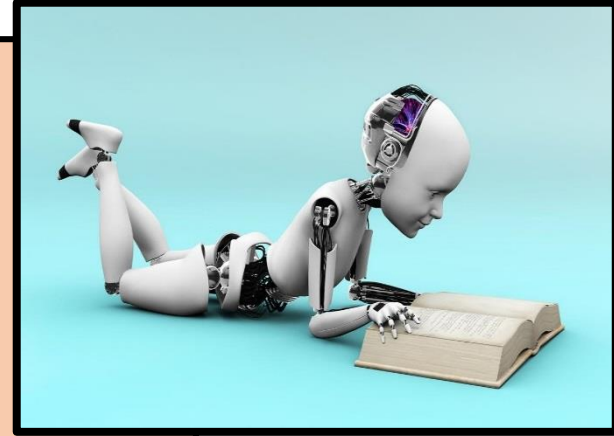# Machine learning from data
# Class1:
# Linear Regression



Zohar Yakhini

IDC

# Example: House Pricing

- We want to know the price of a house as a function of its size (in sqft).

- We want to learn a function from the size of the house x to the price $y = f(x)$ so we would be able to answer the above question

- Training set: 10 house instances with <u>feature</u> values and <u>labels</u>

| Square Feet (x) | House Price in $1000s (y) |
|---|---|
| 1400 | 245 |
| 1600 | 312 |
| 1700 | 279 |
| 1875 | 308 |
| 1100 | 199 |
| 1550 | 219 |
| 2350 | 405 |
| 2450 | 324 |
| 1425 | 319 |
| 1700 | 255 |

**Statistics:**
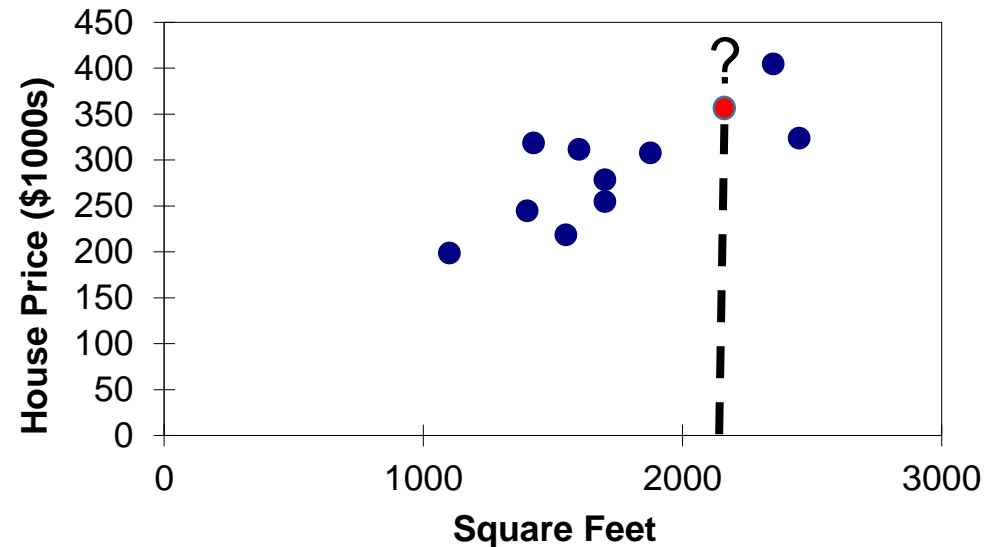**Dependent** variable (y) = house price
**Explaining** variable (x) = square footage

IDC HERZLIYA
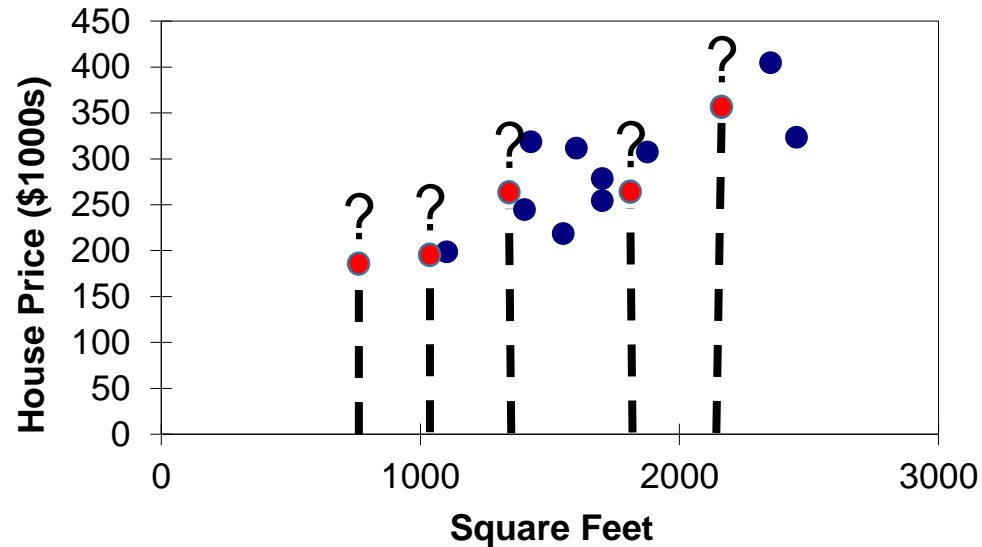
# Graphical Representation

Scatter plot of House Price (y) vs House Size (x)

<u>Prediction:</u>

Given house of size $x$, what would be its price $y = f(x)$?
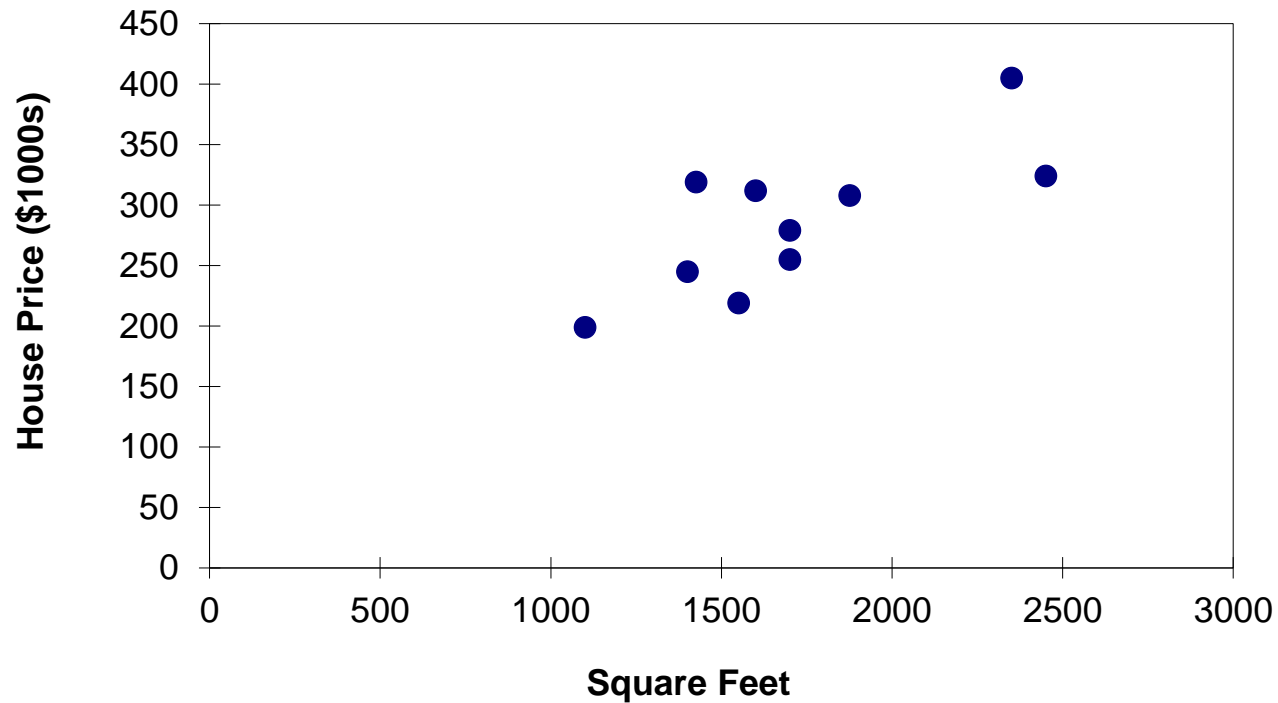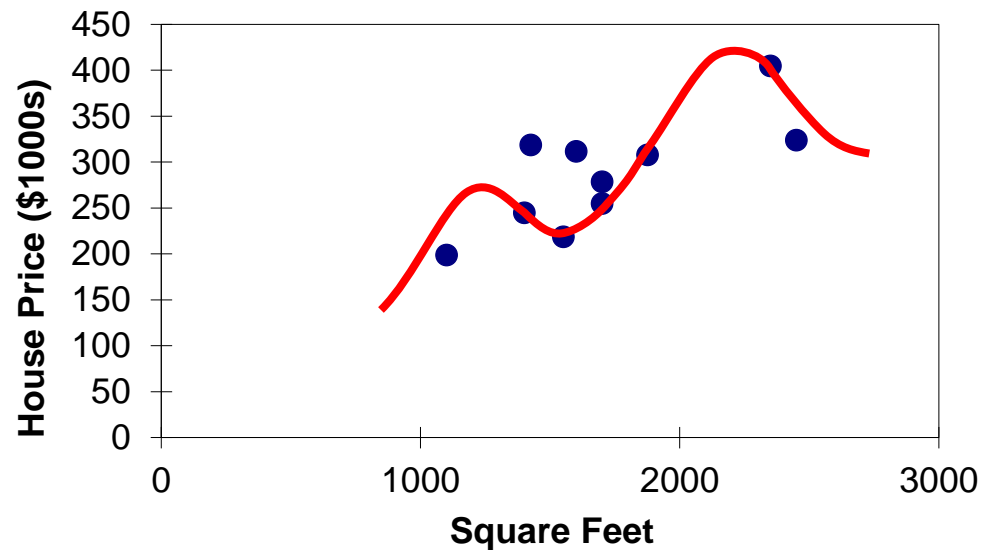
# Memorization?



