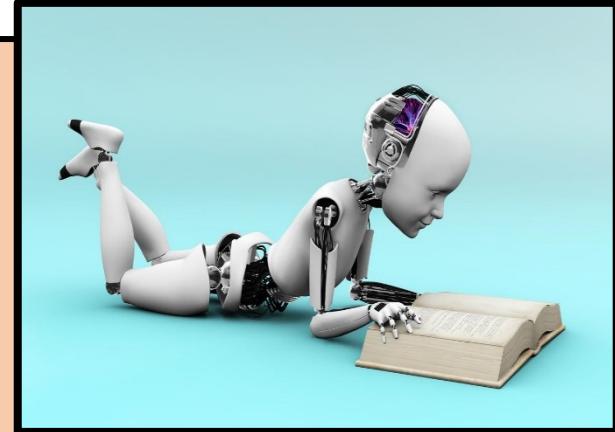


Machine learning

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 feature values and labels

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

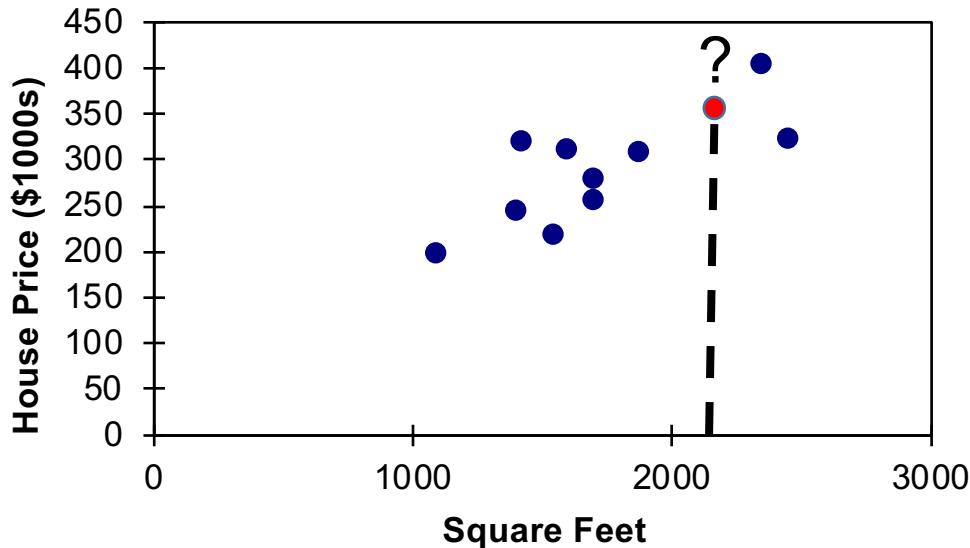
Explaining variable (x) = square footage

Graphical Representation

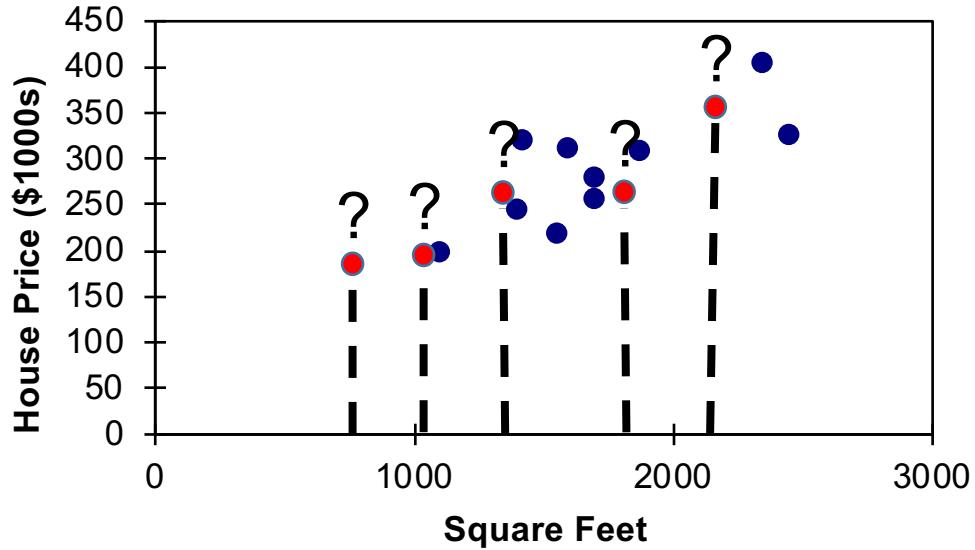
Scatter plot of
House Price (y) vs
House Size (x)

Prediction:

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

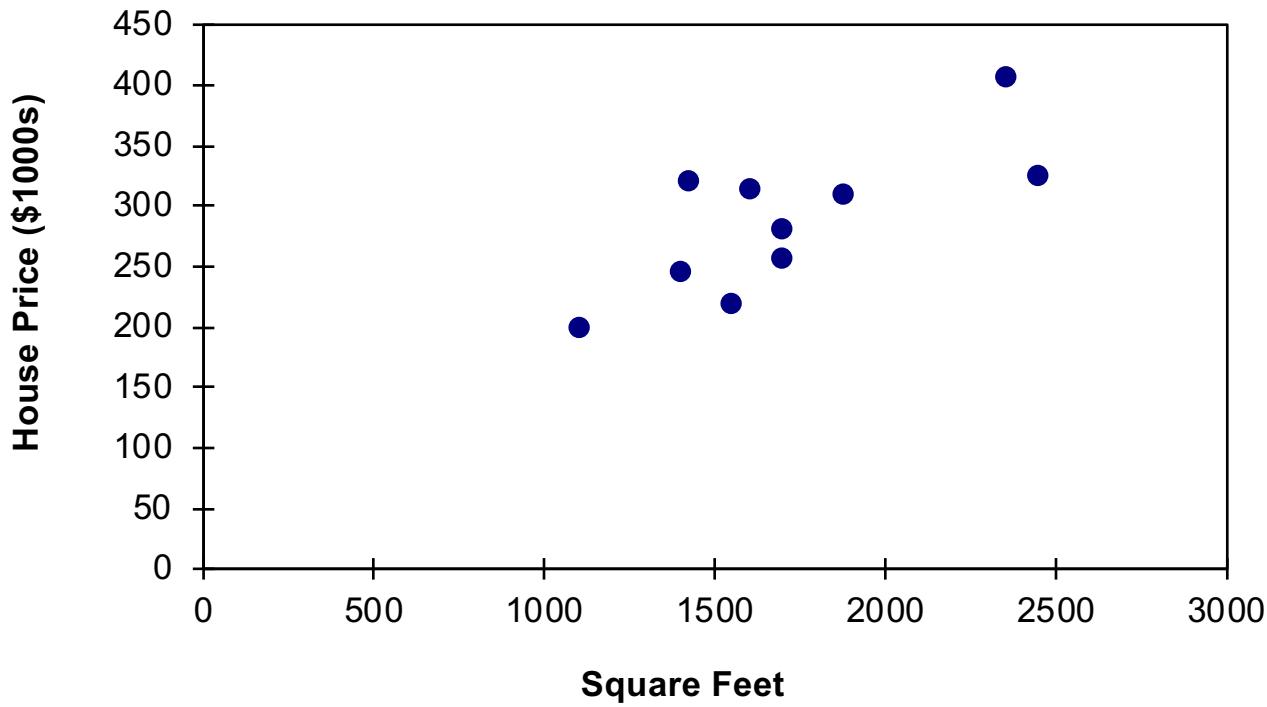


Memorization?

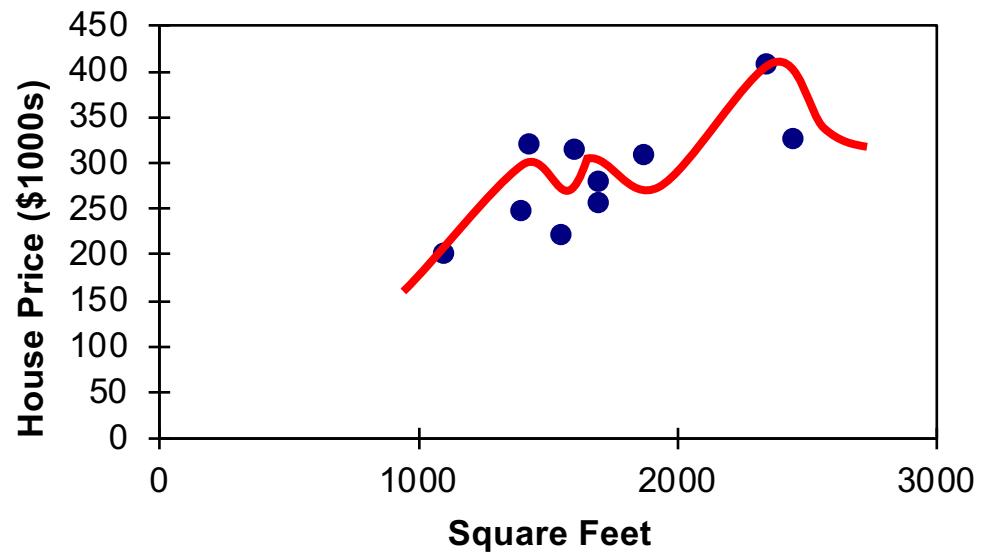


- Store all sizes?
- Our data doesn't cover all sizes ... What shall we do w a house of size 1750 sqft?

OK ... a function?

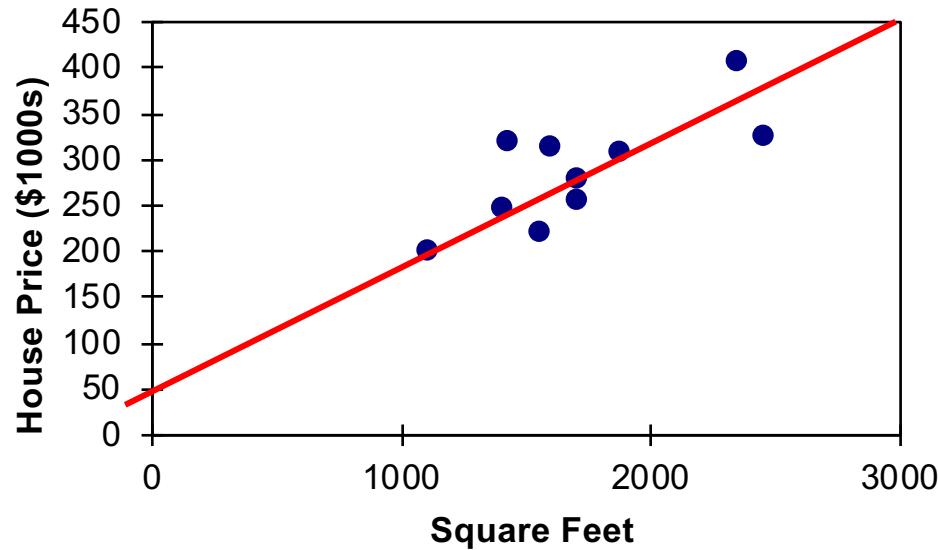


Generalization: Learn a Function



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...

