# Linear Regression

Lecture notes by Kevyn Collins-Thompson
Applied Machine Learning (Coursera)

Lecture notes by Ethem Alpaydın
Introduction to Machine Learning (Boğaziçi Üniversitesi)

Lecture notes by Andrew NG
Machine Learning by Stanford University (Coursera)

# Linear regression with one variable
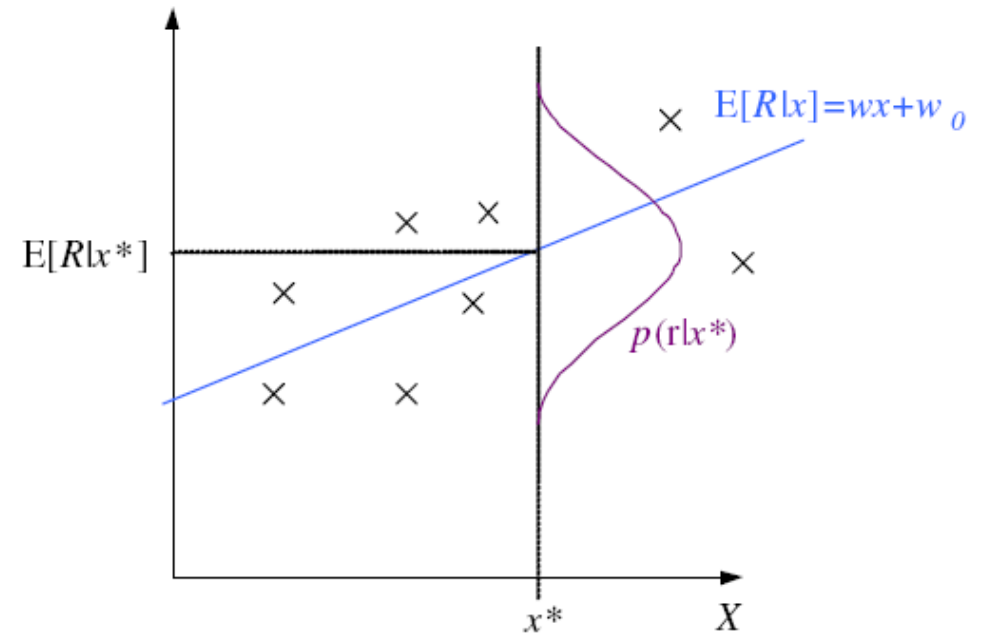
Machine Learning

# Regression

$$r = f(x) + \varepsilon$$

$$\text{estimator}: g(x|\theta)$$

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

$$p(r|x) \sim \mathcal{N}\left(g(x|\theta), \sigma^2\right)$$

$$\mathcal{L}(\theta|\mathcal{X}) = \log \prod_{t=1}^{N} p\left(x^t, r^t\right)$$

$$= \log \prod_{t=1}^{N} p\left(r^t|x^t\right) + \log \prod_{t=1}^{N} p\left(x^t\right)$$



$E[R|x] = wx + w_0$

$E[R|x*]$

$p(r|x*)$

$x*$

$X$

# Regression: From LogL to Error

$$\mathcal{L}(\theta|\mathcal{X}) = \log \prod_{t=1}^{N} \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{\left[r^t - g\left(x^t|\theta\right)\right]^2}{2\sigma^2}\right]$$

$$= -N\log\sqrt{2\pi}\sigma - \frac{1}{2\sigma^2}\sum_{t=1}^{N}\left[r^t - g\left(x^t|\theta\right)\right]^2$$

$$E(\theta|\mathcal{X}) = \frac{1}{2}\sum_{t=1}^{N}\left[r^t - g\left(x^t|\theta\right)\right]^2$$

Most frequently used error function
$$E = -logl$$

$\theta$ minimize the error function are called the least squares estimates.

