

---

# Algorithm: MBGD Linear Regression Learner (with constant stepsize)

---

input:  $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$ ,  $\gamma$ ,  $T$ ,  $b$

$\vec{w} \leftarrow$  random vector in  $\mathbb{R}^{d+1}$

$$M \in \text{floor}\left(\frac{n}{b}\right)$$

for  $t = 1, \dots, T$

randomly shuffle  $D$

for  $m = 1, \dots, M$

$$\nabla \hat{L}(\vec{w}) \leftarrow \frac{2}{b} \sum_{i=mb+1}^{(m+1)b} (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i$$

$$\vec{w} \leftarrow \vec{w} - \gamma \nabla \hat{L}(\vec{w})$$

return  $\hat{f}(\vec{x}) = \vec{x}^T \vec{w}$

---

Setting  $b = n \Rightarrow M = 1 \Rightarrow$  "gradient descent"

$b = 1 \Rightarrow M = n \Rightarrow$  "stochastic gradient  
descent (SGD)"

in this course

in general  $\rightarrow b$  is small

computation:

$$\text{BGD} : O\left(\underbrace{d n T}_{\nabla \hat{L}(\vec{\omega})}\right)$$

# epochs

$$\text{MBGD} : O\left(\underbrace{d b M T}_{\nabla \hat{L}(\vec{\omega})}\right)$$