Multiple Features

- Let $\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$ 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 $\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

for which we have:

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

- We say that $\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$ is the vector of parameters that defines our function (or our **model**)

The Hypothesis (or model) is the execution algorithm

- How does a specific set of values, $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$, which defines a specific hypothesis (model) help us?
- How can we use it?
- Given a new house instance $x' = (x'_0, x'_1, \dots, x'_n)$ we can estimate its price by an inner product with the vector Θ :

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

Finding the Best Hypothesis/Model

- There are many possible parameter vectors

$$\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$$

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?
-
- The training set: m instances where, for each, we know the feature values $\mathbf{x}^{(i)} = (x_0^{(i)}, x_1^{(i)}, \dots, x_n^{(i)})$

as well as the value

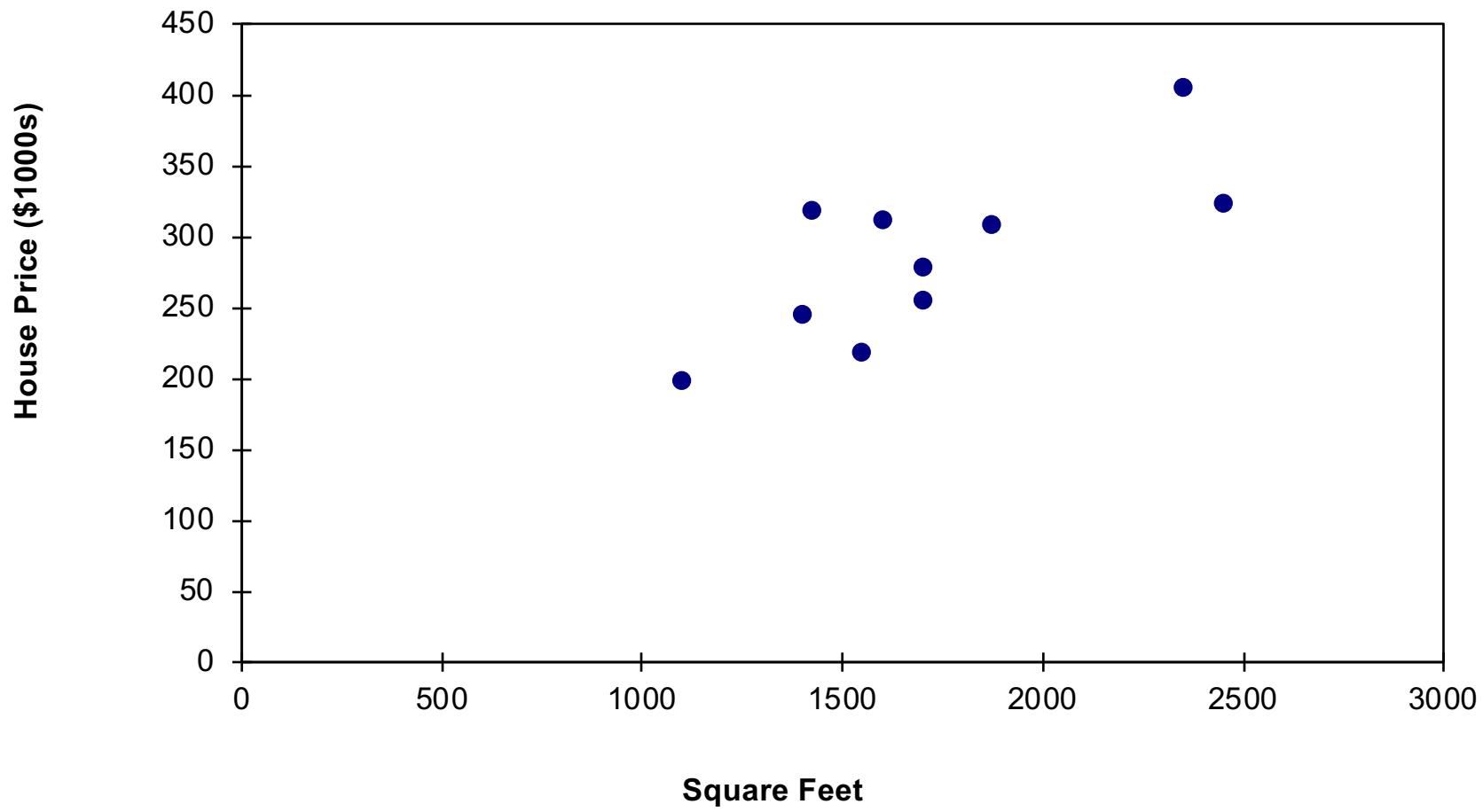
$$y^{(i)}$$

“Training” or “Learning”

- 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 θ such that for all instances i in the training set we will have:

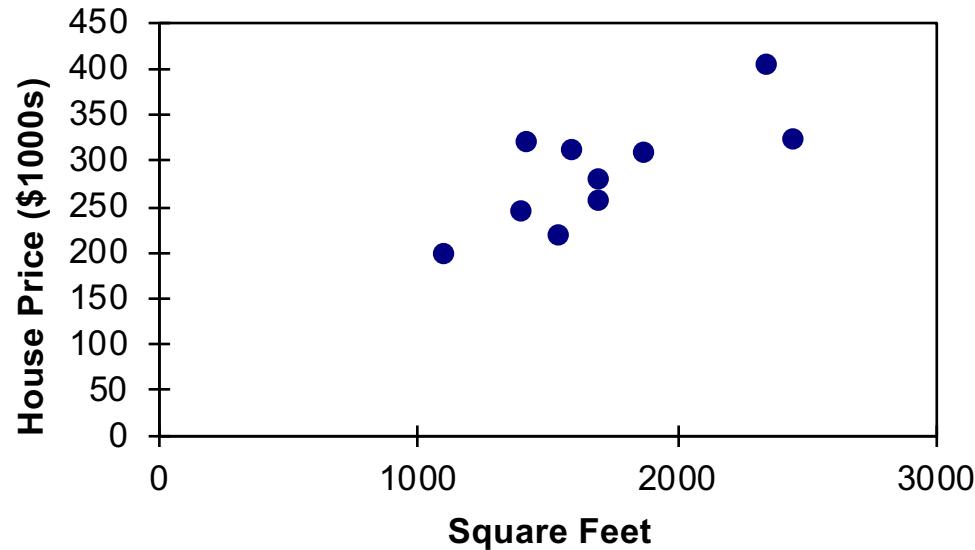
$$y^{(i)} = \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)} = \boldsymbol{\theta} \cdot \mathbf{x}^{(i)}$$

- 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

- Consistent learning may be impossible.

Prediction

Actual value

- Still, we can try to reduce the error.

Per instance the error is:

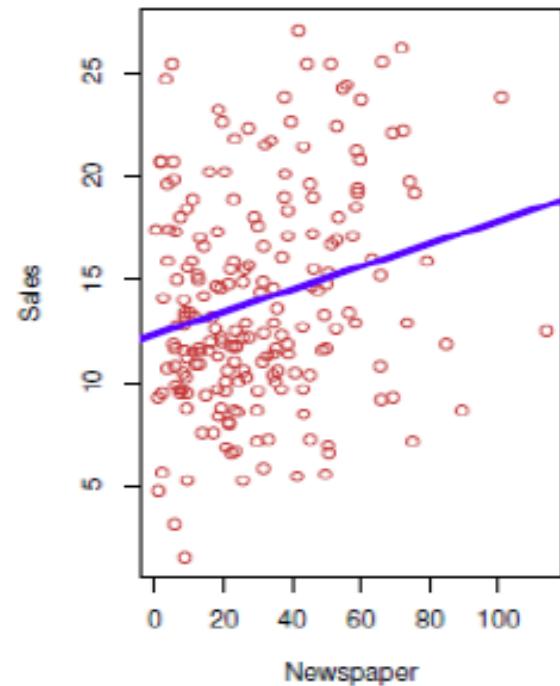
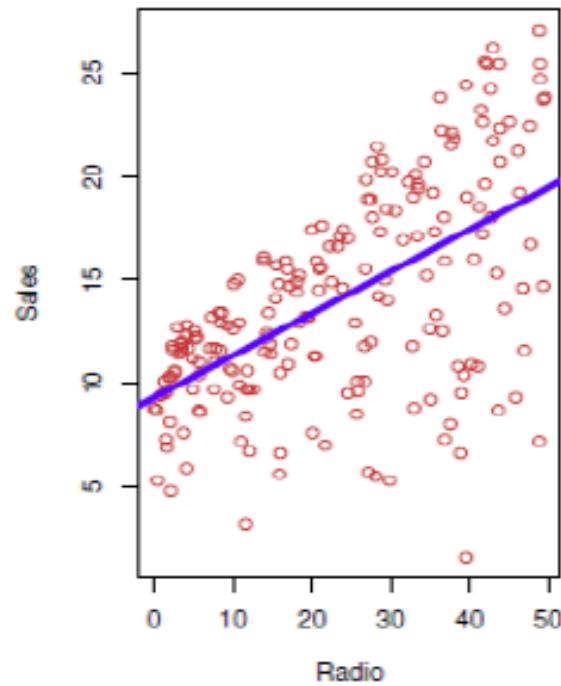
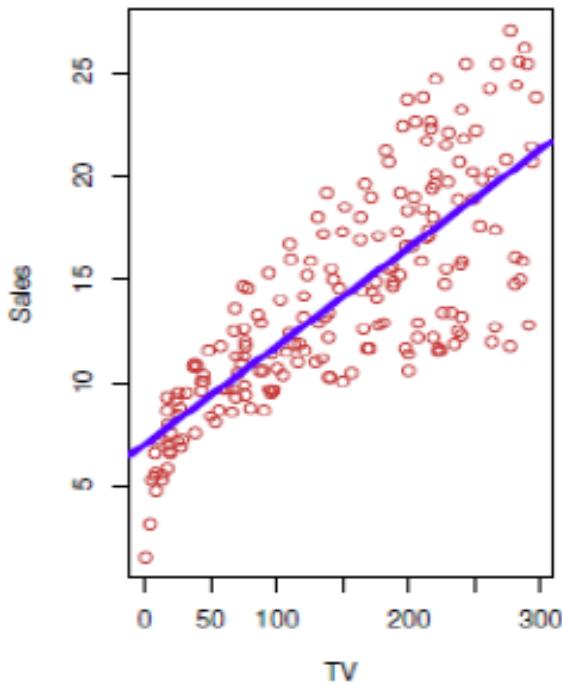
$$(\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) = \boxed{\boldsymbol{\theta} \cdot \mathbf{x}^{(i)}} - \boxed{y^{(i)}}$$