- Store all sizes?
- Our data doesn't cover all sizes … What shall we do w a house of size 1750 sqft?
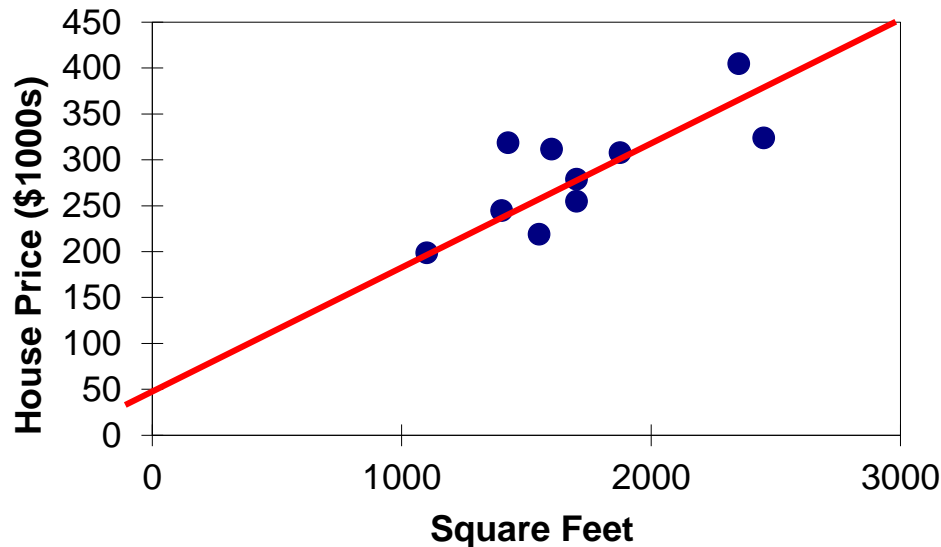
# OK … a function?

# Generalization: Learn a Function



## But which function y=f(x)?

# Simplest: Linear Model



- We assume/hypothesize that the relationship between the observed and the independent/explained variable is linear and thereby conduct our search

- This is our **Hypotheses Space –** all linear functions

# Linear Function Hypothesis

- How do we represent a hypothesis h in this space?

$$y = \theta_0 + \theta_1 x$$

- What if we have many features and not just house size?

- Such as:
  - Number of rooms
  - Distance to shopping center
  - Neighborhood crime rate
  - Distance to IDC
  - More…

9

# Multiple Features; Higher Dimension

- Let $$\mathbf{x^{(i)}} = \left( x_1^{(i)}, x_2^{(i)}, \ldots, x_n^{(i)} \right)$$

  be the vector of feature values for each instance $i$

- For simplicity we add another "constant" feature for every $i$

$$x_0^{(i)} = 1$$

- For $n$ features our linear hypothesis will be represented by the parameters
  $$\boldsymbol{\theta} = \left( \theta_0, \theta_1, \theta_2, \ldots, \theta_n \right)$$

  with which we construct:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- We say that $$\boldsymbol{\theta} = \left( \theta_0, \theta_1, \theta_2, \ldots, \theta_n \right)$$

  is the vector of parameters that defines our function (or our **model/execution-algorithm/hypothesis**)

## The Hypothesis (or model) is the execution algorithm

- How does a specific set of values,

$$\mathbf{\theta} = \left( \theta_0, \theta_1, \theta_2, \ldots, \theta_n \right)$$

  which defines a specific hypothesis (model) help us?

- How can we use it?

- Given a new house instance $\mathbf{x}' = \left( x'_0, x'_1, \ldots, x'_n \right)$
  we can estimate its price by an inner product with the vector $\mathbf{\theta}$ :

$$\text{value of house} = y = \theta_0 + \theta_1 x'_1 + \theta_2 x'_2 + \cdots + \theta_n x'_n$$

# Finding the Best Hypothesis/Model

- There are many possible parameter vectors

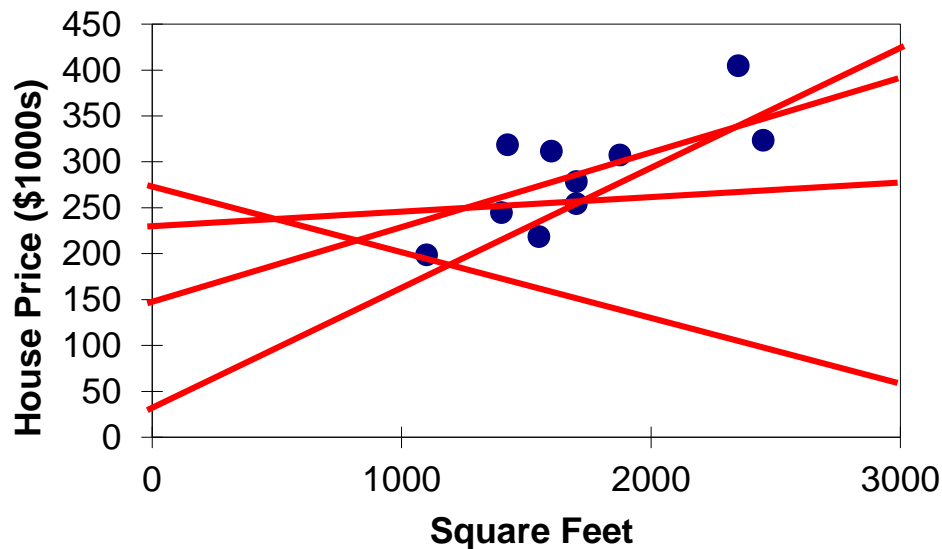$$\boldsymbol{\theta} = \left(\theta_0, \theta_1, \theta_2, \ldots, \theta_n\right)$$

  and each one defines a different hypothesis in our hypotheses space

- How do we find the best one?

- What can we use to help us find it?

- <u>The training set</u>: $m$ instances where, for each, we know the <u>feature values</u> $\mathbf{x^{(i)}} = \left(\mathrm{x}_0^{(i)}, \mathrm{x}_1^{(i)}, \ldots, \mathrm{x}_n^{(i)}\right)$

  as well as the <u>label value</u> $\mathrm{y}^{(i)}$

# The Hypotheses Space

- For the simple case of a single feature we get different possible straight lines when we change $\theta_0 \text{ and } \theta_1$ in the hypothesis (model)
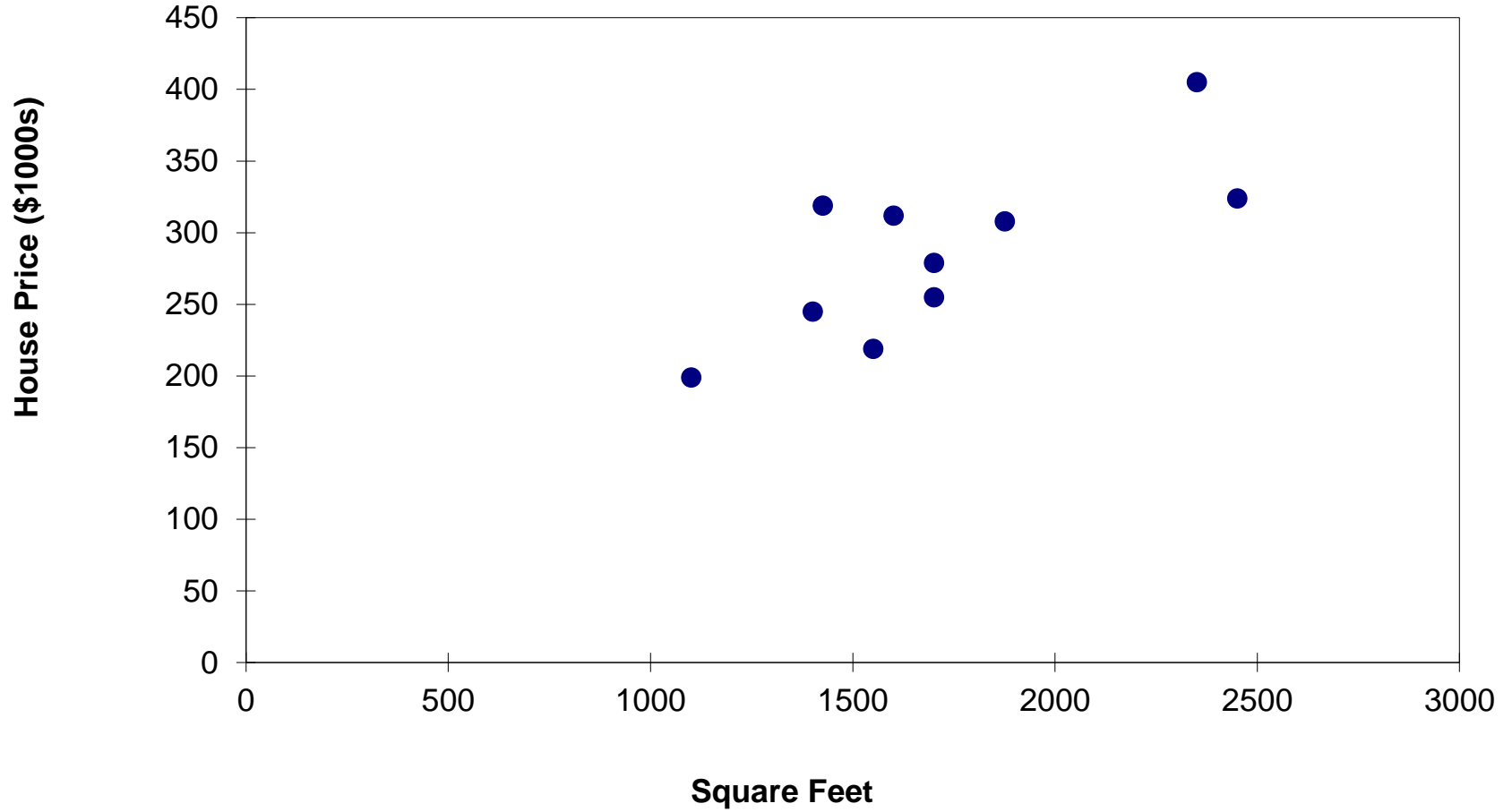
$$y = \theta_0 + \theta_1 x$$