batch size

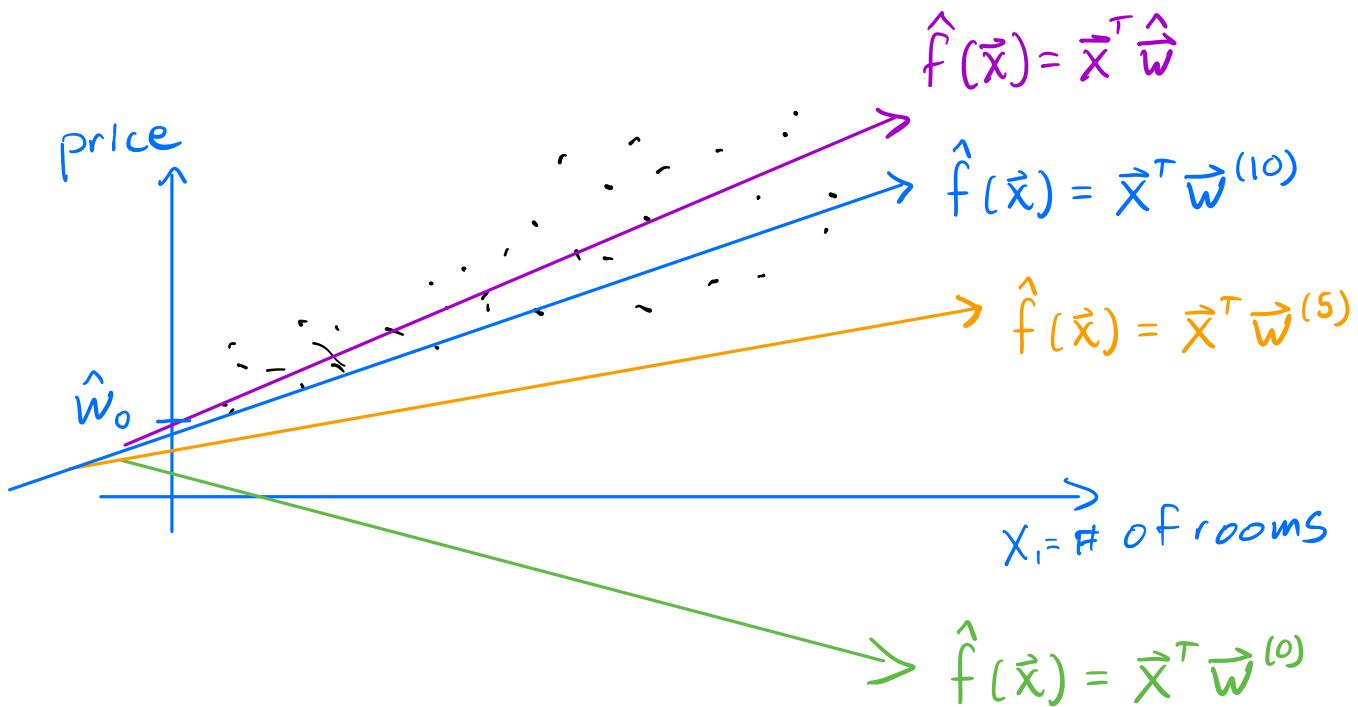
# batches

For same  $T$ : computation identical for BGD & MBGD

but

In practice, for the same  $T$ , MBGD usually finds better  $\vec{\omega}$  (when  $b < n$ )

Ex:

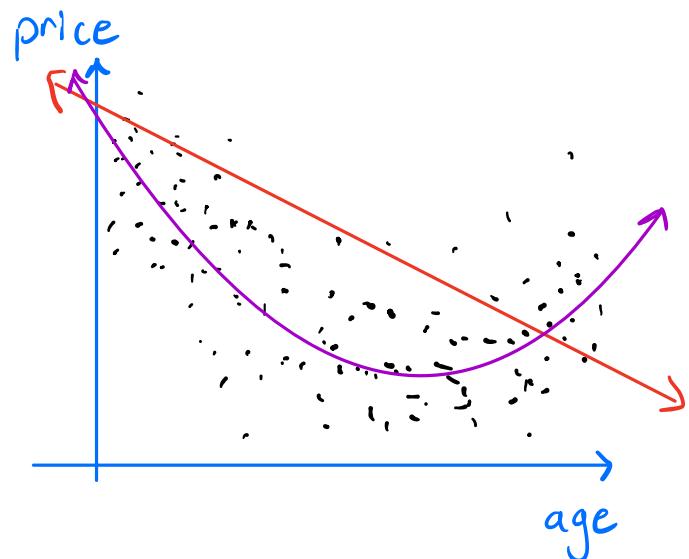
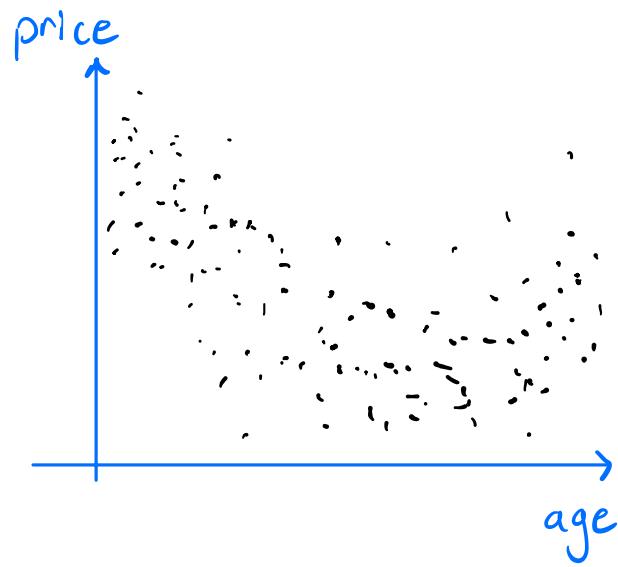


Should we always use a linear function class  $F$ ?

