

$$y^* \in \mathbb{R}^m$$

$$y \in \mathbb{R}^m$$



$$C(\Phi) = \{ \Phi w \mid w \in \mathbb{R}^n \}$$

$y$  will be closest  
to  $y^*$  when  
 $(y - y^*) \perp C(\Phi)$

This plot depicts  $C(\cdot)$  for 1 feature  $\phi_1$  and 2 data points:  $x_1$  &  $x_2$  ( $m=2$ )

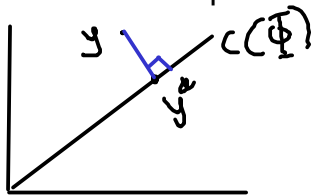


# Geometric Interpretation of Least Square Solution

④ • Let  $\mathbf{y}^*$  be a solution in the column space of  $\Phi$  [ $\mathbf{y}^* = \Phi \mathbf{w}^* \in c(\Phi)$ ]

• The least squares solution  $\mathbf{w}^*$  is such that the distance between  $\mathbf{y}^*$  and  $\mathbf{y}$  is minimized

⑤ { • Therefore, the line joining  $\mathbf{y}^*$  to  $\mathbf{y}$  should be orthogonal to the column space



$$\Phi \mathbf{w}^* = \mathbf{y}^*$$

④

$$(\mathbf{y} - \mathbf{y}^*)^T \Phi = 0$$

⑤

$$\begin{aligned} (\mathbf{y} - \mathbf{y}^*) \perp c(\Phi) &\equiv (\mathbf{y} - \mathbf{y}^*) \perp \text{span}(\text{columns of } \Phi) \\ &\equiv (\mathbf{y} - \mathbf{y}^*) \perp \Phi \end{aligned}$$

$$(\mathbf{y}^*)^T \Phi = (\mathbf{y})^T \Phi$$

⑥

$$(y^*)^T \Phi = y^T \Phi$$

Substitute  
using (4)



From (4),  $y^* = \Phi w$

$$(\Phi w)^T \Phi = y^T \Phi$$

$$w^T \Phi^T \Phi = y^T \Phi$$

$$\Phi^T \Phi w = \Phi^T y$$

$$w = (\Phi^T \Phi)^{-1} \Phi^T y \quad (10)$$

$$(AV)^T = V^T A^T \quad (7)$$

$$(AB)^T = B^T A^T \quad (8)$$

Taking transposes  
on both sides (9)

- Here  $\Phi^T \Phi$  is invertible if and only if  $\Phi$  has full column rank

ie. no feature set for one data point can be obtained using the feature sets of other data points

$$(\Phi \mathbf{w})^T \Phi = \mathbf{y}^T \Phi \quad (7)$$

$$\mathbf{w}^T \Phi^T \Phi = \mathbf{y}^T \Phi \quad (8)$$

$$\Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{y} \quad (9)$$

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad (10)$$

- Here  $\Phi^T \Phi$  is invertible if and only if  $\Phi$  has full column rank

**PROOF?**

Proof:

**Claim** :  $\Phi^T \Phi$  is invertible if and only if  $\Phi$  is full column rank

Proof :

$\Leftarrow$  Given that  $\Phi$  has full column rank and hence columns are linearly independent, we have that  $\Phi \mathbf{w} = 0 \Rightarrow \mathbf{w} = 0$  (only vector in null space is 0)

Assume on the contrary that  $\Phi^T \Phi$  is non invertible. Then  $\exists \mathbf{w} \neq 0$  such that  $\Phi^T \Phi \mathbf{w} = 0$

$$\begin{aligned} & \Rightarrow \mathbf{w}^T \Phi^T \Phi \mathbf{w} = 0 \quad (\text{just premultiplying } 0 \text{ vec by } 0) \\ & \Rightarrow (\Phi \mathbf{w})^T \Phi \mathbf{w} = 0 \\ & \Rightarrow \Phi \mathbf{w} = 0 \\ & \quad \|\Phi \mathbf{w}\|^2 = 0 \text{ iff } \Phi \mathbf{w} = 0 \end{aligned}$$

$(\Phi \mathbf{w})^T$

This is a contradiction. Hence  $\Phi^T \Phi$  is invertible if  $\Phi$  is full column rank

**Claim** :  $\Phi^T \Phi$  is invertible if and only if  $\Phi$  is full column rank

Proof :

$\implies$  If  $\Phi^T \Phi$  is invertible and  $\Phi \mathbf{w} = 0$ , then  $(\Phi^T \Phi \mathbf{w}) = 0$ , which in turn implies  $\mathbf{w} = 0$ . This implies  $\Phi$  has full column rank if  $\Phi^T \Phi$  is invertible. The converse can also be proved similarly.