# "Training" or "Learning". ERM

- We require that on our training data the values of our prediction function ($f$) would be similar to the known value of the house.

- So, we want to find $\boldsymbol{\theta}$ such that <u>for all</u> instances $i$ in the training set we will have:
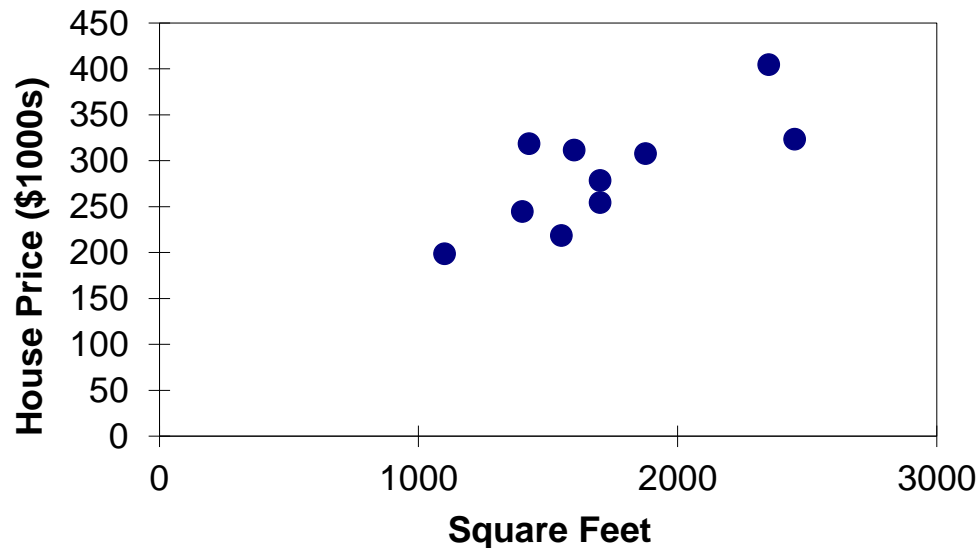
$$y^{(i)} = \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \cdots + \theta_n x_n^{(i)} = \boldsymbol{\theta} \cdot \mathbf{x^{(i)}}$$

- However, this may not always be possible – (why?)

# Consistent Learners

- A learning algorithm that can achieve 0 error on the training set is called a "consistent learner"
- This can not be done here.

# Cost Function of a Model θ

*Prediction*

*Actual value*

- As consistent learning may be impossible.
- Still, we can try to reduce the <u>error</u>.
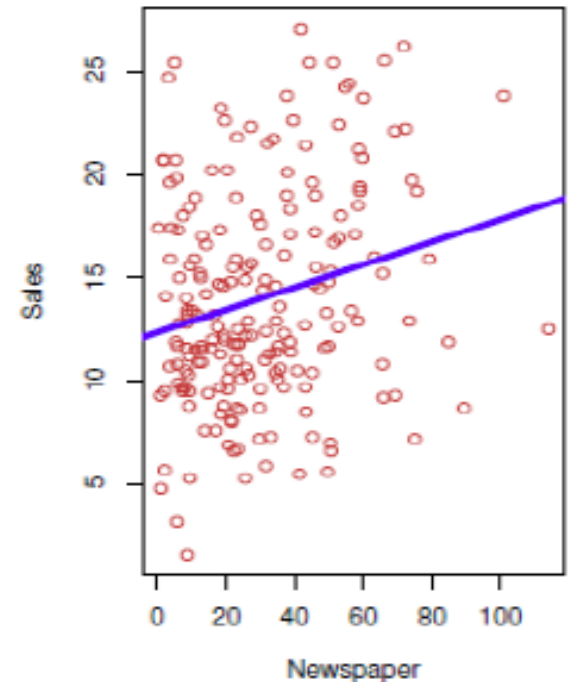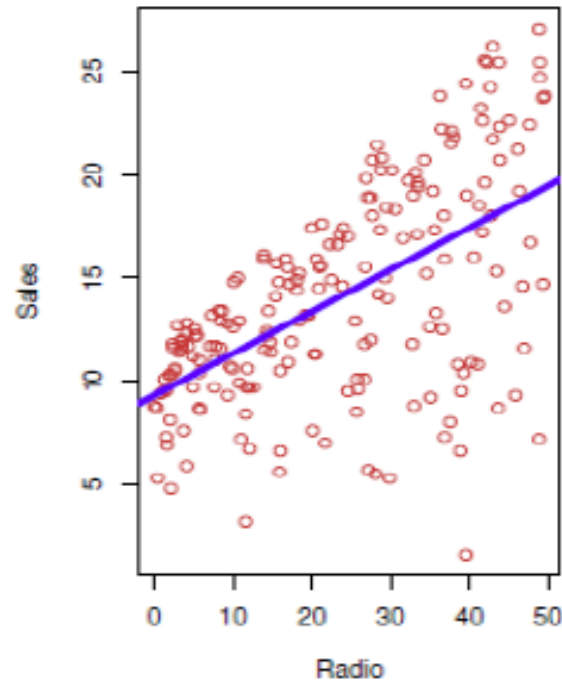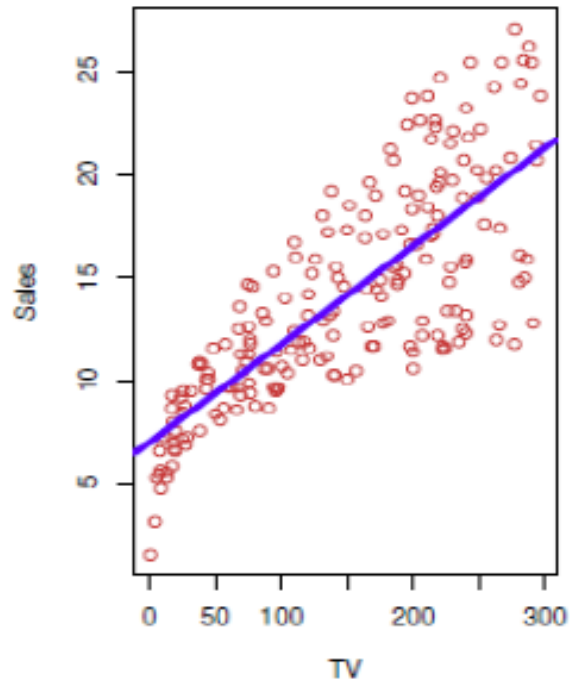  Per instance the error is:

$$(\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \cdots + \theta_n x_n^{(i)} - y^{(i)}) = \boxed{\boldsymbol{\theta} \cdot \mathbf{x}^{(i)}} - \boxed{y^{(i)}}$$