→ No, only if it fits your data

# Polynomial Regression

Ex: Predicting car price based on age of the car



Let  $\tilde{F}$  contain polynomial functions instead of just linear functions

Ex:  $d=1$ ,  $X = \mathbb{R}^{d+1} = \mathbb{R}^2$ ,  $y = \mathbb{R}$

linear function:  $\vec{x}^T \vec{w} = (1, x_1) (w_0, w_1)^T = w_0 + x_1 w_1, \vec{w} \in \mathbb{R}^2$

$\tilde{F}_1 = \{f \mid f: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \vec{x}^T \vec{w} = w_0 + x_1 w_1, \vec{w} \in \mathbb{R}^2\}$

degree 2 polynomial:  $w_0 + x_1 w_1 + x_1^2 w_2, \vec{w} \in \mathbb{R}^3$

"feature map"  $= \phi_2(\vec{x})^T \vec{w} \quad \text{Linear in } \phi_2(\vec{x})$

where  $\phi_2(\vec{x}) = (x_0=1, x_1, x_1^2)^T \in \mathbb{R}^3$  transformed set of features

$\tilde{F}_2 = \{f \mid f: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \phi_2(\vec{x})^T \vec{w}, \vec{w} \in \mathbb{R}^3\}, \tilde{F}_1 \subset \tilde{F}_2$

degree 3 polynomial:  $w_0 + x_1 w_1 + x_1^2 w_2 + x_1^3 w_3 \quad \vec{w} \in \mathbb{R}^4$   
 $= \phi_3(\vec{x})^\top \vec{w}$  linear in  $\phi_3(\vec{x})$

where  $\phi_3(\vec{x}) = (1, x_1, x_1^2, x_1^3)^\top \in \mathbb{R}^4$

$$\widetilde{F}_3 = \{f \mid f: \mathbb{R}^2 \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \phi_3(\vec{x})^\top \vec{w}, \vec{w} \in \mathbb{R}^4\}$$

$$\widetilde{F}_1 \subset \widetilde{F}_2 \subset \widetilde{F}_3$$

Ex:  $d=2, \mathcal{X} = \mathbb{R}^{d+1} = \mathbb{R}^3, \mathcal{Y} = \mathbb{R}$

linear function:  $\vec{x}^\top \vec{w} = (1, x_1, x_2) (w_0, w_1, w_2)^\top$   
 $= w_0 + x_1 w_1 + x_2 w_2 \quad \vec{w} \in \mathbb{R}^3$

$$\widetilde{F}_1 = \{f \mid f: \mathbb{R}^3 \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \vec{x}^\top \vec{w}, \vec{w} \in \mathbb{R}^3\}$$

degree 2 polynomial:  $w_0 + x_1 w_1 + x_2 w_2 + x_1^2 w_3 + x_2^2 w_4 + x_1 x_2 w_5$   
 $= \phi_2(\vec{x})^\top \vec{w}, \quad \vec{w} \in \mathbb{R}^6 \quad \text{linear in } \phi_2(\vec{x})$

where  $\phi_2(\vec{x}) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)^\top \in \mathbb{R}^6$

$$\widetilde{F}_2 = \{f \mid f: \mathbb{R}^3 \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \phi_2(\vec{x})^\top \vec{w}, \vec{w} \in \mathbb{R}^6\}$$

$$\widetilde{F}_1 \subset \widetilde{F}_2$$

In general  $\phi_p: \mathcal{X} \rightarrow \mathbb{Z}$  is a degree  $p$  polynomial  
 "feature map" binomial coefficient  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

For  $\mathcal{X} = \mathbb{R}^{d+1}$ ,  $\mathbb{Z} = \mathbb{R}^{\bar{P}}$  where  $\bar{P} = \binom{(d+1)+p-1}{p} = \binom{d+p}{p}$

$$\underline{\text{Ex}}: d=2, p=2 \Rightarrow \bar{P} = \binom{2+2}{2} = \binom{4}{2} = \frac{4 \cdot 3}{2 \cdot 1} = 6$$

$$\underline{\text{Ex}}: d=2, p=3 \Rightarrow \bar{P} = \binom{2+3}{2} = \binom{5}{2} = \frac{5 \cdot 4}{2} = 10$$