- And we now average on ALL m training instances to get our **cost function**:

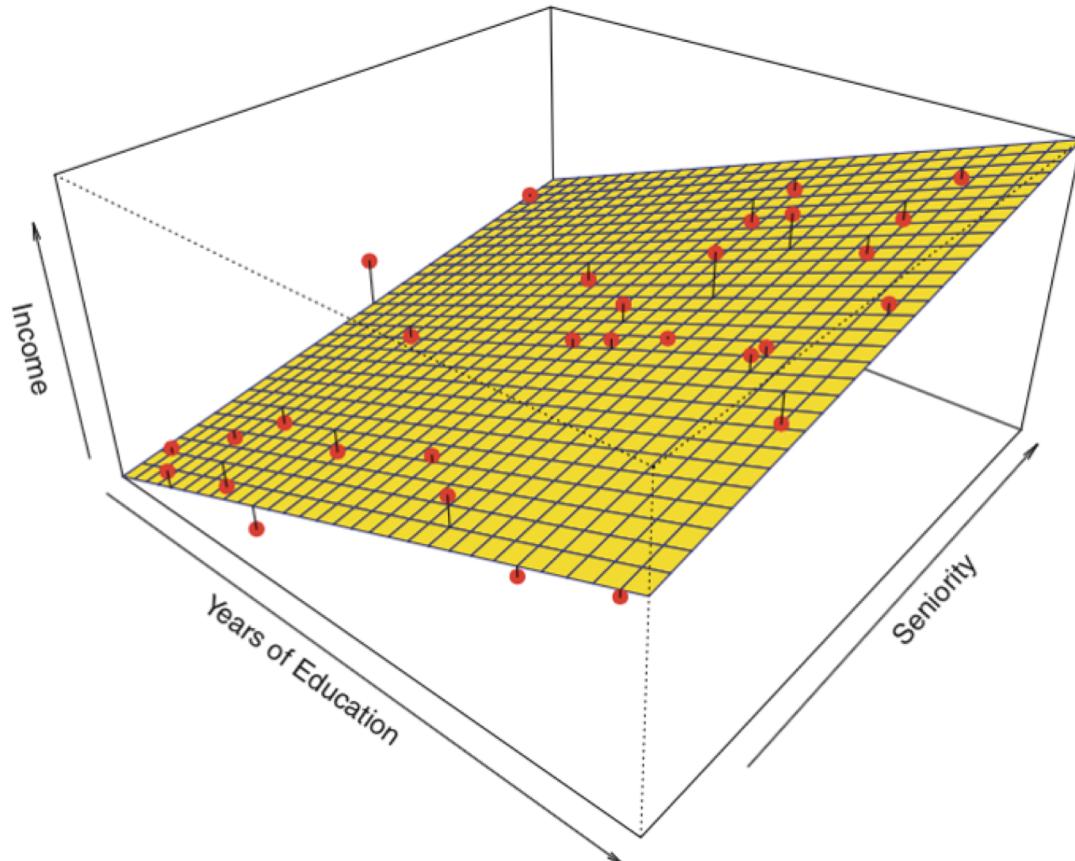
$$J(\boldsymbol{\theta}) = \frac{1}{2} \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta} \cdot \mathbf{x}^{(i)} - y^{(i)})^2$$

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

Example: advertising and sales



In higher dimensions



Minimizing the Cost Function

- Hence, our best hypothesis θ^* would be the one that minimizes the cost function:

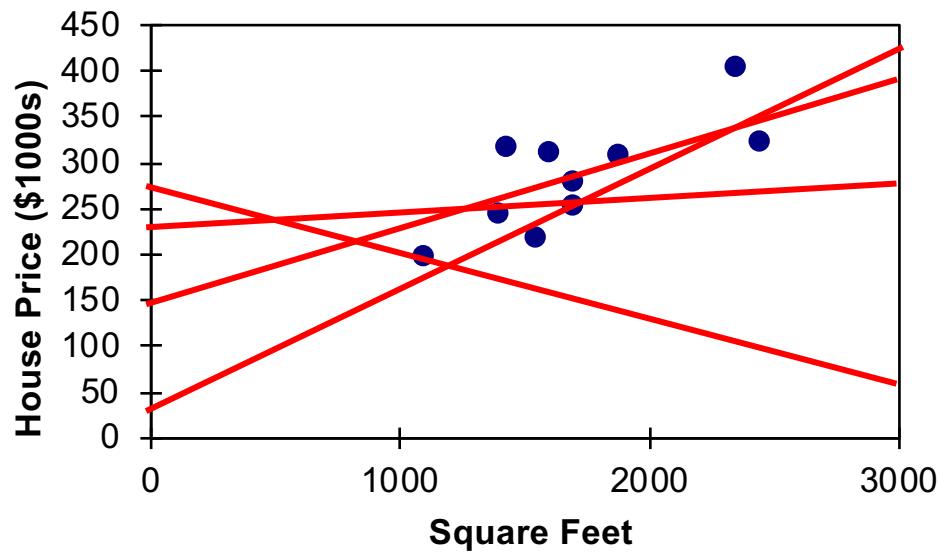
$$\theta^* = \arg \min_{\theta} [J(\theta)] = \arg \min_{\theta} \left[\frac{1}{2m} \sum_{i=1}^m (\theta \cdot \mathbf{x}^{(i)} - y^{(i)})^2 \right]$$

- How can we find it?

The Hypotheses Space

- For the simple case of a single feature we get different possible straight lines when we change θ_0 and θ_1 in the hypothesis (model)

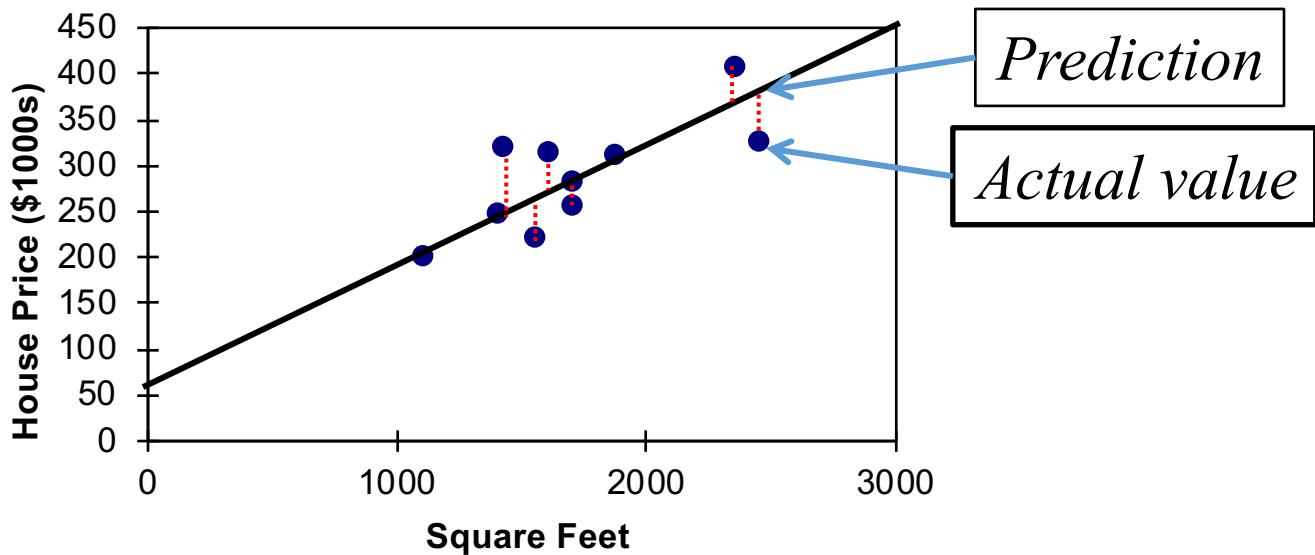
$$y = \theta_0 + \theta_1 x$$



Geometric interpretation of the error

We search for θ_0 and θ_1 that minimize the sum of (squared) errors:

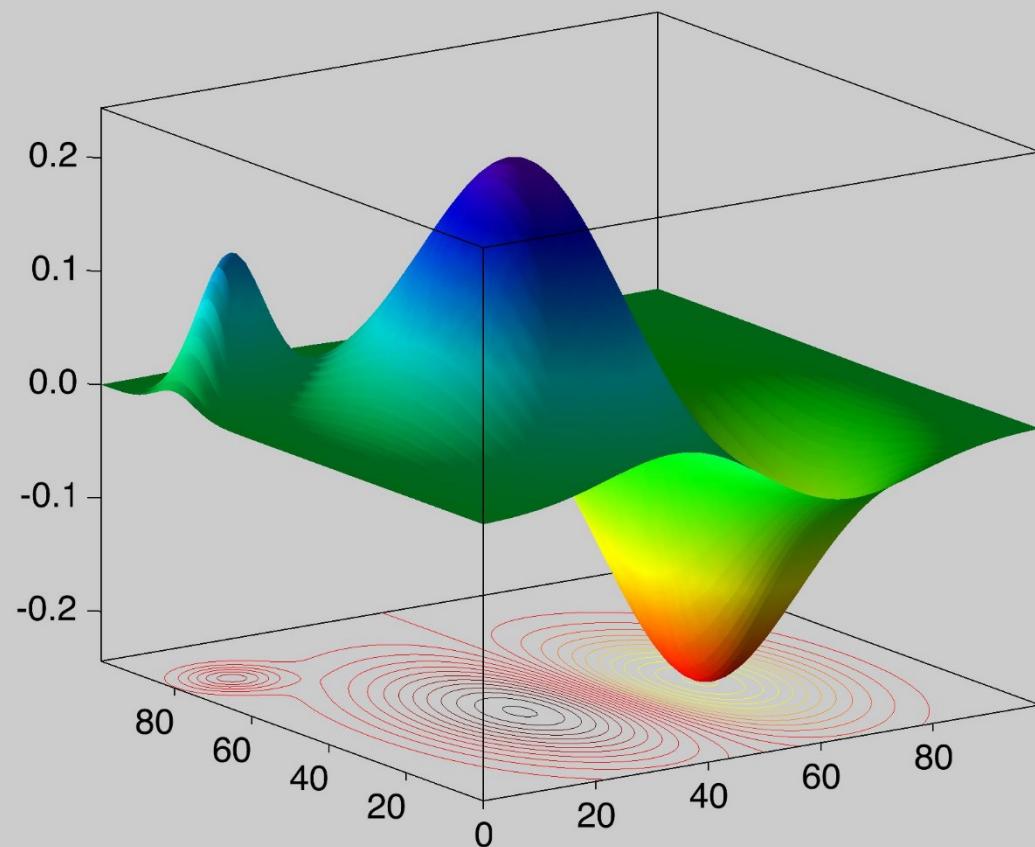
$$J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