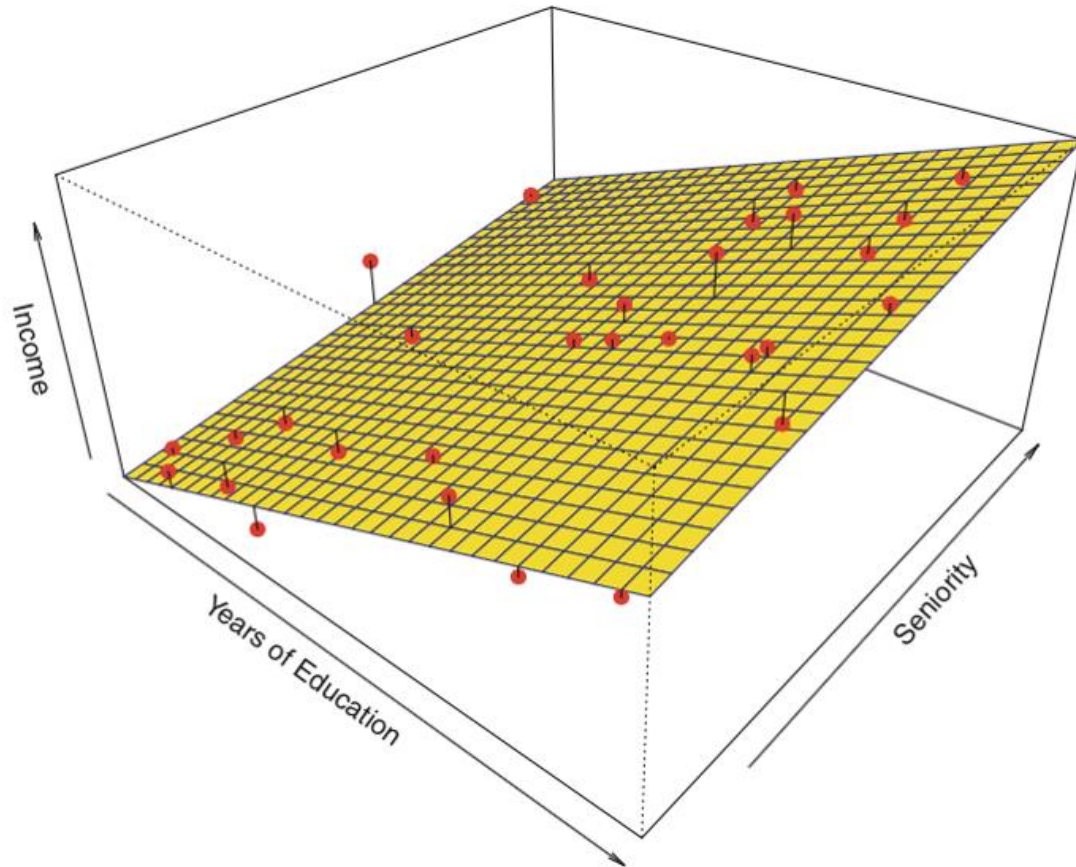- And we now average on <u>ALL</u> $m$ training instances to get our
  ***cost function***:

$$J(\boldsymbol{\theta}) = \frac{1}{2} \frac{1}{m} \sum_{i=1}^{m} \left( \boldsymbol{\theta} \cdot \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

- Square errors are used so that errors in different directions don't cancel out …
  Its also a smoother function compared to $|x|$

# Example: advertising and sales

# In higher dimensions
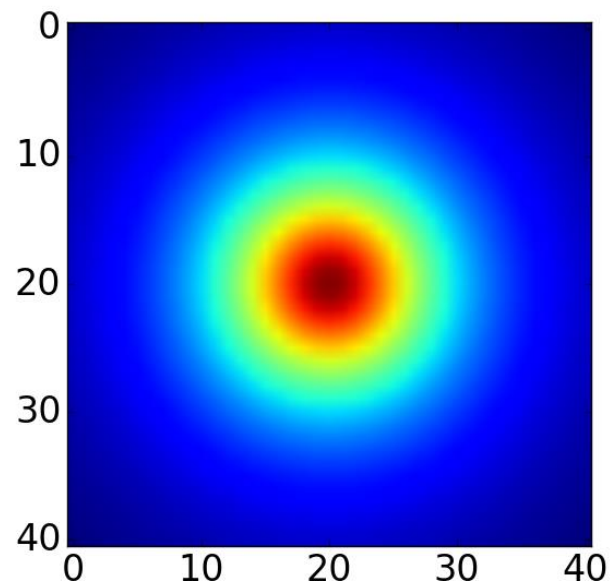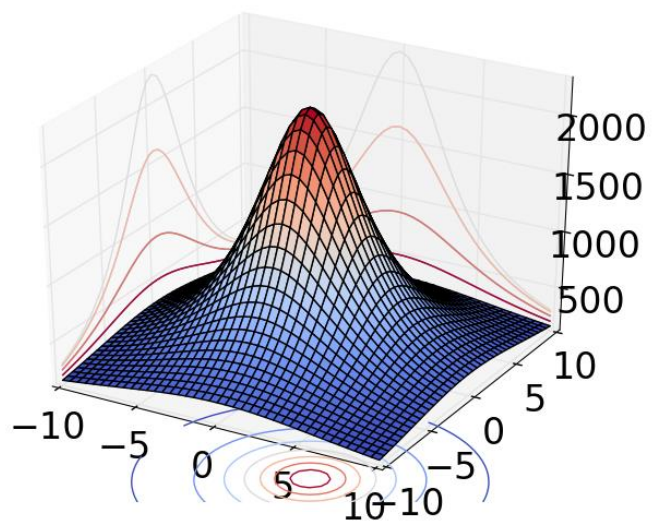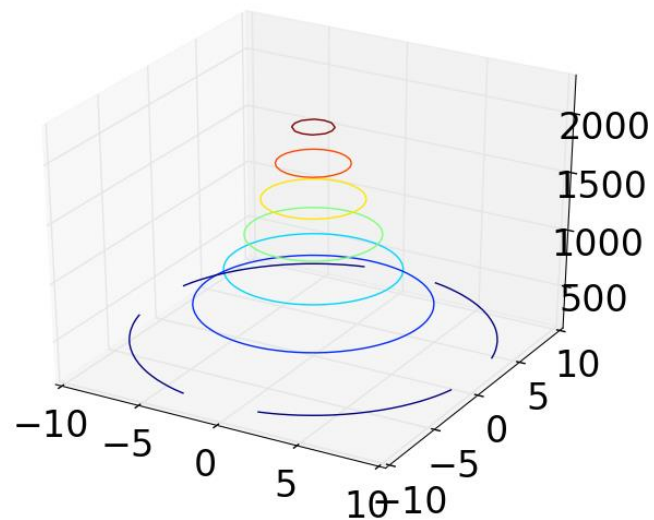
# Minimizing the Cost Function
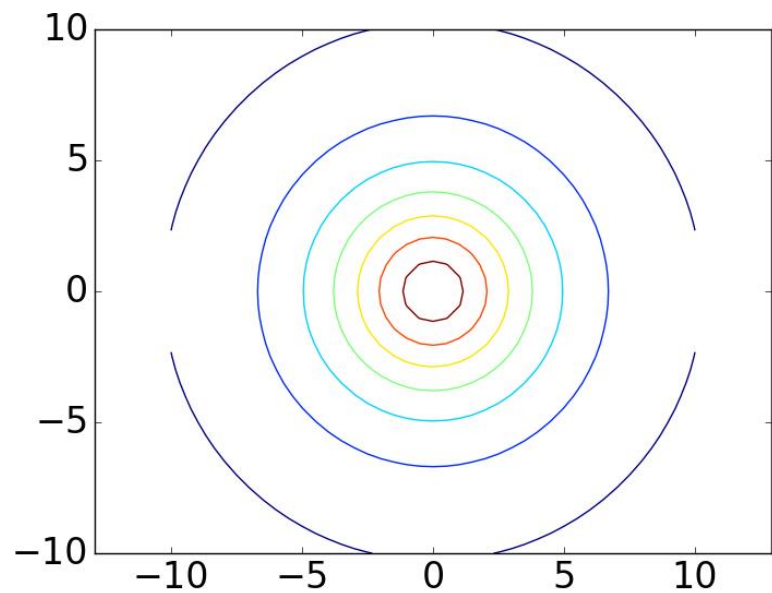
- In the ERM approach, our best hypothesis $\theta^*$ would be the one that <u>minimizes the cost function</u>. Formally:
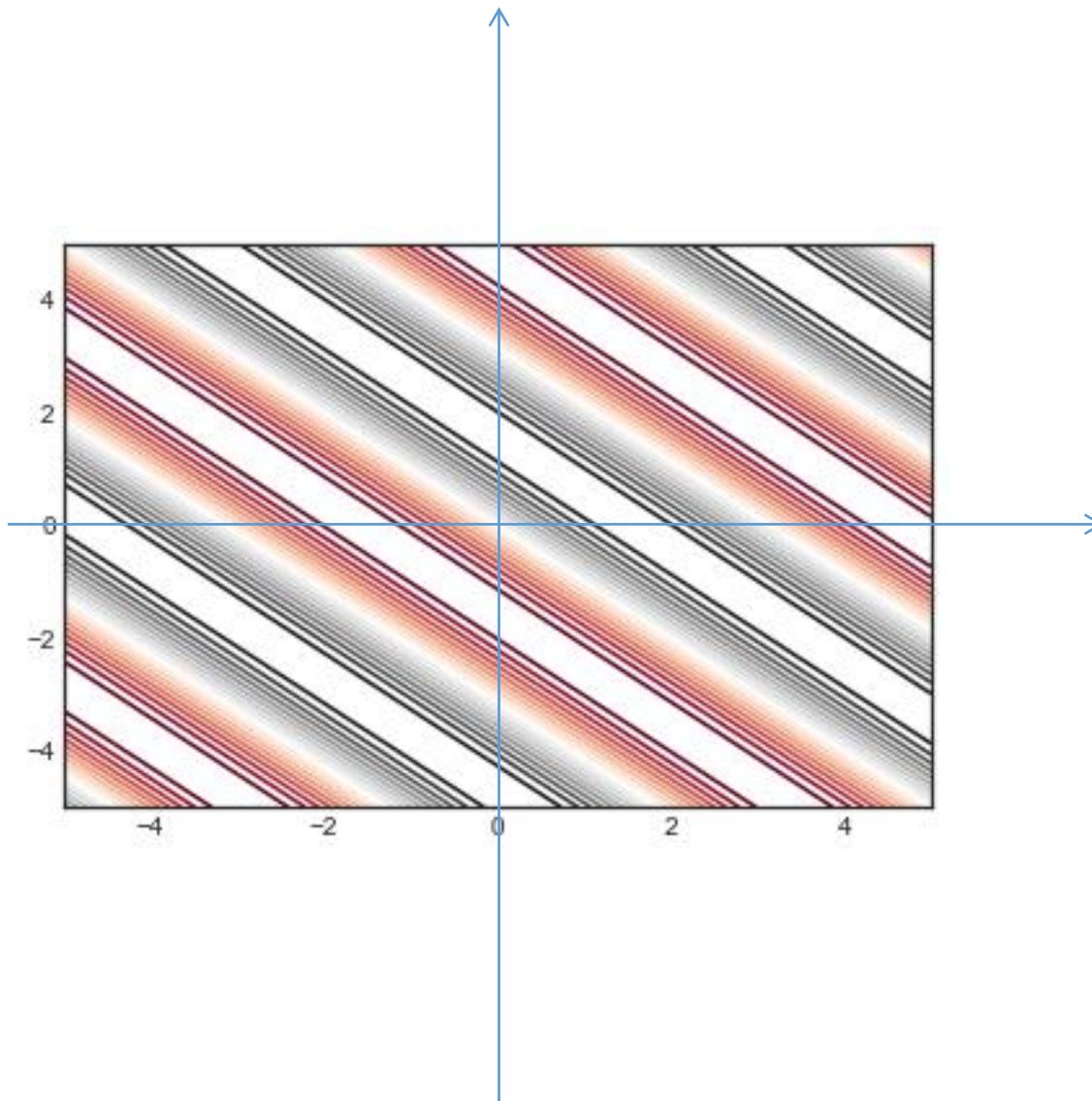
$$\theta^* = \arg\min_{\theta}\left[J(\theta)\right] = \arg\min_{\theta}\left[\frac{1}{2m}\sum_{i=1}^{m}\left(\theta \cdot \mathbf{x}^{(i)} - y^{(i)}\right)^2\right]$$