The function class becomes:

$$\widetilde{F}_p = \left\{ f \mid f: \mathbb{R}^{d+1} \rightarrow \mathbb{R} \text{ and } f(\vec{x}) = \phi_p(\vec{x})^T \vec{w}, \vec{w} \in \mathbb{R}^{\bar{P}} \right\}$$

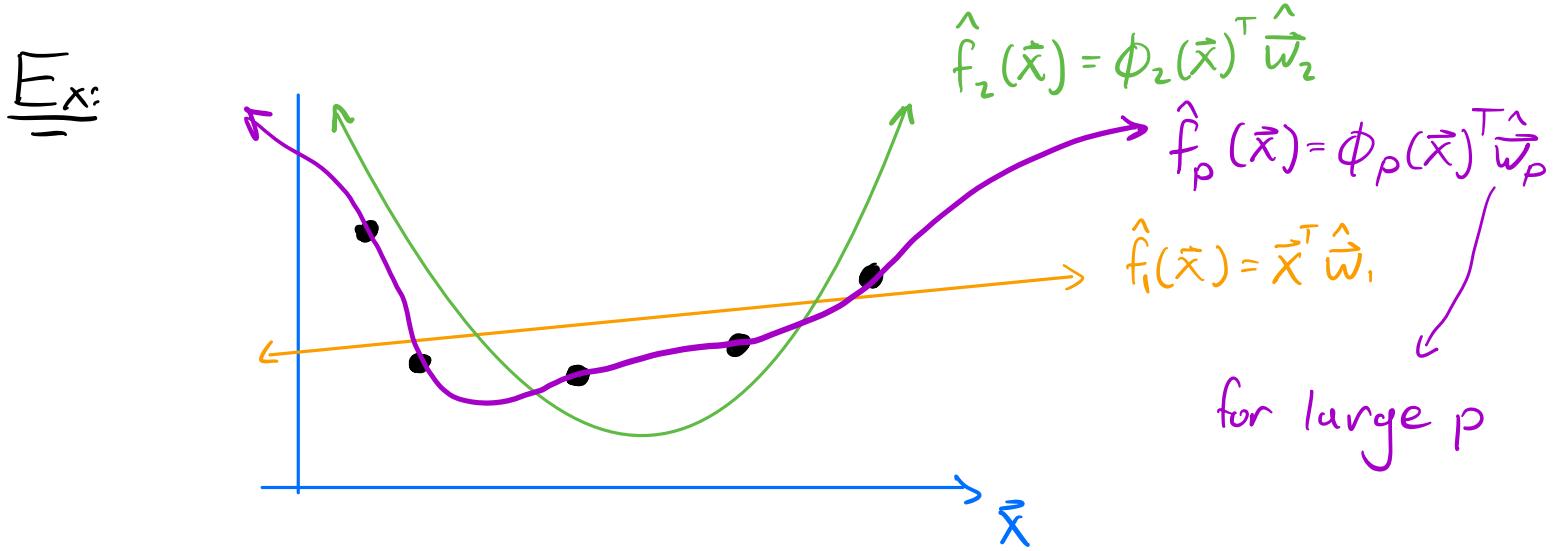
$$\widetilde{F}_1 \subset \widetilde{F}_2 \subset \dots \subset \widetilde{F}_p$$

You can then solve for  $\hat{\vec{w}}_p = \arg \min_{\vec{w} \in \mathbb{R}^{\bar{P}}} \hat{L}_p(\vec{w})$

$$\text{where } \hat{L}_p(\vec{w}) = \frac{1}{n} \sum_{i=1}^n l(\phi_p(\vec{x}_i)^T \vec{w}, y_i)$$

using the closed form solution

or gradient descent as before



$$\hat{L}(\hat{f}_1) \geq \hat{L}(\hat{f}_2) \geq \dots \geq \hat{L}(\hat{f}_p) \approx 0$$

Why not set  $p$  as large as possible?

what does  $\hat{L}(\hat{f}_p)$  look like?