# Linear Regression

$$g\left(x^t \mid w_1, w_0\right) = w_1 x^t + w_0$$

$$\sum_t r^t = N w_0 + w_1 \sum_t x^t$$

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t \left(x^t\right)^2$$

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t \left(x^t\right)^2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix} \quad \Longrightarrow \quad \mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$$
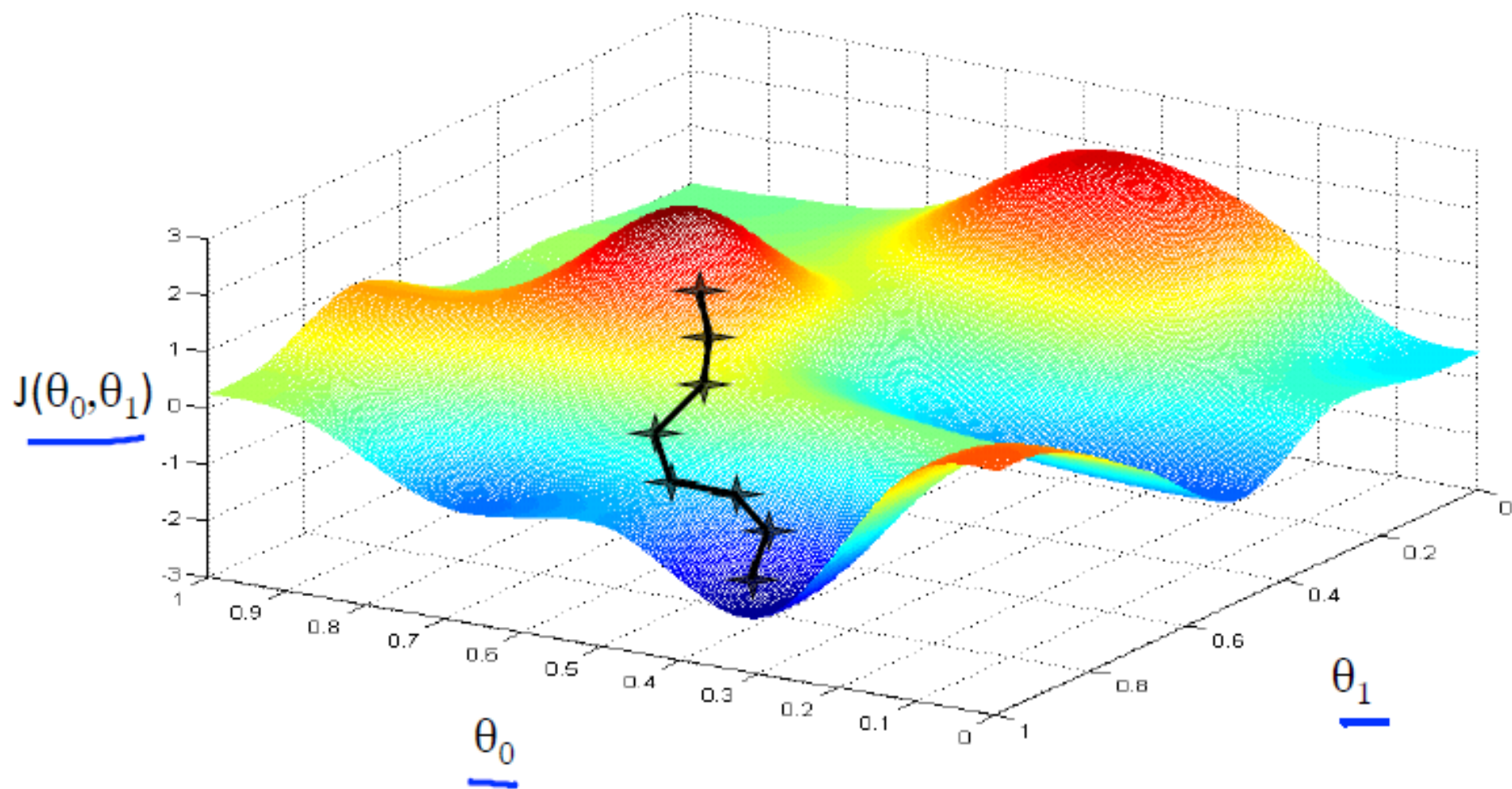
# Gradient Descent

Have some function $J(\theta_0, \theta_1)$

Want $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

**Outline**

Start with some $\theta_0, \theta_1$

Keep chaning $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum

Machine Learning

# Gradient Descent

```
repeat until convergence {
```

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

```
                    (for j=0 and j=1)
}
```