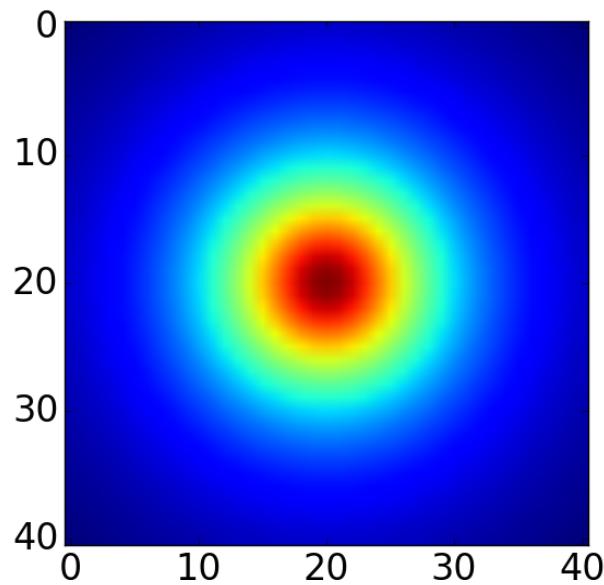
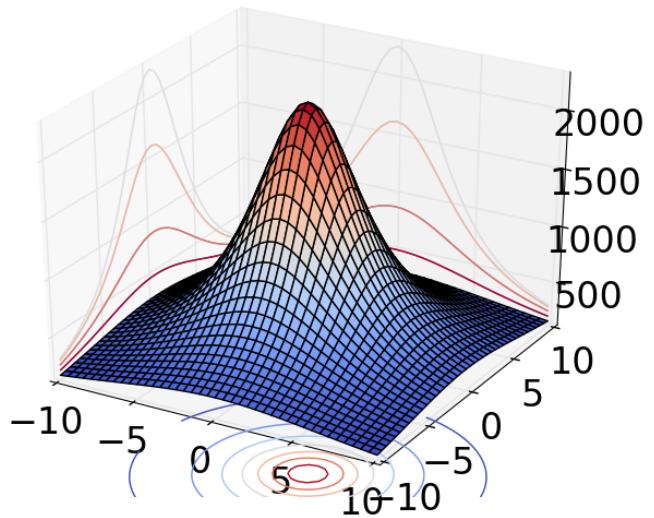
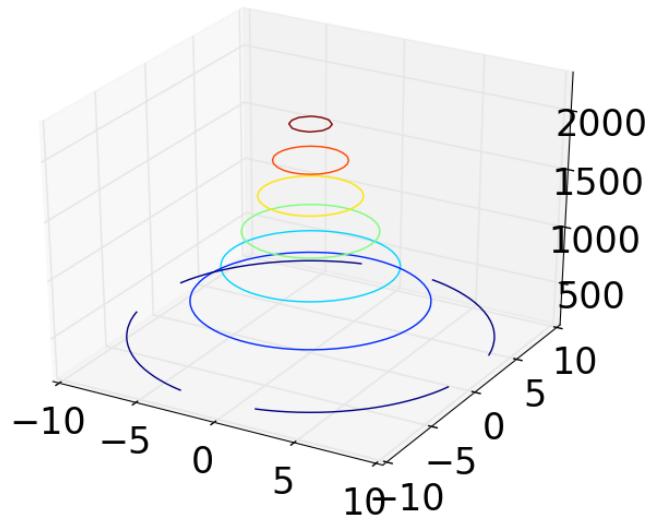
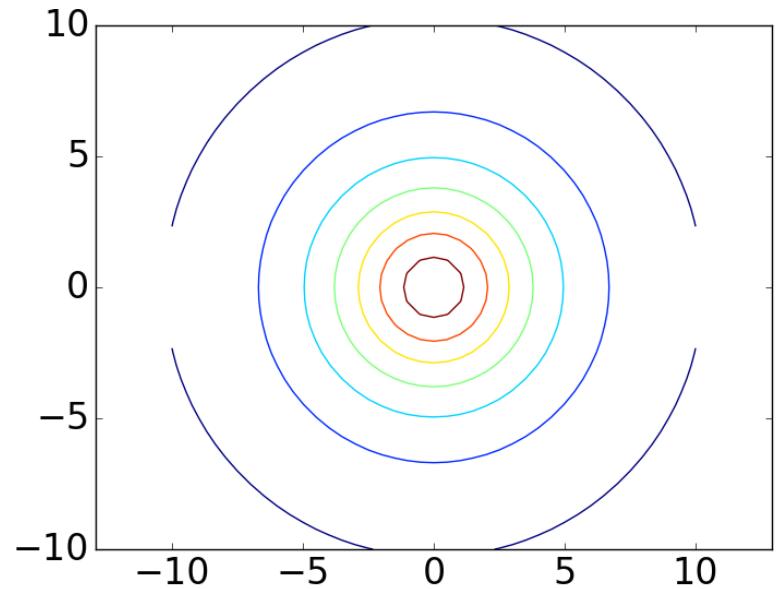


The Cost Function

$$\theta_0, \theta_1$$

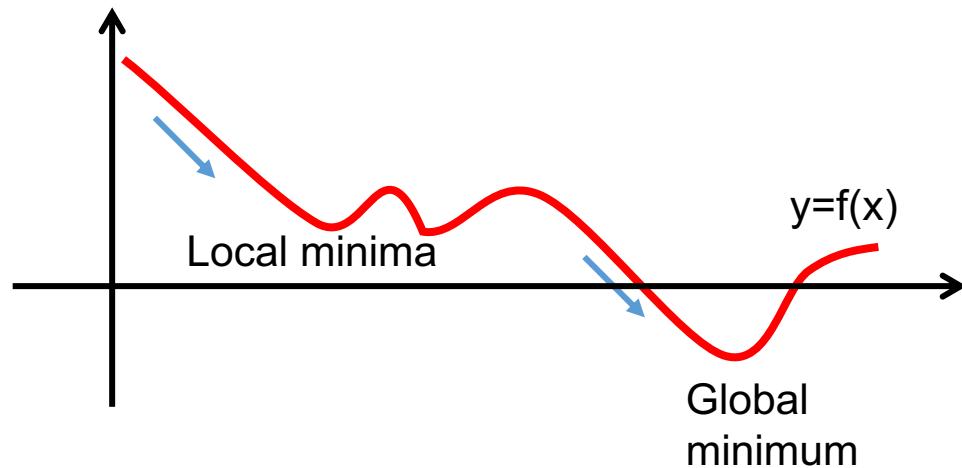
- In our simple case we have 2 parameters:
- For each value of these two parameters we can calculate the cost (the training error).
- This is a function from \mathbb{R}^2 to \mathbb{R} .





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 in the negative direction
- Can be trapped in local minimum!



Directional Derivatives

- For a differentiable 2D function $f(x_1, x_2)$ the partial derivatives in the directions x_i are

$$\frac{\partial}{\partial x_1} f, \frac{\partial}{\partial x_2} f$$

- For example ...

$$f(x, y) = 2x + 13y$$

$$f(x, y) = y \exp(x)$$

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

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

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

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

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

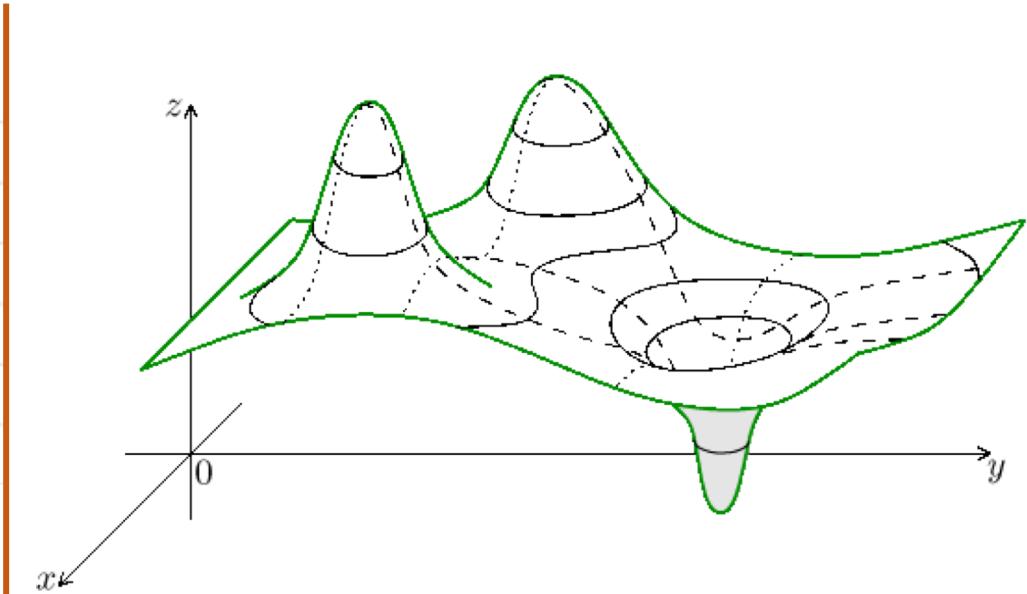
$$\frac{\partial f}{\partial x}(x, y) = y^4 e^x$$