- How can we find it?

# The Cost Function

- In our simple case (house prices) we have 2 parameters: $\theta_0, \theta_1$

- For each value of these two parameters we can calculate the cost $J(\theta_0, \theta_1)$ over the entire training data (the training error).

- This is a function from $\mathbb{R}^2$ to $\mathbb{R}$ .

```python
def f(x, y):
    return np.sin(x+y)

x = np.linspace(-5, 5, 50)
y = np.linspace(-5, 5, 50)

mx, my = np.meshgrid(x, y)
z = f(mx, my)

plt.contour(mx, my, z, 20, cmap='RdGy');
```
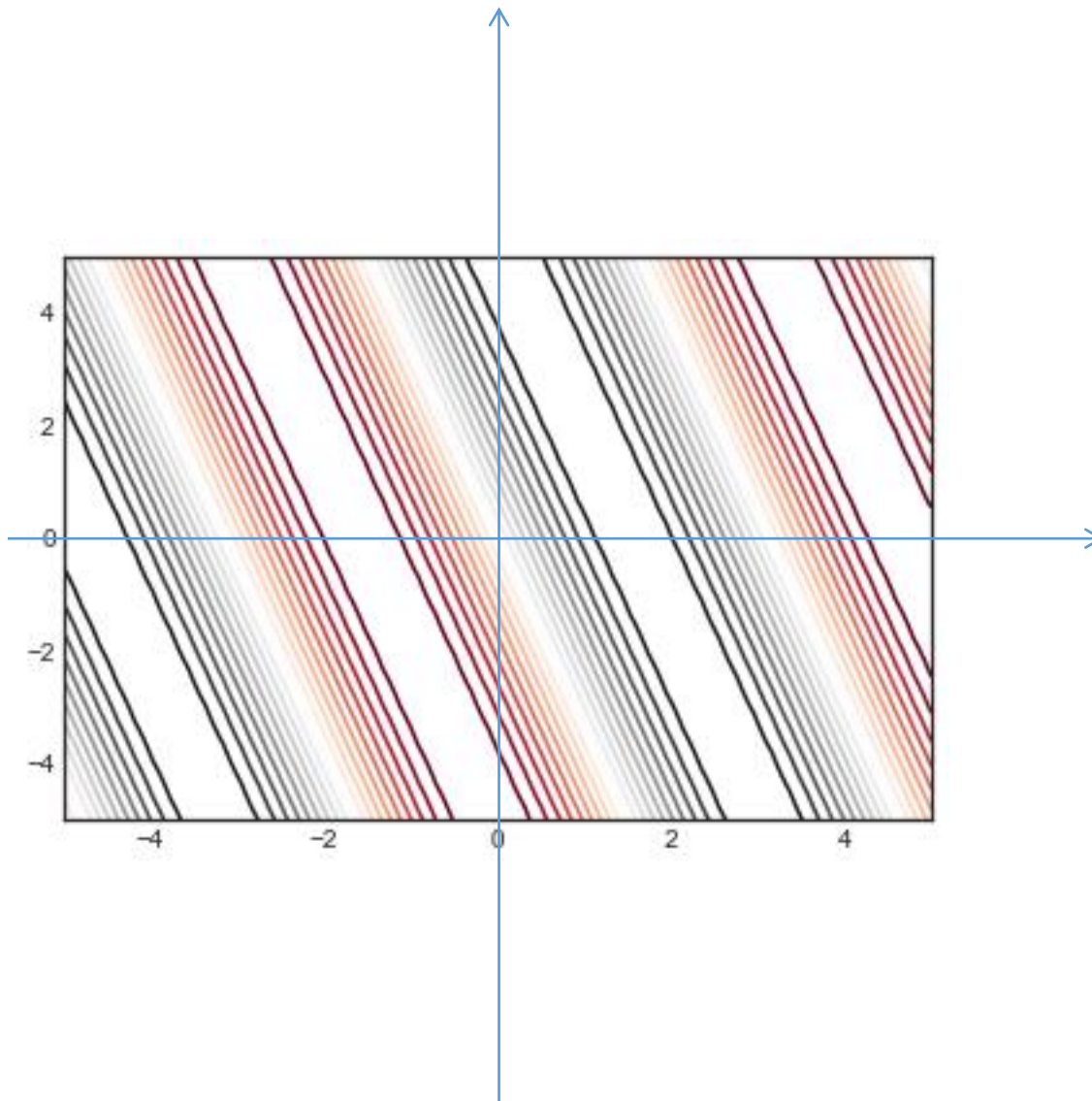
y

x

3D plot.
Not a contour plot

**IDC
HERZLIYA**

Recall:
We want to find the argmin of the cost function $J(\theta)$

# How to Find the Minimum of a 1D function?

- In high school:
  derivative at the point where an extremum is attained should be 0
- We can also follow the "downward" direction.
- How is this done?
  1. Find the derivative
  2. Move against its sign