Simultaneously update

# Gradient Descent

Simultaneously update

$$temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$
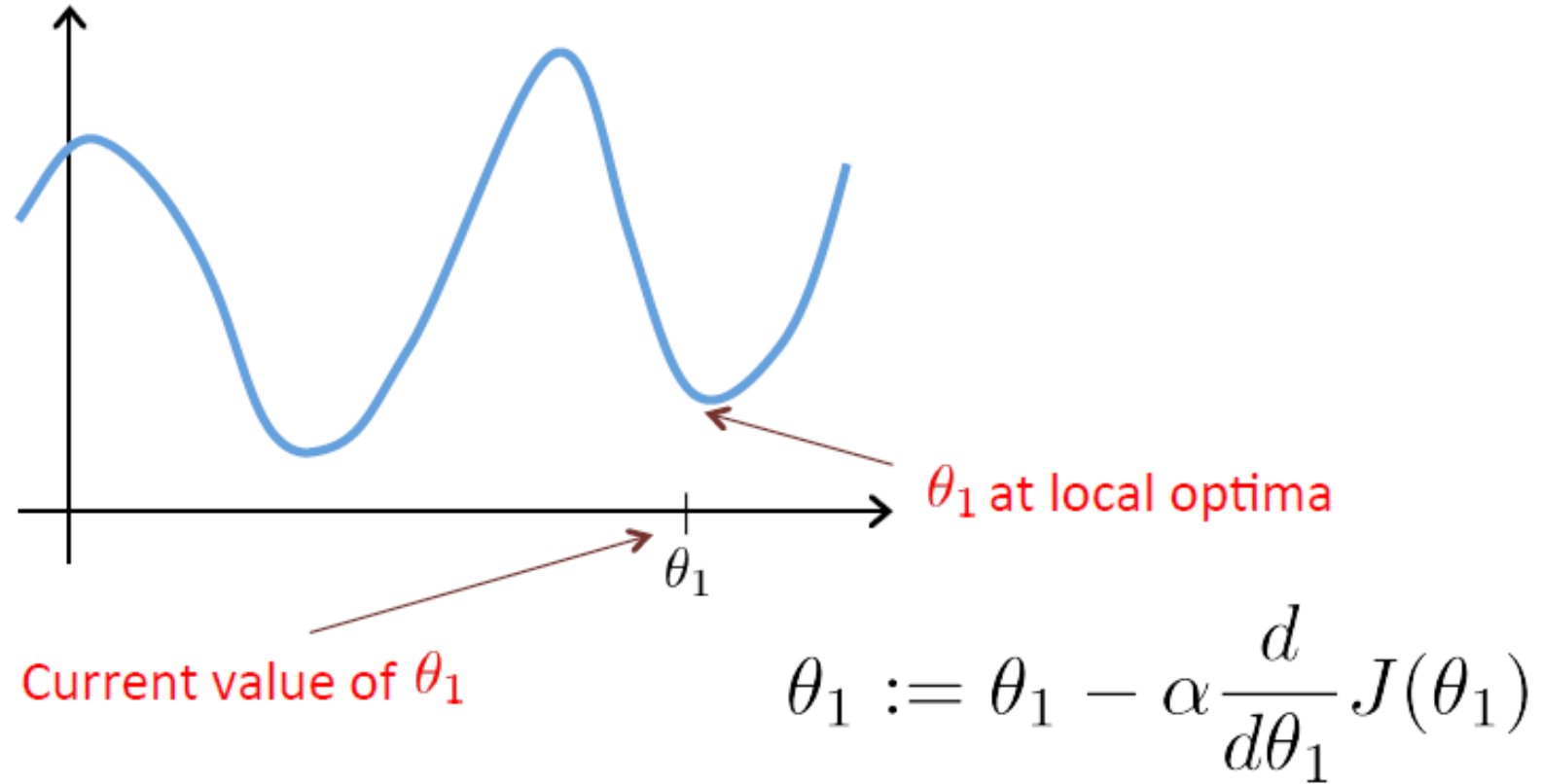
$$\theta_0 := temp0$$

$$\theta_1 := temp1$$

Machine Learning

# Gradient Descent

If $\alpha$ (learning rate) is too small, gradient descent can be slow.

If $\alpha$ (learning rate) is too large, gradient descent can overshoot the minimum. It may fail to converge.

# Gradient Descent



$\theta_1$ at local optima

Current value of $\theta_1$

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

# Gradient Descent

Gradient descent can converge to a local minimum, even with the learning rate $\alpha$ fixed.

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease $\alpha$ over time.