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

Extrema and zero derivatives



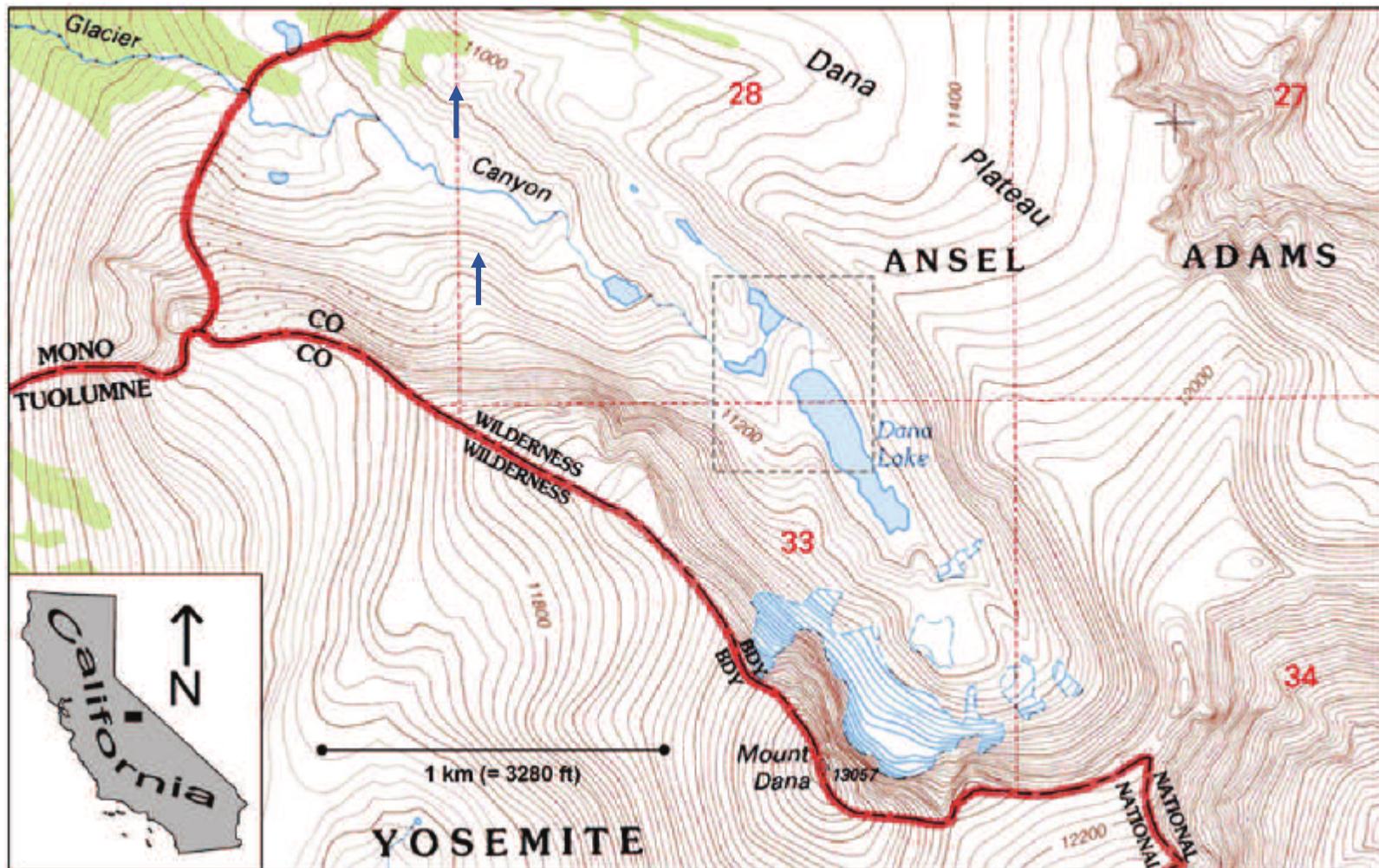
1D



2D

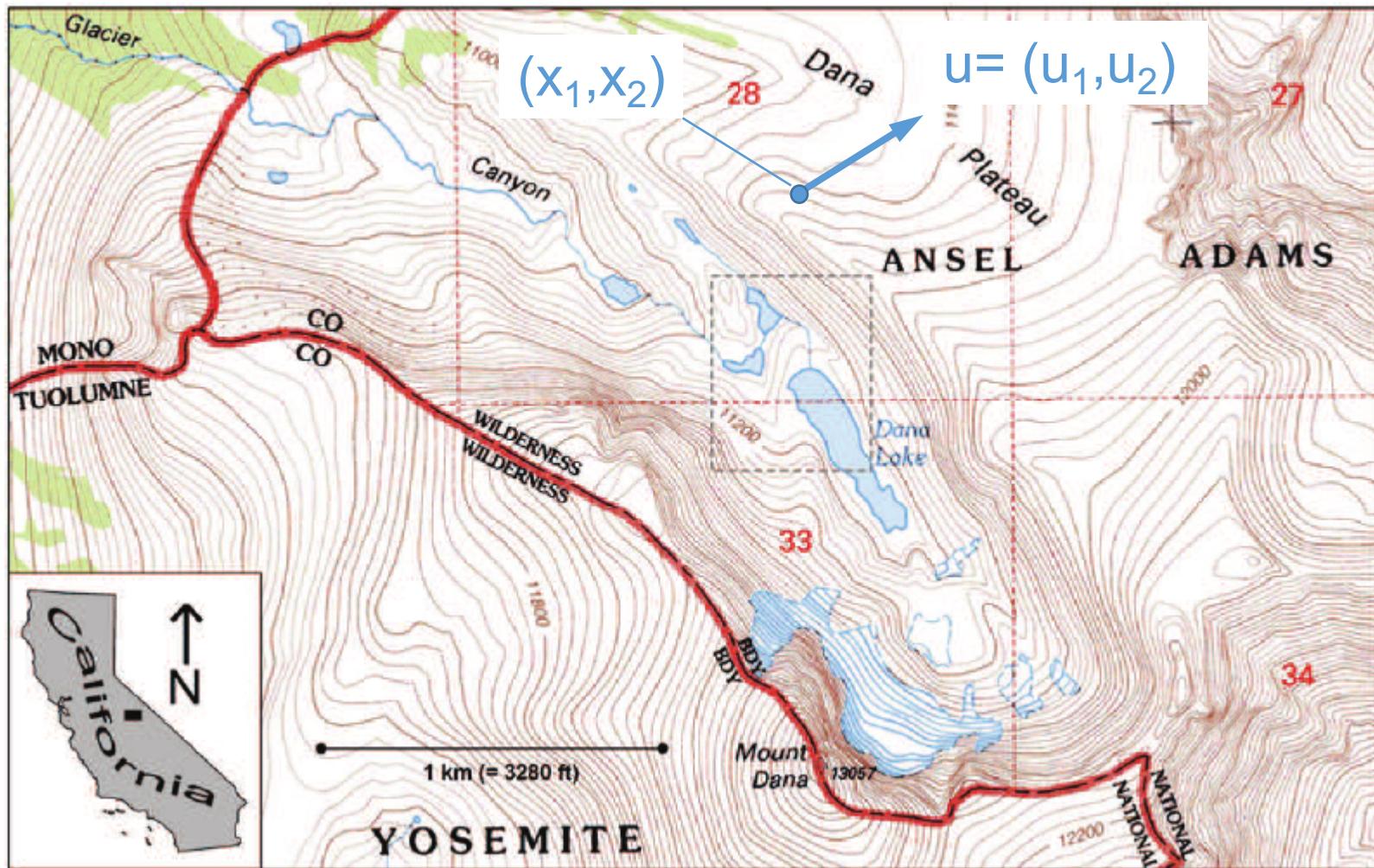
Example

What is the sign of the y directional derivative?



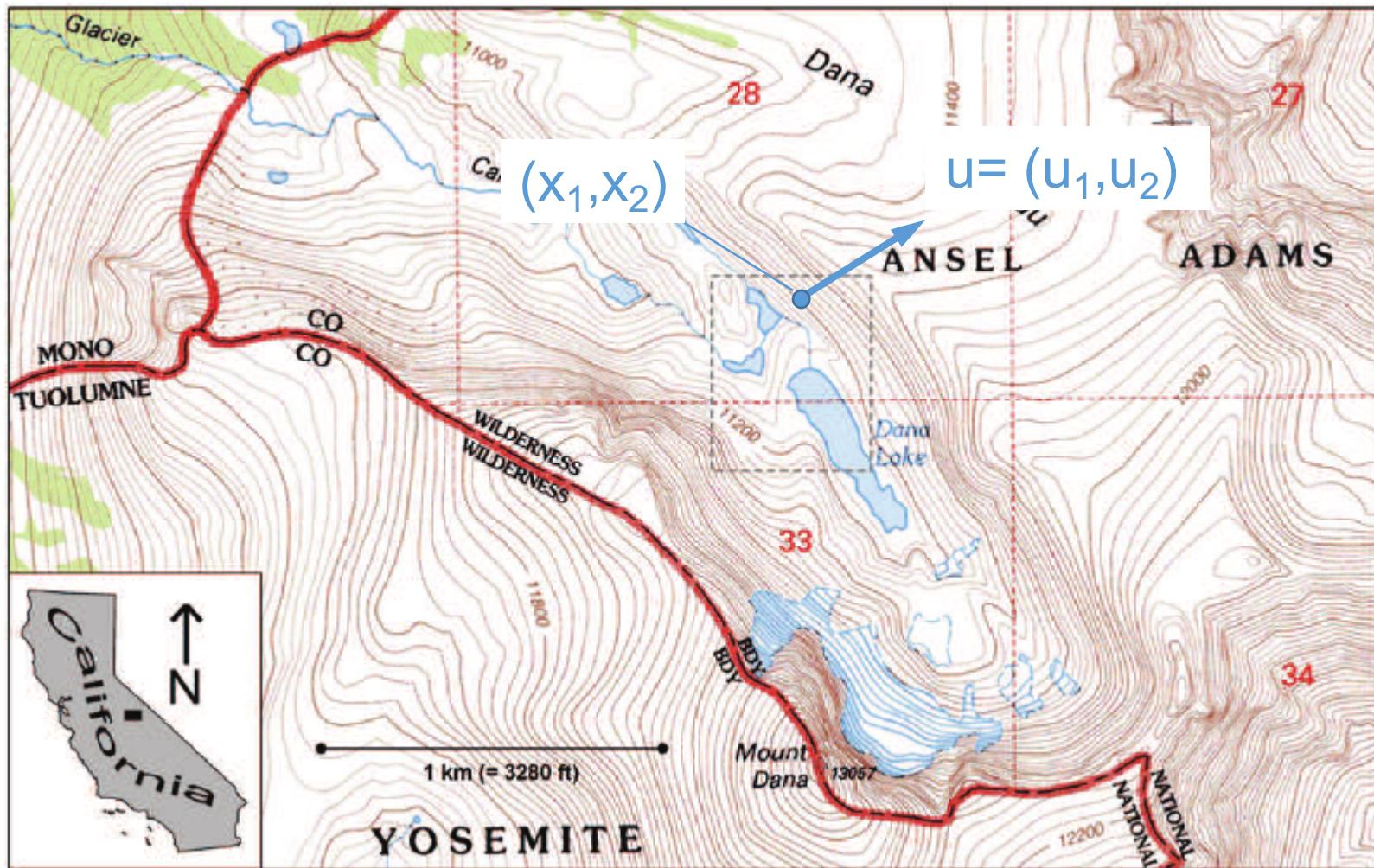
Example

Up or down?