# Directional Derivatives

Consider a differentiable 2D function $f(x_1, x_2)$ The derivative in a <u>general direction</u> $u = (u_1, u_2)$ (unit 2D vector) is called the directional derivative $D_u f$ and is defined as:
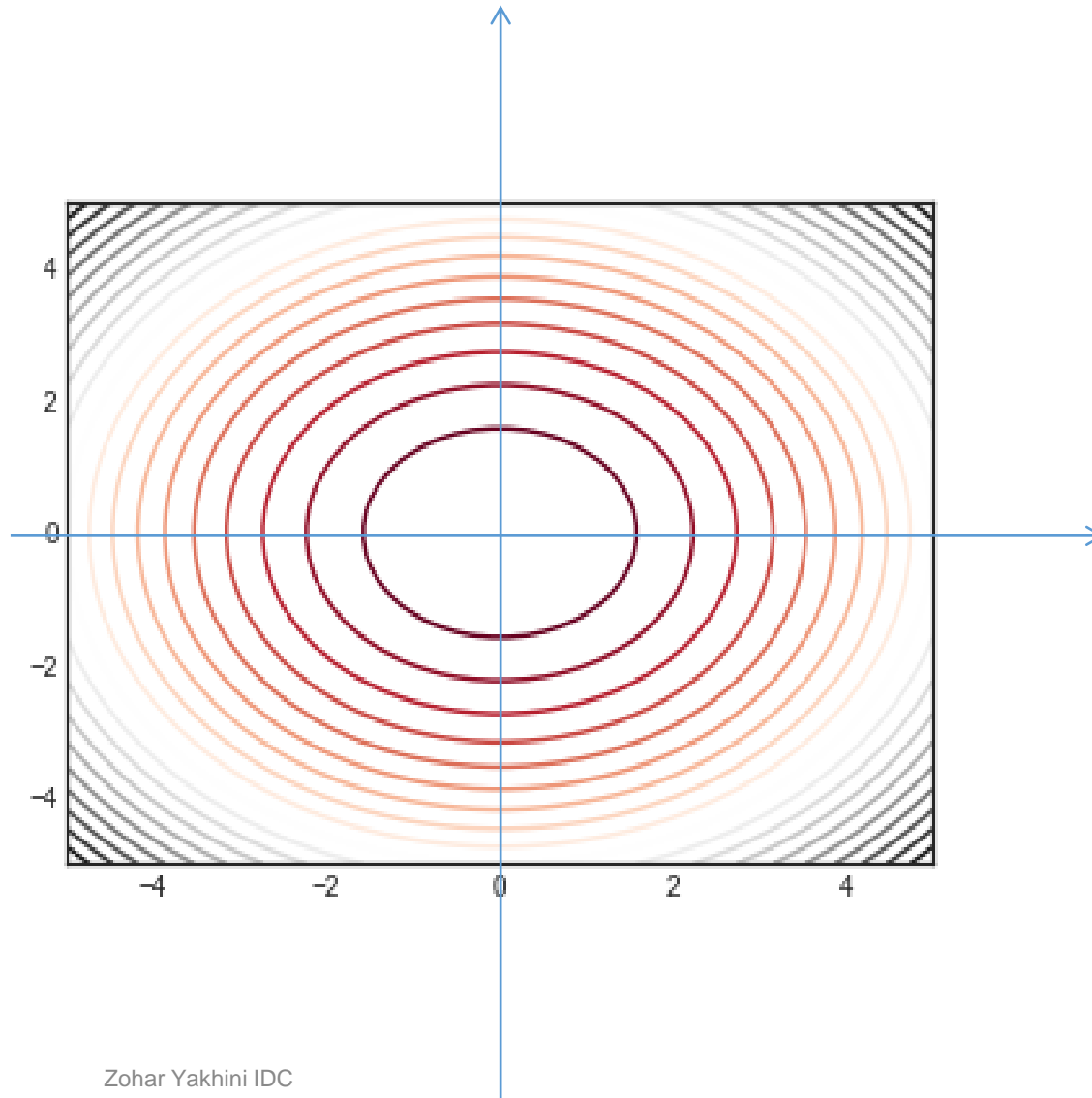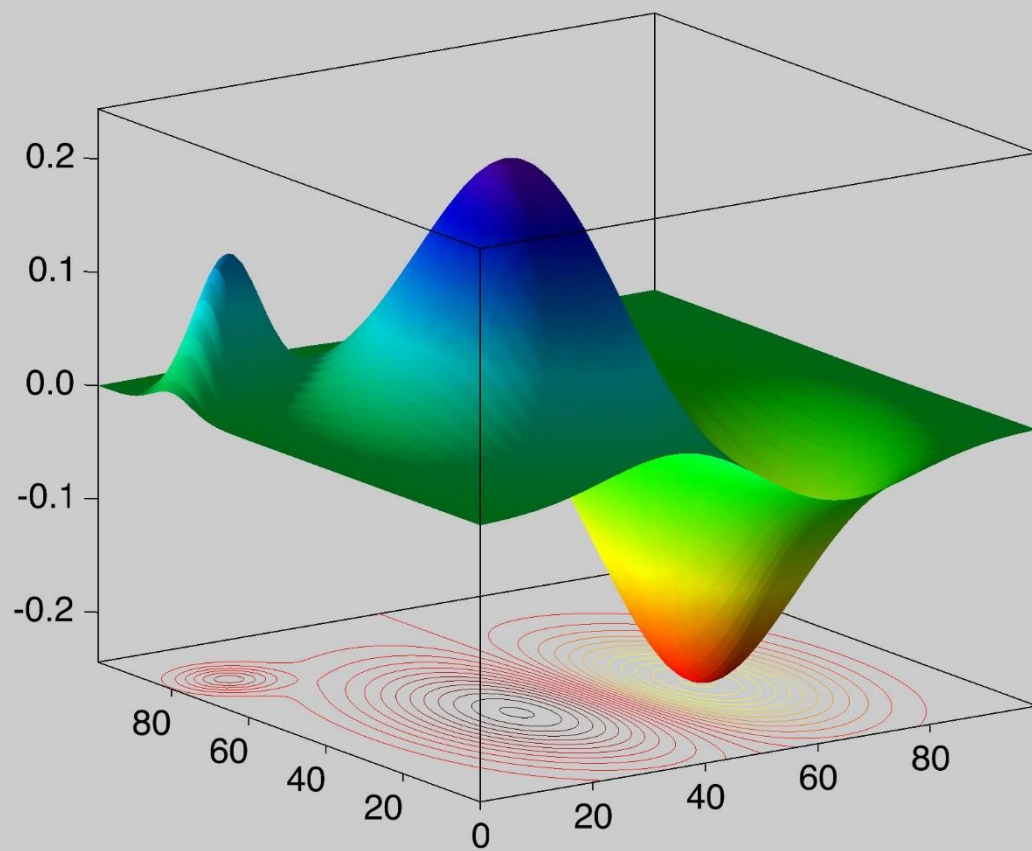
$$D_u f(x_1, x_2) = \lim_{s \to 0} \frac{f(x_1 + su_1, x_2 + su_2) - f(x_1, x_2)}{s} = \left( \frac{df}{ds} \right)_u$$

For a differentiable 2D function $f(x_1, x_2)$ the principal partial derivatives, those in direction $u = x_i$ are denoted

$$\frac{\partial f}{\partial x_1}(x_1, x_2) \text{ and } \frac{\partial f}{\partial x_2}(x_1, x_2)$$

# Partial Derivatives

- Again, for a differentiable 2D function $f(x, y)$, the (principal) partial derivatives in the directions $x$ and $y$ are denoted

$$\frac{\partial f}{\partial x}(x, y) \,, \frac{\partial f}{\partial y}(x, y)$$

- They can be computed by keeping one variable constant and differentiating by the other.
  For example ...

$$f(x, y) = 2x + 13y$$
$$f(x, y) = y \exp(x)$$
$$f(x, y) = 2x + 3xy + 5y^2$$

https://mathinsight.org/directional_derivative_gradient_introduction

$$f(x, y) = y\, e^x$$

$$\frac{\partial f}{\partial x}(x, y) = ye^x$$

$$\frac{\partial f}{\partial y}(x, y) = e^x$$

$$f(x, y) = y^4 e^x$$

$$\frac{\partial f}{\partial x}(x, y) = y^4 e^x$$

$$\frac{\partial f}{\partial y}(x, y) = 4y^3 e^x$$

# Extrema and zero derivatives



1D

2D

# The Gradient of a function

Define the
**GRADIENT of $f$:**

$$\nabla f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)$$

Thm:

For any direction $u = (u_1, u_2)$ and any point $\vec{x} = (x_1, x_2)$ we have:

$$D_u f(\vec{x}) = \langle \nabla f(\vec{x}), u \rangle$$

# Most Rapid Increase at a point $\vec{x}$ ?



The pt $\vec{x}$

## Most Rapid Increase at a point $\vec{x}$

- The directional derivative in the direction of the vector $u = (u_1, u_2)$ (a scalar!) can also be written as:

$$D_u f(x_1, x_2) = \nabla f \cdot \vec{\mathbf{u}} = \left|\nabla f\right|\left|\vec{\mathbf{u}}\right| \cos \beta = \left|\nabla f\right| \cos \beta$$

- Where β is the angle between $u$ and $\nabla f$

- However, $\cos \beta \leq 1$.

- Therefore:
  1. The greatest increase in the function happens in the direction of the gradient (i.e. β =0)
  2. The greatest decrease is in the direction $-\nabla f$ (i.e β =180⁰)

# Steepest ascent and Iso-Contours

- Steepest ascent follows the gradient direction

- Using same argument, when the direction $u$ is perpendicular to the gradient, $\cos\beta$ =0 and there is no change in the function – this is exactly the direction of the iso-contours, walking with no change in altitude.



$\nabla f$
$\nabla f$
$\nabla f$
$\nabla f$

ERWIN R. AISZ

# The gradient of n-dimensional functions

- For an n-D function $y = f(x_1, x_2, \ldots, x_n)$ the gradient is defined as :

$$\nabla f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \ldots, \frac{\partial f}{\partial x_n} \right)$$

- <u>The gradient vector at any given point is the direction of greatest increase of $f$ at this point.</u>

- The direction of greatest decrease of $f$ is opposite to the gradient: $-\nabla f(x_1, x_2 \ldots x_n)$

## How to find a minimum of a (reasonably smooth) n-dimensional function?

- Take steps in the direction of maximum decrease – i.e. opposite the gradient direction!

- This is called **Gradient Descent**

# Gradient Descent Steps

- For a small enough $\alpha > 0$ we have
$f(x - \alpha \nabla f(x)) < f(x)$ ,
so we take one step of size $\alpha$ in the opposite direction of the gradient.

- This represents greedy local descent

- Note:
such steps can zig-zag

# Back to Minimizing the Cost Function in Linear Regression

- Recall that our best model, $\theta^*$, is the one that minimizes the cost function. That is:

$$\boldsymbol{\theta}^* \; attains \; \min_{\boldsymbol{\theta}}\big[\mathrm{J}(\boldsymbol{\theta})\big] = \min_{\boldsymbol{\theta}}\left[ \tfrac{1}{2m} \sum_{i=1}^{m} \left( \boldsymbol{\theta} \cdot \mathbf{x}^{(i)} - \mathrm{y}^{(i)} \right)^2 \right]$$

- $J(\vec{\theta})$ is a real valued function of $\vec{\boldsymbol{\theta}} \in \mathbb{R}^n$

- So - we can start with some initial guess $\vec{\boldsymbol{\theta}_0^*}$ and then use gradient descent

IDC
HERZLIYA

# Gradient Descent Algorithm

- Start with some value $\theta(0) = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$
- Repeat until you reach a minimum:

  For all $j = 0 \dots n,$

  Update $\theta_j\,(t+1) \leftarrow \ \theta_j(t) - \alpha \frac{\partial}{\partial \theta_j} J(\theta(t))$

- $\alpha > 0$ is a parameter of the algorithm called the <u>learning rate</u>
- Updates are simultaneous (in all $n+1$ directions)
- In the general case this process can still be trapped in local minima!

