# EECS 545 → Machine Learning

## Lecture 2 → Linear Regression (Part 1)

- Supervised Learning → given data $x$ in feature space and corresponding labels $y$, learn to predict $y$ from $x$
- Classification → discrete-valued labels
- Regression → continuous-valued labels

- Notation → $x \in \mathbb{R}^D$ = data (scalar or vector)

  $\phi_j(x) \in \mathbb{R}$ = j-th feature for $x$ (scalar), $j = 0, \ldots, M-1$

  $\phi(x) \in \mathbb{R}^M$ = features for $x$ (vector)

  $y \in \mathbb{R}$ = continuous-valued label

  $x^{(n)}$ = n-th training example

  $y^{(n)}$ = n-th training label

- We want to learn a function $h(x, w) \approx y$ to predict future values
- 0-th order polynomial → $h(x, w) = w_0$
- 1-st order polynomial → $h(x, w) = w_0 + w_1 x$     } 1-Dimensional features
- 3-rd order polynomial → $h(x, w) = w_0 + w_1 x + w_2 x^2 + w_3 x^3$

- General case → $h(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x)$
- $h(x, w)$ is linear in parameters $w$
- For simplicity, we convert $w_0$ to a bias term

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

$$w = (w_0, \ldots, w_{M-1})^T$$
$$\phi(x) = (\phi_0(x), \ldots, \phi_{M-1}(x))^T$$
$$\phi_0(x) = 1$$

## Error Functions

- Sum of squares → $E(w) = \frac{1}{2} \sum_{n=1}^{N} \{ \underbrace{h(x^{(n)}, w)}_{\text{prediction}} - \underbrace{y^{(n)}}_{\text{target}} \}^2$

- Want to find $w$ that minimizes $E(w)$ over the training data

- Gradient of SoS → $\dfrac{\partial E(w)}{\partial w_k} = \dfrac{\partial}{\partial w_k} \frac{1}{2} \sum_{n=1}^{N} \left( \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - y^{(n)} \right)^2$

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

$$= \sum_{n=1}^{N} \left( \underbrace{\sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - y^{(n)}}_{\text{error}} \right) \phi_k(x^{(n)})$$

- Batch Gradient Descent → Repeat until convergence,

$$w := w - \eta \, \nabla_w E(w)$$

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

· Stochastic Gradient Descent → instead of computing batch gradient descent which is over the entire dataset, compute gradient for individual examples and update

· Repeat until convergence,

$$w := w - \alpha \nabla_w E(w | x^{(n)}) \qquad \nabla E(w | x^{(n)}) = \sum_{j=0}^{M-1} w_j \phi_j(x^{(n)}) - y^{(n)}) \phi(x^{(n)})$$
$$= w^T \phi(x^{(n)}) - y^{(n)}) \phi(x^{(n)})$$

· In SGD, iterate over entire dataset with multiple epochs until convergence

## Closed form solution

· Compute gradient, then set equal to 0
· Solve the equation in a closed form

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

$$= \frac{1}{2} \sum_{n=1}^{N} \left( w^T \phi(x^{(n)}) - y^{(n)} \right)^2$$

$$= \frac{1}{2} \sum_{n=1}^{N} \left( w^T \phi(x^{(n)}) \right)^2 - \sum_{n=1}^{N} y^{(n)} w^T \phi(x^{(n)}) + \frac{1}{2} \sum_{n=1}^{N} (y^{(n)})^2$$

$$= \frac{1}{2} w^T \Phi^T \Phi w - w^T \Phi^T y + \frac{1}{2} y^T y$$

where we define $\Phi$ to be the design matrix, $N \times M$ containing the M basis functions (columns) and N data points (rows)

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

· We now need to compute the derivative in matrix form,

$$\nabla_w E(w) = \nabla_w \left( \frac{1}{2} w^T \Phi^T \Phi w - w^T \Phi^T y + \frac{1}{2} y^T y \right)$$

$$= \Phi^T \Phi w - \Phi^T y \qquad (\text{set equal to 0})$$

· Solving this equation gives us,

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

$$w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T y$$

· We call this the Moore-Penrose pseudo inverse : $\Phi^\dagger = (\Phi^T \Phi)^{-1} \Phi^T$ applied to $\Phi w \approx y$

<u>Aside → Calculating Matrix Gradients</u>

· Suppose $f: \mathbb{R}^{M \times n} \to \mathbb{R}$ is a function that takes an input matrix $A$ of size $M \times n$ and returns a real value
· The gradient of $f$ w.r.t. $A \in \mathbb{R}^{M \times n}$ is,

$$\nabla_A f(A) \in \mathbb{R}^{M \times n} = \begin{bmatrix} \dfrac{\partial f(A)}{\partial A_{11}} & \dfrac{\partial f(A)}{\partial A_{12}} & \cdots & \dfrac{\partial f(A)}{\partial A_{1n}} \\ \dfrac{\partial f(A)}{\partial A_{21}} & \dfrac{\partial f(A)}{\partial A_{22}} & \cdot & \dfrac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial f(A)}{\partial A_{M1}} & \dfrac{\partial f(A)}{\partial A_{M2}} & \cdots & \dfrac{\partial f(A)}{\partial A_{Mn}} \end{bmatrix} \qquad (\nabla_A f(A))_{ij} = \dfrac{\partial f(A)}{\partial A_{ij}}$$

· E.g. if $A$ is just a vector, $x \in \mathbb{R}^n$, then $\nabla_x f(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\ \dfrac{\partial f(x)}{\partial x_2} \\ \vdots \\ \dfrac{\partial f(x)}{\partial x_n} \end{bmatrix}$

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

· Linear function → $f(x) = \sum_{i=1}^{n} b_i x_i = b^T x$
· Gradient → $\dfrac{\partial f(x)}{\partial x_k} = \dfrac{\partial}{\partial x_k} \sum_{i=1}^{n} b_i x_i = b_k$

· Compact form → $\nabla_x f(x) = b$

· Quadratic function → $f(x) = \sum_{i,j=1}^{n} x_i A_{ij} x_j = x^T A x$
· Gradient → $\dfrac{\partial f(x)}{\partial x_k} = 2 \sum_{j=1}^{n} A_{kj} x_j = 2(Ax)_k$

· Compact form → $\nabla_x f(x) = 2Ax$