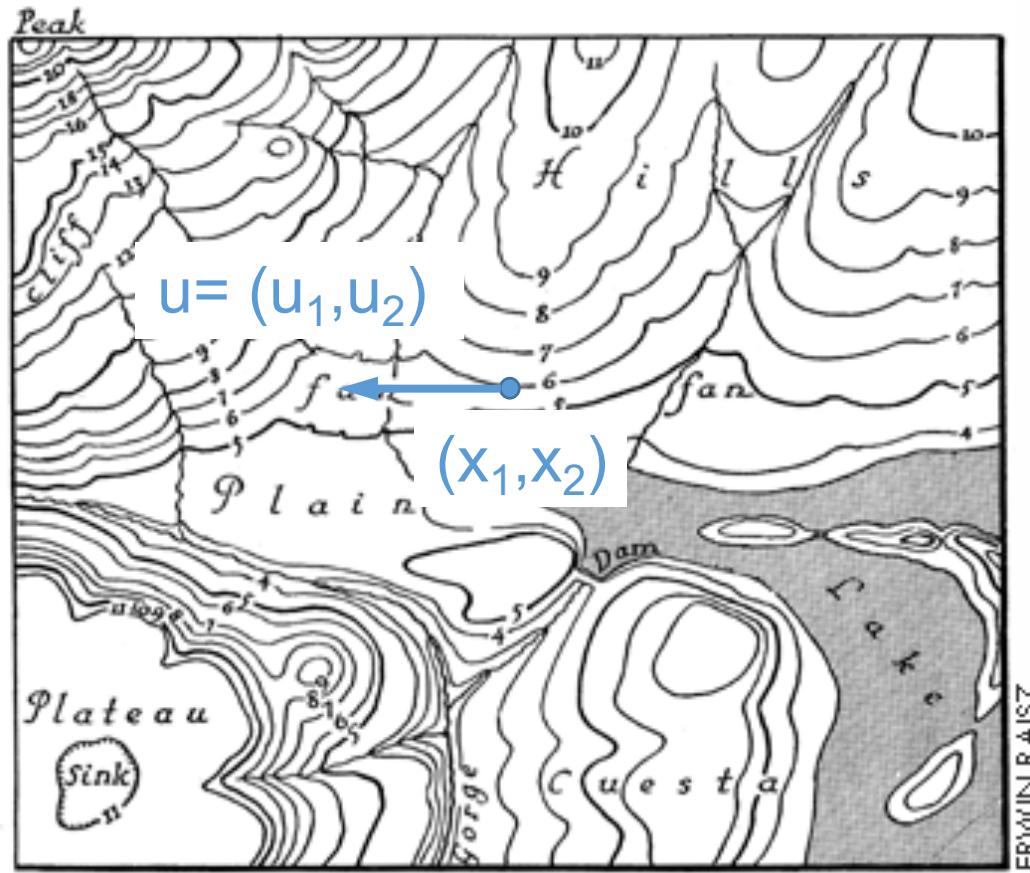
Example

Up or down?



Another ...

Up or down?



Directional Derivatives

- For a differentiable 2D function $f(x_1, x_2)$ the partial derivatives in direction x_i are $\frac{\partial}{\partial x_1} f, \frac{\partial}{\partial x_2} f$
- The derivative in a general direction $u = (u_1, u_2)$ (unit 2D vector) is called the directional derivative $D_u f$ and is defined as:

$$D_u f(x_1, x_2) = \lim_{s \rightarrow 0} \frac{f(x_1 + su_1, x_2 + su_2) - f(x_1, x_2)}{s} = \left(\frac{df}{ds} \right)_u$$

The Gradient of a function

$$\begin{aligned} D_u f(x_1, x_2) &= \left(\frac{df}{ds} \right)_u = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s} = \\ &= \frac{\partial f}{\partial x_1} u_1 + \frac{\partial f}{\partial x_2} u_2 = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right) \cdot (u_1, u_2) = \nabla f(x_1, x_2) \cdot \mathbf{u} \end{aligned}$$

This is the dot product of the direction \mathbf{u} with a vector that is called the **GRADIENT** and consists of the principle directional derivatives evaluated at (x_1, x_2) :

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

The Gradient of a function

Define the
GRADIENT of f :

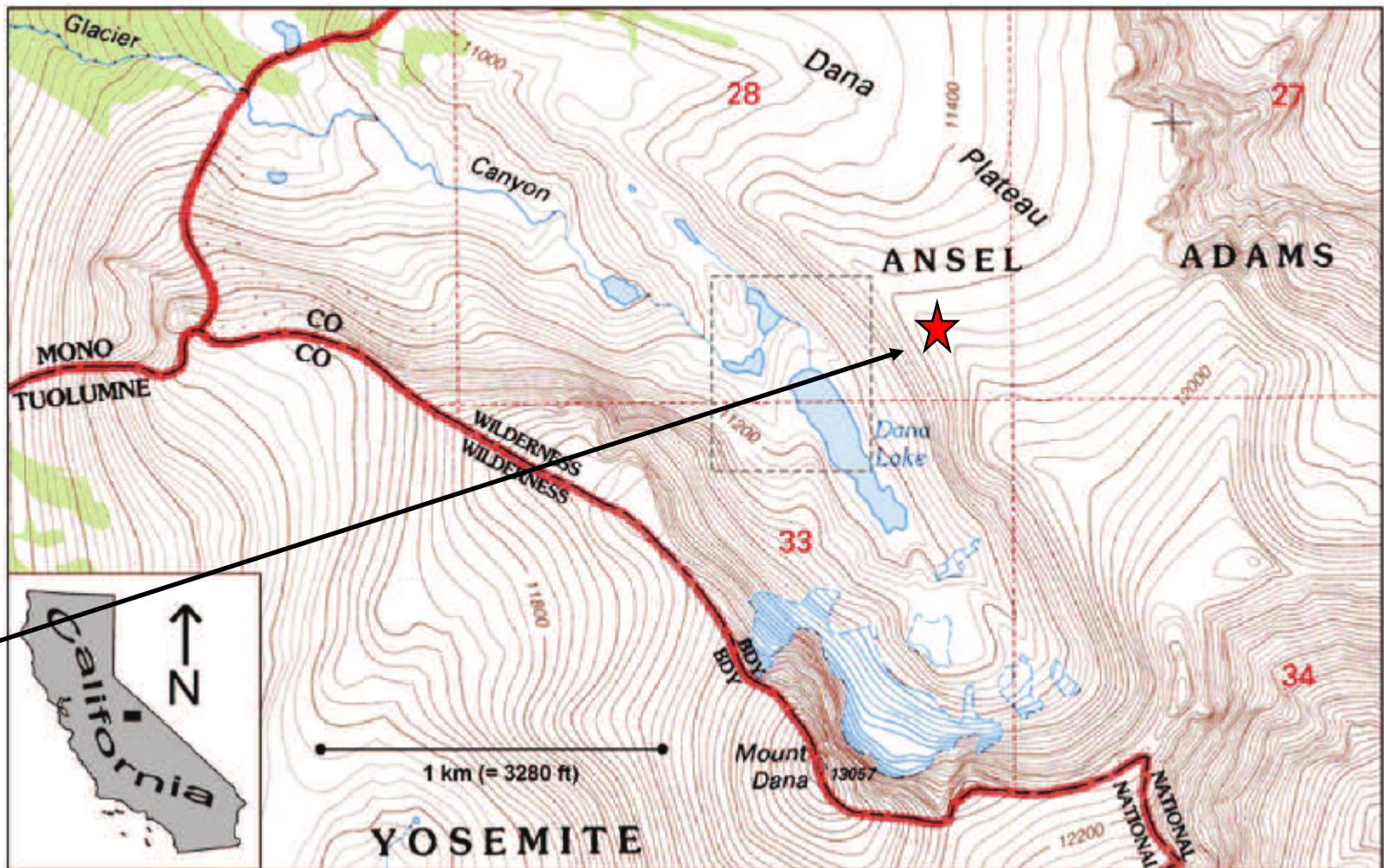
$$\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)$$

Thm:

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

$$D_u f(x) = \nabla f(x) \cdot u$$

Most Rapid Increase at a point x



Most Rapid Increase at a point x

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

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

- Where β is the angle between \mathbf{u} and ∇f
- However, $\cos \beta \leq 1$.
- Therefore:
 1. The greatest increase in the function happens in the direction of the gradient (i.e. $\beta = 0$)
 2. The greatest decrease is in the direction $-\nabla f$ (i.e $\beta = 180^\circ$)

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.



The gradient of n-dimensional functions

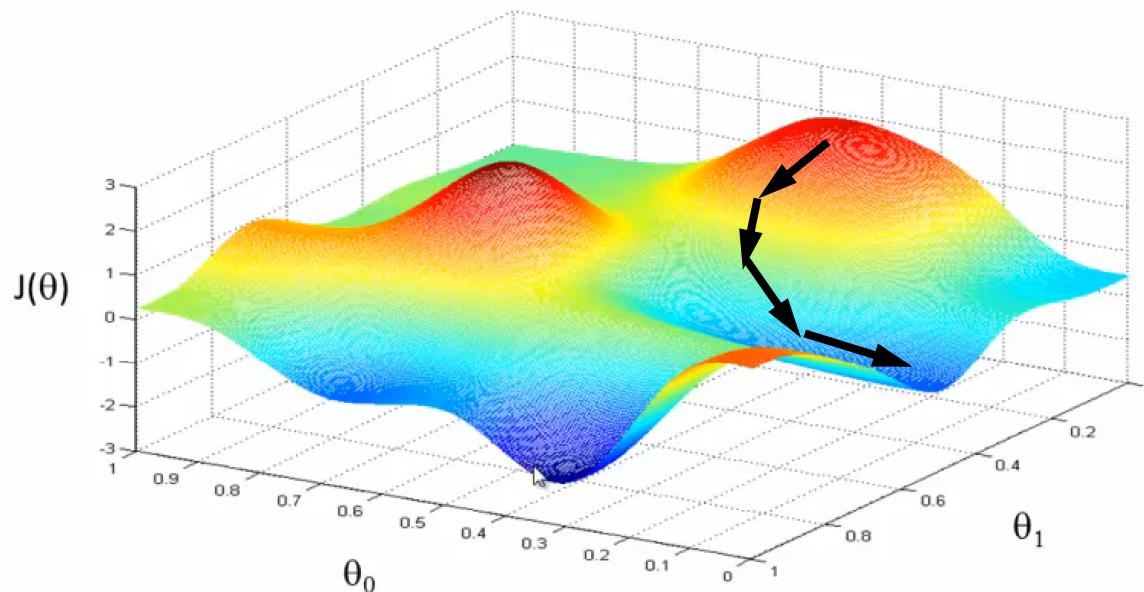
- For an n-D function $y = f(x_1, x_2, \dots, x_n)$ the gradient is defined as :