IDC
HERZLIYA

# Minimizing the Cost Function

- This gradient descent process is seeking:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}}\left[J(\boldsymbol{\theta})\right] = \arg\min_{\boldsymbol{\theta}}\left[\frac{1}{2m}\sum_{i=1}^{m}\left(\boldsymbol{\theta}\cdot\mathbf{x}^{(i)} - y^{(i)}\right)^2\right]$$

- Missing details?

# Calculating the Partial Derivatives

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1 \ldots, \theta_n) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \frac{\partial}{\partial \theta_j} (\theta_0 + \theta_1 x_1^{(i)} + \cdots + \theta_j x_j^{(i)} + \cdots \theta_n x_n^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} 2(\theta_0 + \theta_1 x_1^{(i)} + \cdots + \theta_j x_j^{(i)} + \cdots \theta_n x_n^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

IDC
HERZLIYA

# Gradient Descent for Linear Regression

- Initialize $\mathbf{\theta} = \left( \theta_0, \theta_1, \theta_2, \ldots, \theta_n \right)$
- Repeat until you reach a minimum (or stop cdn):
  - For all $0 \leq j \leq n$,

$$\theta_j\,(t+1) = \theta_j(t) - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_{\theta(t)}\left(x^{(i)}\right) - y^{(i)} \right) \cdot x_j^{(i)}$$

- In words: set the new $\theta_j$ to the current $\theta_j$ minus the learning rate (α) times the partial derivative of the error function with respect to $\theta_j$, computed at the current θ.

- Also remember that $x_0^{(i)} = 1$

# Closed form solution?

$$X \cdot \theta = y$$

$$
m \begin{pmatrix}
x_0^{(1)} & x_1^{(1)} & & x_n^{(1)} \\
x_0^{(2)} & x_1^{(2)} & \cdots & x_n^{(2)} \\
& & & \\
& \vdots & \ddots & \vdots \\
& & & \\
x_0^{(m)} & x_1^{(m)} & \cdots & x_n^{(m)}
\end{pmatrix}
\begin{pmatrix}
\theta_0 \\
\vdots \\
\theta_n
\end{pmatrix}
=
\begin{pmatrix}
y^{(1)} \\
y^{(2)} \\
\\
\vdots \\
\\
y^{(m)}
\end{pmatrix}
$$

$n+1$  $\quad \theta \in \mathbb{R}^{n+1} \quad Y \in \mathbb{R}^m$

- If $X$ is square and non singular we can write $\vec{\boldsymbol{\theta}} = X^{-1}y$
- In ML typically $X$ is overdetermined $(m \gg n)$.
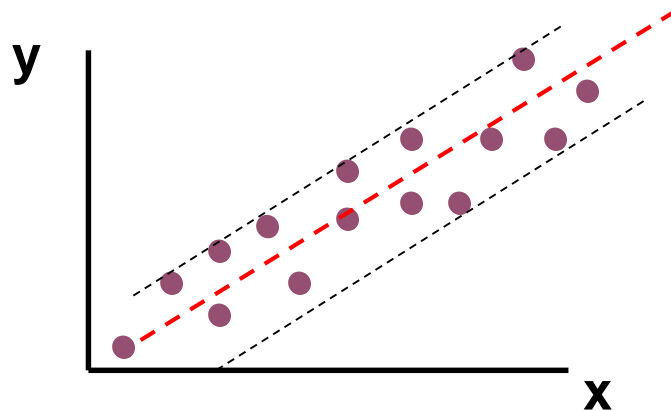
IDC
HERZLIYA

## Recitation

- Description of HWA1
- Feature scaling
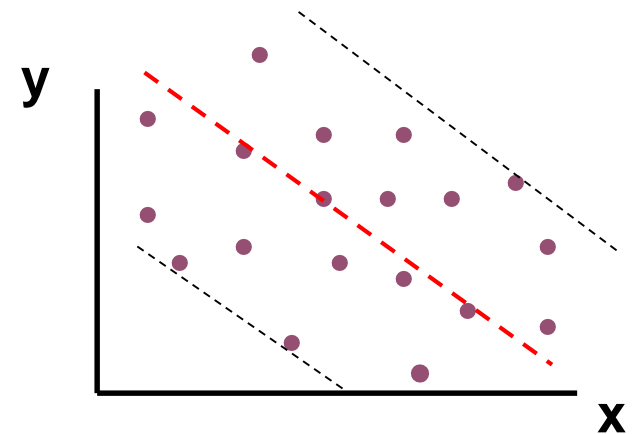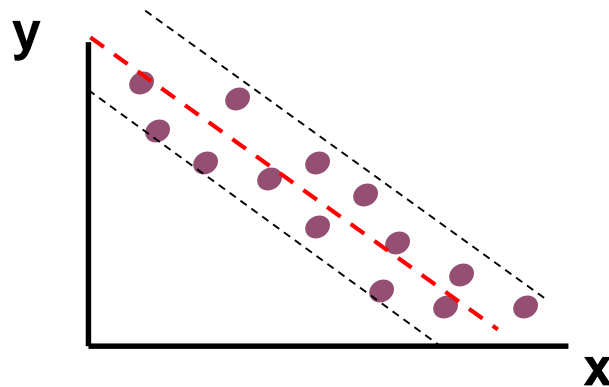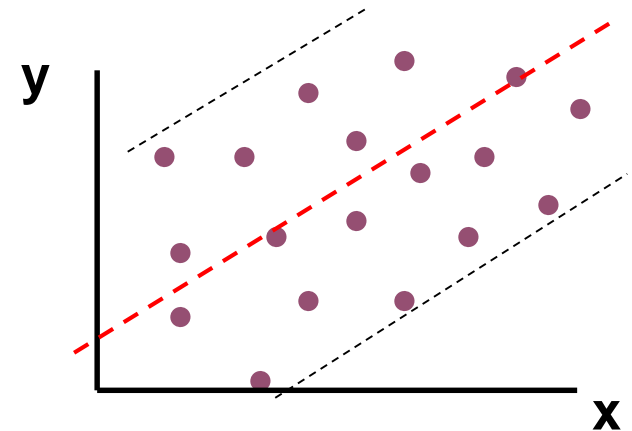- Learning rate
- Python

## Extra slides/notes

- Correlations
- Pseudo-inverse of a matrix (closed form solution)
- Going beyond Linear (also in the HW)

# Strong vs. Weak Linear Relationship
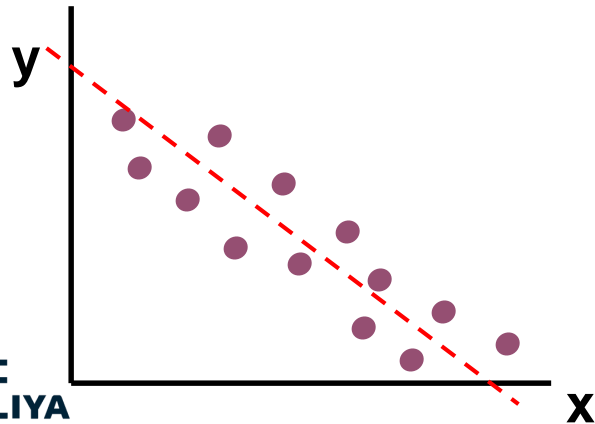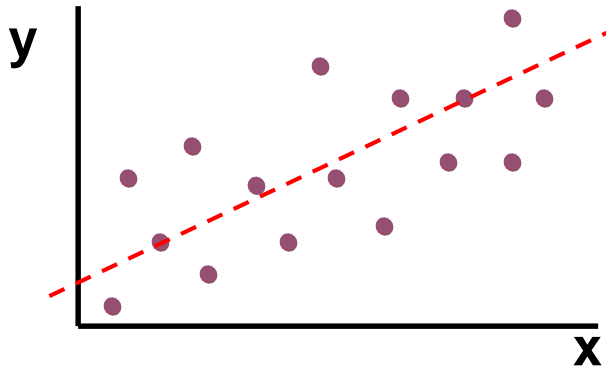
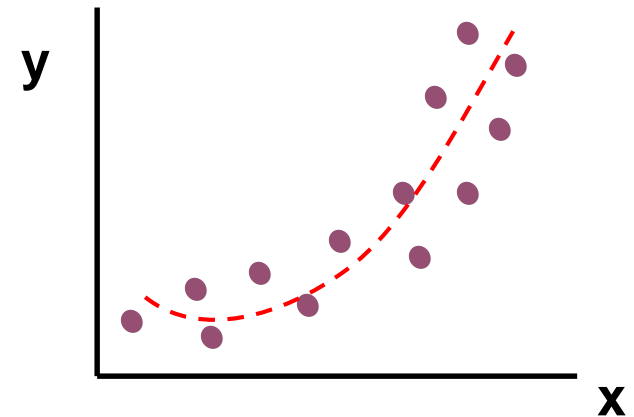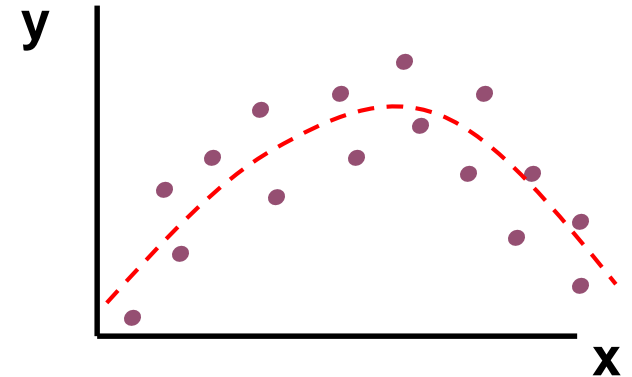**Strong relationships**