# Gradient descent for linear regression

**Gradient Descent**

```
repeat until convergence {
```

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

```
            (for j=0 and j=1)
```

```
}
```

$$g\left(x^t \mid w_1, w_0\right) = w_1 x^t + w_0$$

**Error function – linear regression**

$$E(w|X) = \frac{1}{2} \sum_{t=1}^{N} [r^t - g(x^t|w)]^2$$

# Gradient descent for linear regression

```
repeat until convergence {
```
$$w_0 := w_0 - \alpha \sum_{t=1}^{N} (g(x^t|w) - r^t)$$

$$w_1 := w_1 - \alpha \sum_{t=1}^{N} (g(x^t|w) - r^t).x^t$$
```
}
```

Simultaneously update

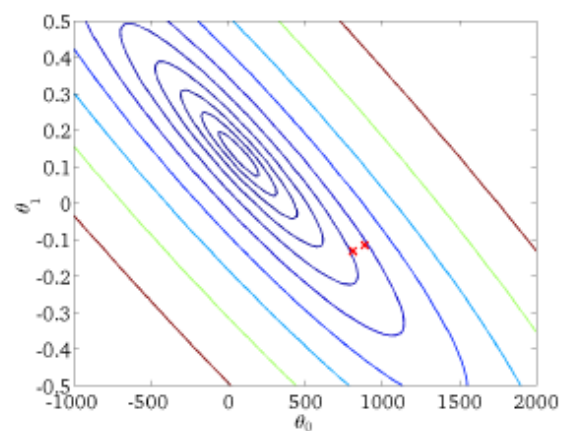Machine Learning

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

Andrew Ng

•••

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

Andrew Ne

# Linear regression with multiple variables

# Multivariate Regression

$$g(x^t|w_0, w_1, \ldots, w_d) = w_0 + w_1 x_1^t + w_2 x_2^t + \cdots + w_d x_d^t = w^T x^t$$

$$E(w|X) = \frac{1}{2} \sum_{t=1}^{N} [r^t - g(x^t|w)]^2$$

# Multivariate Regression

- Normal Equation

$$\sum_t r^t = N w_0 + w_1 \sum_t x_1^t + w_2 \sum_t x_2^t + \cdots + w_d \sum_t x_d^t$$

$$\sum_t x_1^t r^t = w_0 \sum_t x_1^t + w_1 \sum_t (x_1^t)^2 + w_2 \sum_t x_1^t x_2^t + \cdots + w_d \sum_t x_1^t x_d^t$$

$$\sum_t x_2^t r^t = w_0 \sum_t x_2^t + w_1 \sum_t x_1^t x_2^t + w_2 \sum_t (x_2^t)^2 + \cdots + w_d \sum_t x_2^t x_d^t$$

$$\vdots$$

$$\sum_t x_d^t r^t = w_0 \sum_t x_d^t + w_1 \sum_t x_d^t x_1^t + w_2 \sum_t x_d^t x_2^t + \cdots + w_d \sum_t (x_d^t)^2$$

# Multivariate Regression

- Normal Equation

$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \cdots & x_d^1 \\ 1 & x_1^2 & x_2^2 & \cdots & x_d^2 \\ \vdots & & & & \\ 1 & x_1^N & x_2^N & \cdots & x_d^N \end{bmatrix}, \boldsymbol{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d, \end{bmatrix}, \boldsymbol{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

$$X^T X w = X^T r \;\; \rightarrow w = (X^T X)^{-1} X^T r$$

# Example

$$w = (X^T X)^{-1} X^T r$$

| | Size (feet$^2$) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|---|---|---|---|---|---|
| $x_0$ | $x_1$ | $x_2$ | $x_3$ | $x_4$ | r |
| 1 | 2104 | 5 | 1 | 45 | 460 |
| 1 | 1416 | 3 | 2 | 40 | 232 |
| 1 | 1534 | 3 | 2 | 30 | 315 |
| 1 | 852 | 2 | 1 | 36 | 178 |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad r = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

# Gradient Descent

```
repeat until convergence {
```
$$w_j := w_j - \alpha \sum_{t=1}^{N} (g(x^t|w) - r^t).x_j^t$$
```
}
```

# Feature Scaling

Idea: Make sure features are on a similar scale

**_Mean normalization_**

$$x' = \frac{x - mean(x)}{\max(x) - \min(x)}$$

$x'$ : normalized value

# Debugging

- How to make sure that gradient descent is working correctly

- How to choose learning rate

Convergence Plot (number of iteration vs $J(\theta_0, \theta_1)$)

$$\min_\theta J(\theta)$$