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

- The gradient vector at any given point is the direction of greatest increase of f at this point.
- The direction of greatest decrease of f is the negative of the gradient: $-\nabla f(x_1, x_2, \dots, x_n)$

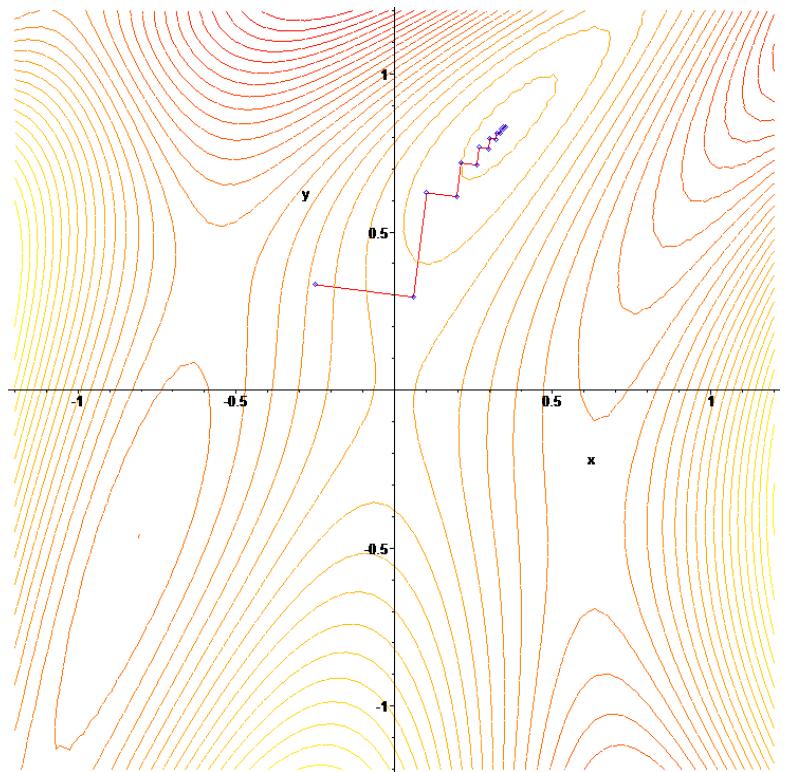
How to find a minimum of a (reasonably smooth) n-dimensional function?

- Take steps in the direction of maximum decrease – i.e. the negative gradient direction!
- This is called **Gradient Descent**



Gradient Descent Steps

- For a small enough $\alpha > 0$ we have
$$f(x) - \alpha \|\nabla f\| < f(x),$$
so we take one step of size α in the opposite direction of the gradient
- Note: such steps can zig-zag



Back to Minimizing the Cost Function

- Recall that our best hypothesis θ^* is the one that minimizes the cost function:

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

- This is a real valued function of $\vec{\theta} \in \mathbb{R}^n$
- So - we can start with some initial guess $\vec{\theta}_0^*$ and then use gradient descent

Gradient Descent Algorithm

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

For all j ,

$$\text{Update } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

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

Minimizing the Cost Function

- Hence, our best hypothesis θ^* would be the one that minimizes the cost function:

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

- How can we find it?

Calculating the Partial Derivative

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

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

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

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

Gradient Descent for Linear Regression

- Initialize $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)$
- Repeat until you reach a minimum (or stop cdn):
 - For all $0 \leq j \leq n$,
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$
- In words: set the new θ_j to the current θ_j minus the learning rate (α) times the partial derivative of the error function with respect to θ_j , computed at the current θ .
- Also remember that $x_0^{(i)} = 1$

Recitation

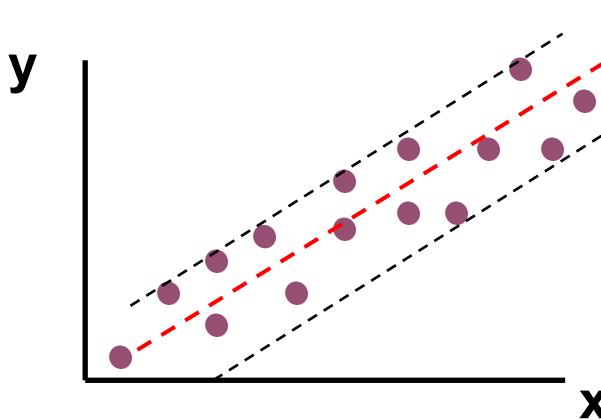
- Description of HWA1
- Feature scaling
- Learning rate
- Python

Extra slides/notes

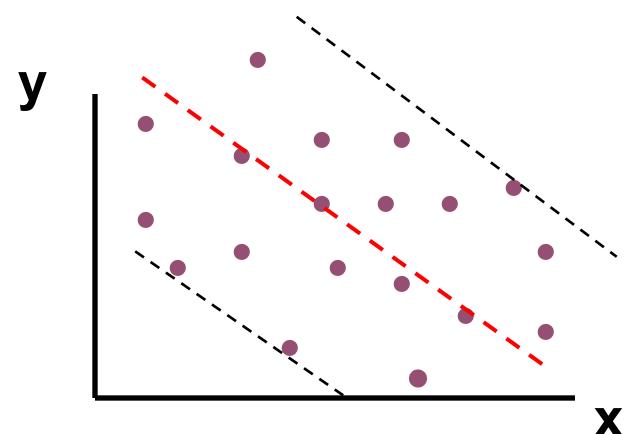
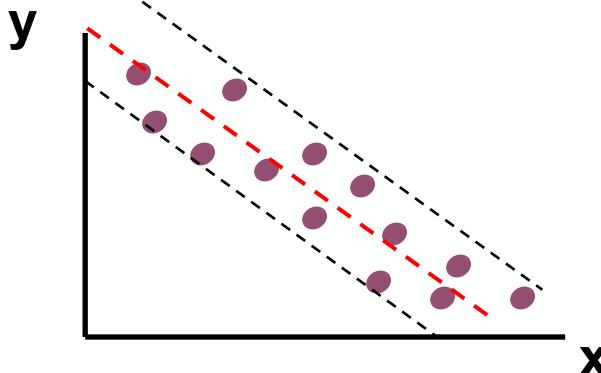
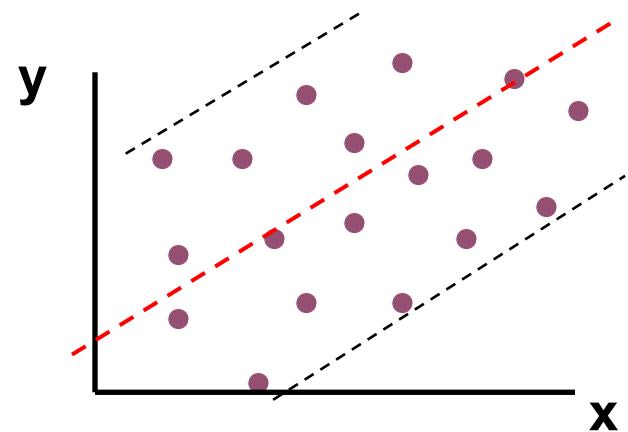
- Correlations
- Pseudo-inverse of a matrix
- Going beyond Linear

Strong vs. Weak Linear Relationships

Strong relationships

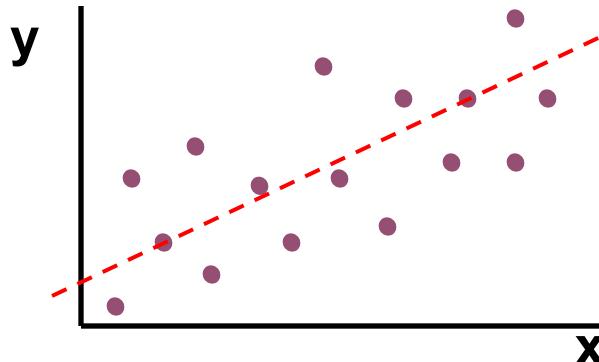


Weak relationships

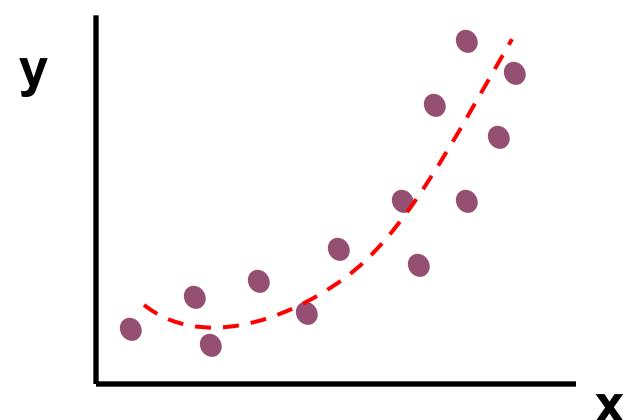
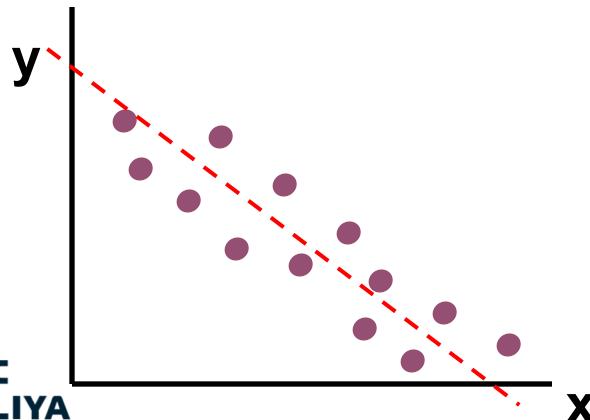
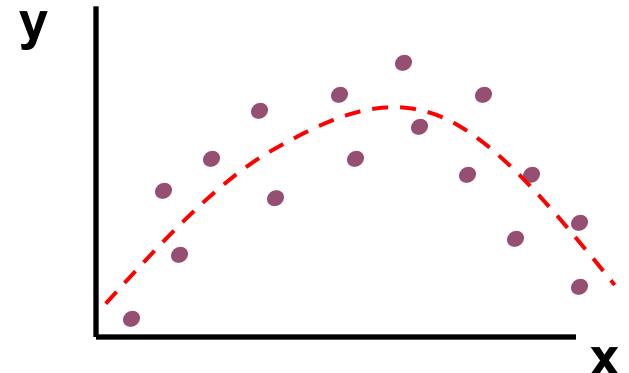


Not Everything is Linear!