# Later: More on Optimization

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y \text{ minimizes } E(\cdot) = \|\Phi w - y\|_2^2$$

- More generally: How to minimize a function?
  - Level Curves and Surfaces
  - Gradient Vector
  - Directional Derivative
  - Hyperplane
  - Tangential Hyperplane
- Iterative Algorithms such as (Stochastic) Gradient Descent Algorithm

# Building on questions on Least Squares Linear Regression

- ✓ 1 Is there a probabilistic interpretation?
  - Gaussian Error, Maximum Likelihood Estimate
- 2 Addressing overfitting
  - Bayesian and Maximum A posteriori Estimates, Regularization
- 3 How to minimize the resultant and more complex error functions?
  - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

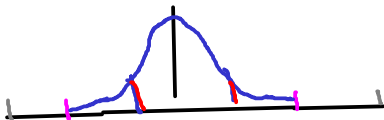


# Building on questions on Least Squares Linear Regression

- 1 Is there a probabilistic interpretation?
  - Gaussian Error, Maximum Likelihood Estimate
- 2 Addressing overfitting
  - Bayesian and Maximum A posteriori Estimates, Regularization
- 3 How to minimize the resultant and more complex error functions?
  - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

# Probabilistic Modeling of Linear Regression

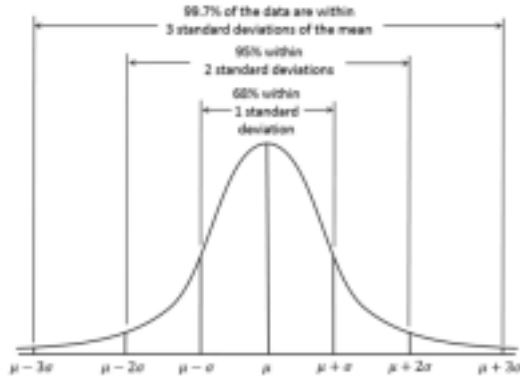
- Linear Model:  $Y$  is a linear function of  $\phi(x)$ , subject to a random noise variable  $\varepsilon$  which we believe is 'mostly' bounded by some threshold  $\sigma$ :



$$Y = \underbrace{w^T \phi(x)} + \underbrace{\varepsilon}_{\text{Random noise}}$$
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- Motivation:  $\mathcal{N}(\mu, \sigma^2)$ , has maximum entropy among all real-valued distributions with a specified variance  $\sigma^2$
- 3 -  $\sigma$  rule: About 68% of values drawn from  $\mathcal{N}(\mu, \sigma^2)$  are within one standard deviation  $\sigma$  away from the mean  $\mu$ ; about 95% of the values lie within  $2\sigma$ ; and about 99.7% are within  $3\sigma$ .

$$H(p) = - \int p(x) \ln p(x) dx = E[-\ln p(x)]$$



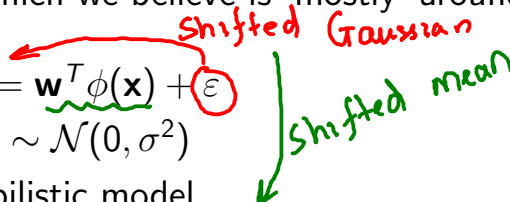
**Figure 2:** 3 –  $\sigma$  rule: About 68% of values drawn from  $\mathcal{N}(\mu, \sigma^2)$  are within one standard deviation  $\sigma$  away from the mean  $\mu$ ; about 95% of the values lie within  $2\sigma$ ; and about 99.7% are within  $3\sigma$ . Source: [https://en.wikipedia.org/wiki/Normal\\_distribution](https://en.wikipedia.org/wiki/Normal_distribution)

# Probabilistic Modeling of Linear Regression

- Linear Model:  $Y$  is a linear function of  $\phi(\mathbf{x})$ , subject to a random noise variable  $\varepsilon$  which we believe is 'mostly' around some threshold  $\sigma$ :

$$Y = \underbrace{\mathbf{w}^T \phi(\mathbf{x})}_{\text{shifted mean}} + \underbrace{\varepsilon}_{\text{shifted Gaussian}}$$

$\varepsilon \sim \mathcal{N}(0, \sigma^2)$



- This allows for the Probabilistic model

$$P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \mathcal{N}(\underbrace{\mathbf{w}^T \phi(\mathbf{x}_j)}_{\text{shifted mean}}, \sigma^2)$$

$$P(y | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \prod_{j=1}^m P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2)$$

- Note:  $E[Y(\mathbf{w}, \mathbf{x}_j)] = \mathbf{w}^T \phi(\mathbf{x}_j)$

# Probabilistic Modeling of Linear Regression

- Linear Model:  $Y$  is a linear function of  $\phi(\mathbf{x})$ , subject to a random noise variable  $\varepsilon$  which we believe is 'mostly' around some threshold  $\sigma$ :

$$Y = \mathbf{w}^T \phi(\mathbf{x}) + \varepsilon$$
$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

- This allows for the Probabilistic model

$$P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \mathcal{N}(\mathbf{w}^T \phi(\mathbf{x}_j), \sigma^2)$$

$$P(y | \mathbf{w}, \mathbf{x}_j, \sigma^2) = \prod_{j=1}^m P(y_j | \mathbf{w}, \mathbf{x}_j, \sigma^2)$$

- Note:  $E[Y(\mathbf{w}, \mathbf{x}_j)] = \mathbf{w}^T \phi(\mathbf{x}_j)$   
 $= \mathbf{w}_0^T + \mathbf{w}_1^T \phi_1(\mathbf{x}_j) + \dots + \mathbf{w}_n^T \phi_n(\mathbf{x}_j)$