**Weak relationships**

# Not Everything is Linear!

**Linear relationships**

**Curvilinear relationships**

# No Evident Relationship

# Correlation Analysis

- Correlation analysis is used to measure strength of the association (**linear relationship**) between two variables

  - Only concerned with strength of the relationship

  - No causal effect is implied

- The sample (Pearson) correlation coefficient **r** is a measure of the strength of the linear relationship between two variables, based on sample observations

# The Pearson Correlation Coefficient

$$r = \frac{\sum_{i=1}^{m}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{m}(x_i - \bar{x})^2 \sum_{i=1}^{m}(y_i - \bar{y})^2}} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

Where:

r  =  Sample Pearson correlation coefficient

m =  Number of samples

x  =  Value of an explaining variable

y  =  Value of the dependent variable
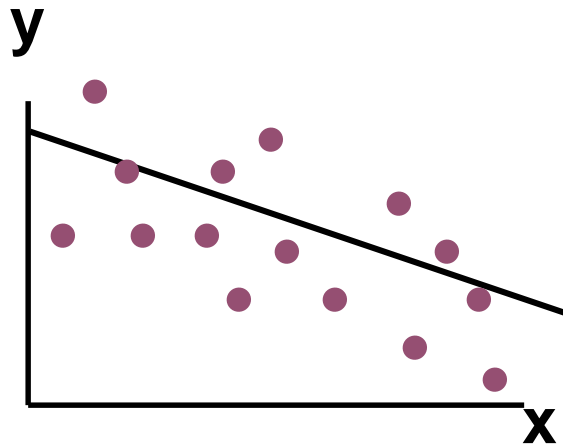
# Features of  r

- Unit free
- Ranges between -1 and 1
- The closer to 0, the weaker the <u>linear</u> relationship
- The closer to -1, the stronger the negative <u>linear</u> relationship
- The closer to 1, the stronger the positive <u>linear</u> relationship

# Examples of r Values



r = -1

r = -.6

r = 0

r = +.3

r = +1

# More Complex Relationships May Not Be Captured Using Pearson r

# Polynomial Regression

- We can expand our feature space by using functions of the original features.

- For example, if we want to use a cubic function feature space we can define:

$$x_0 = 1, x_1 = x, x_2 = x^2, x_3 = x^3$$

then use regular regression and in essence we are learning the function

$$\theta_0 + \theta_1 x + \theta_1 x^2 + \theta_3 x^3$$

# Other Functions

- We can use other functions of features such as $x^{1/q}$

- We can define functions of sets of variables such as $x_1 x_5 x_7$ or $x_1^2 x_5 x_7^3$

- All of these can be represented as new features and increase the dimension of the entire calculation

# Closed form solution?

$$X \cdot \theta = y$$

$$
m \left.
\begin{pmatrix}
x_0^{(1)} & x_1^{(1)} & & x_n^{(1)} \\
x_0^{(2)} & x_1^{(2)} & \cdots & x_n^{(2)} \\
& & & \\
\vdots & & \ddots & \vdots \\
& & & \\
x_0^{(m)} & x_1^{(m)} & \cdots & x_n^{(m)}
\end{pmatrix}
\right.
\begin{pmatrix}
\theta_0 \\
\vdots \\
\theta_n
\end{pmatrix}
=
\begin{pmatrix}
y^{(1)} \\
y^{(2)} \\
\\
\vdots \\
\\
y^{(m)}
\end{pmatrix}
$$

$n+1$     $\theta \in \mathbb{R}^{n+1}$    $Y \in \mathbb{R}^m$

- If $X$ is square and non singular we can write $\vec{\boldsymbol{\theta}} = X^{-1} y$
- Most of the time $X$ is overdetermined ($m \gg n$).

IDC
HERZLIYA

# Schematic set-up

$$n+1 \qquad \theta \in \mathbb{R}^{n+1} \quad y \in \mathbb{R}^m$$

$$m \quad X \qquad \sim$$

# Minimum when $\nabla = 0$

- Remember the (MSE) cost function:

$$E(\boldsymbol{\theta}) = \left\| X \cdot \vec{\boldsymbol{\theta}} - y \right\|_2^2 = \sum_{i=1}^{m} \left( x^{(i)} \theta^T - y^{(i)} \right)^2$$

- This will be minimal when the gradient $\nabla E$ is 0.
- Forming the derivatives yields (see below):

$$\nabla E(\boldsymbol{\theta}) = \sum_{i=1}^{m} 2 \left( x^{(i)} \boldsymbol{\theta}^T - y^{(i)} \right) x^{(i)} = 2 X^T (X\boldsymbol{\theta} - y)$$

# Proof

$$\nabla E(\theta) = 0, \quad \nabla E(\theta) = \left(\frac{\partial E}{\partial \theta_j}(\theta)\right)_{j=1}^{n}$$

Expanding the partial derivatives, we get the following expression (for setting to 0)

$$\left(\nabla E(\theta)\right)_j = \frac{\partial E}{\partial \theta_j}(\theta) = \frac{\partial}{\partial \theta_j}\left(\sum_{i=1}^{m}\left(X(i,:)\cdot\theta - y(i)\right)^2\right)$$

Recall that $X(i,:)$ is the notation for the $i^{th}$ row. Now we expand the dot product

$$\frac{\partial E}{\partial \theta_j}(\theta) = \frac{\partial}{\partial \theta_j}\left(\sum_{i=1}^{m}\left(\left(\sum_{k=1}^{n}X(i,k)\cdot\theta_k\right) - y(i)\right)^2\right)$$

$$\frac{\partial E}{\partial \theta_j}(\theta) = \left(\sum_{i=1}^{m}\frac{\partial}{\partial \theta_j}\left(\left(\sum_{k=1}^{n}X(i,k)\cdot\theta_k\right) - y(i)\right)^2\right)$$

$$\frac{\partial E}{\partial \theta_j}(\theta) = \left(\sum_{i=1}^{m}2\left(\left(\sum_{k=1}^{n}X(i,k)\cdot\theta_k\right) - y(i)\right)\cdot X(i,j)\right)$$

# Proof, cont

For all $j \neq k$ the derivative is 0, so we only consider $\theta_j$
$$\nabla E(\theta) = \left( 2 \sum_{i=1}^{m} \big( X(i,:) \cdot \theta - y(i) \big) \cdot X(i,j) \right)$$
$$\nabla E(\theta) = (2X(:,j)(X \cdot \theta - y))$$
Notice that $\big( X(:,j) \big)^T = X^T(j,:)$ and therefore
$$\nabla E(\theta) = (2X^T(j,:)(X \cdot \theta - y))$$
$$\nabla E(\theta) = (2X^T(X \cdot \theta - y))$$

Recall that we want to solve for the $\theta$ that makes this 0:

$$0 = 2X^T(X \cdot \theta - y)$$
$$0 = (X^T X \cdot \theta - X^T y)$$
$$X^T X \cdot \theta = X^T y$$
$$\theta = (X^T X)^{-1} X^T y$$

- Setting the gradient of E to zero yields the necessary condition for minimum (see notes and slides above): $X^T X \cdot \theta = X^T y$

- Now, $\boldsymbol{X^T X}$ is square and often nonsingular and so we can solve for $\theta$ uniquely as:

$$\theta = pinv(X)y \qquad \text{where} \qquad pinv(X) = (X^T X)^{-1} X^T$$

- The n×m matrix $\boldsymbol{pinv(X)} = (\boldsymbol{X^T X})^{-1} \boldsymbol{X^T}$ is called the **pseudo inverse** of **X** (which is mxn)
- If $\boldsymbol{X}$ is square it is just its inverse.

# Gradient Descent and Pseudo Inverse

- $X^T X$ can still be singular (or not full rank) and not have an inverse. This can be resolved with some more algebra.

- This pinv technique doesn't work for all error functions J. Gradient descent is more general.

- Gradinet descent allows parallelization.

# Summary

- Regression learns a function to predict values based on a vector of features

- Linear Regression uses a Gradient Descent Algorithm or a Pseudo-Inverse solution

- Gradient descent is a general minimization procedure

- There are other (non-linear) relations between variables

- The Pearson correlation coefficient is a measure of the strength of the linear relationship between two variables

- Pseudo inverse solution

# Summary - cont

The execution algorithm:

$$y = f(x) = x\theta^T$$



$X$ (instances by features)
$y$ (instances)

The learning algorithm (regression)

x

θ

y