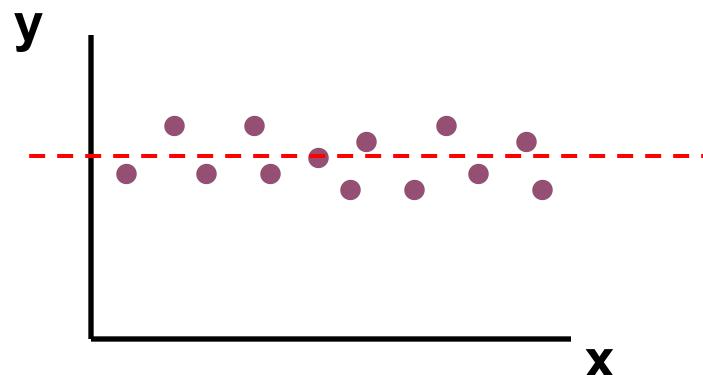
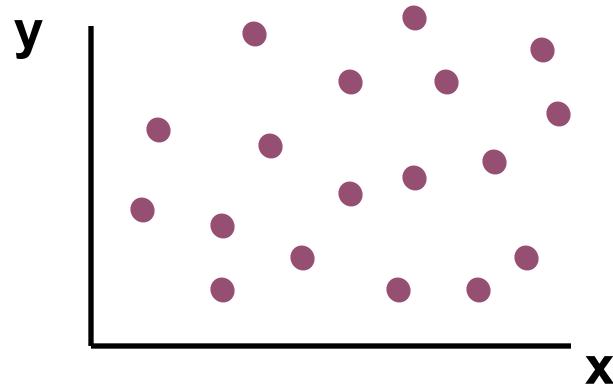
Linear relationships



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

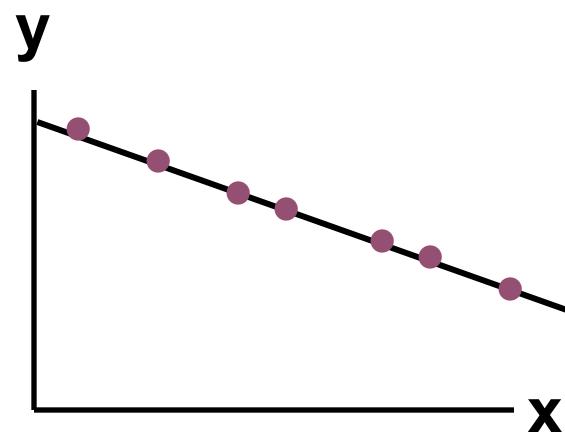
x = Value of an explaining variable

y = Value of the dependent variable

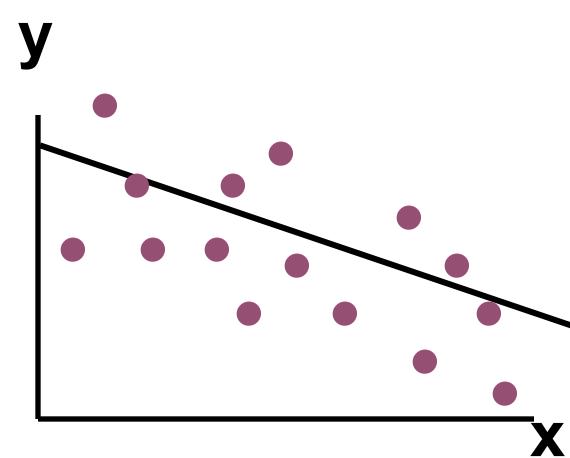
Features of r

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

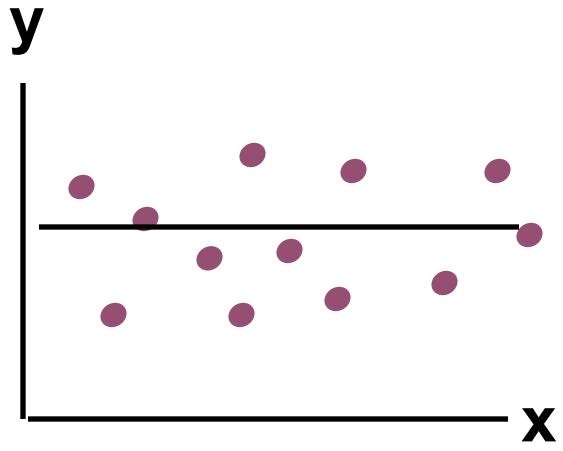
Examples of r Values



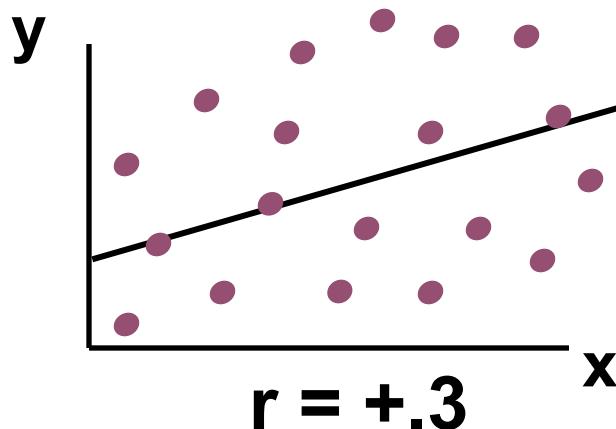
$r = -1$



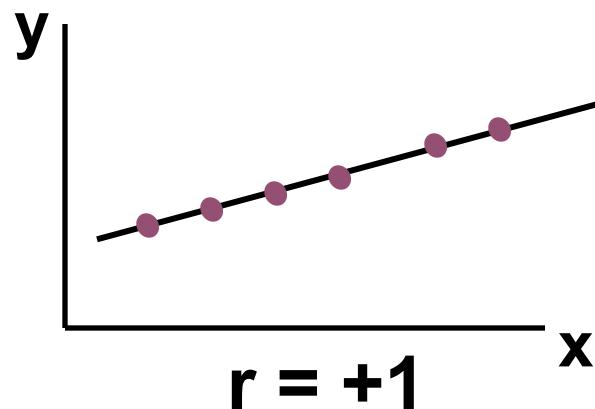
$r = -.6$



$r = 0$

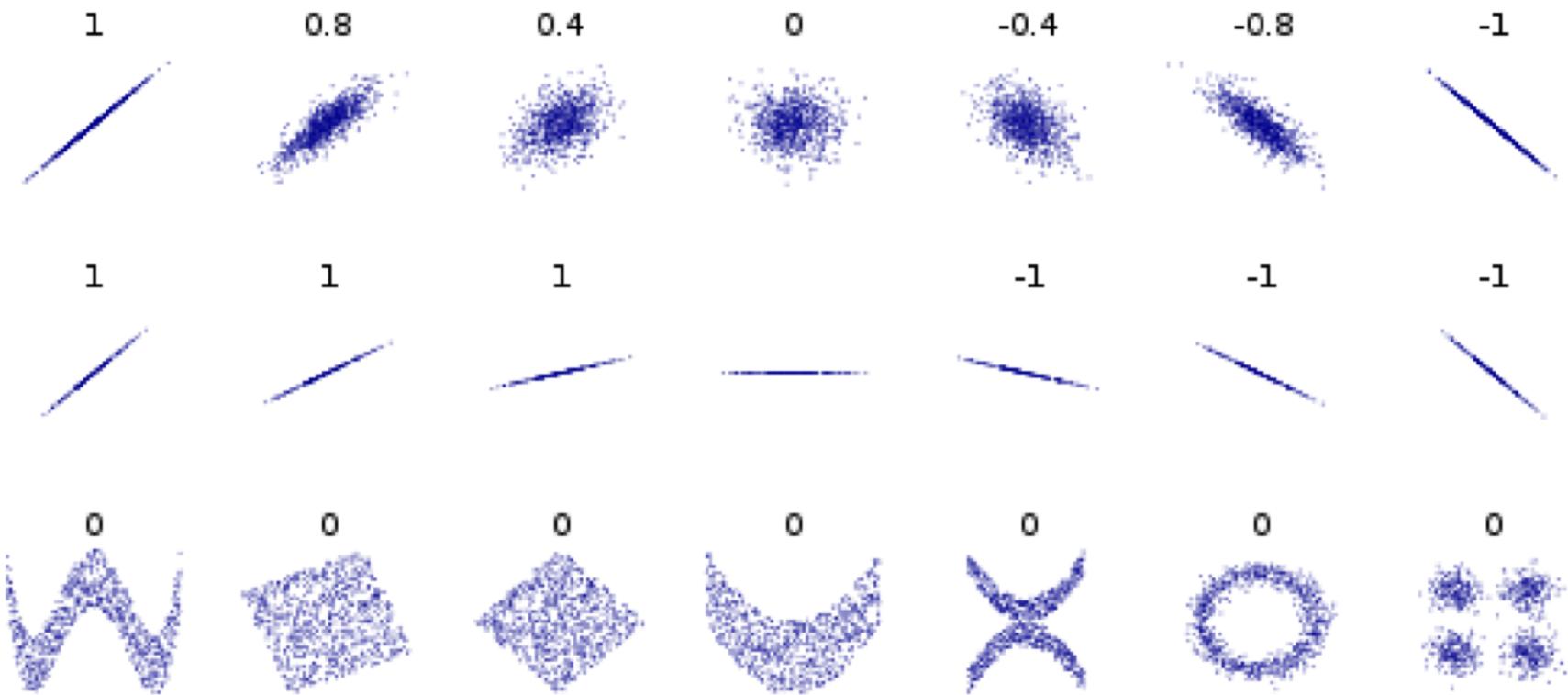


$r = +.3$



$r = +1$

More Complex Relationships May Not Be Captured Using Pearson r



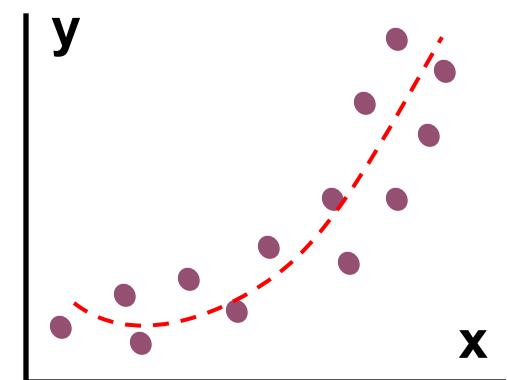
Polynomial Regression

- 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

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$



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?

$$X \cdot \theta = y$$

$$\begin{pmatrix} X_0^{(1)} & X_1^{(1)} & \dots & X_n^{(1)} \\ X_0^{(2)} & X_1^{(2)} & \dots & X_n^{(2)} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ X_0^{(m)} & X_1^{(m)} & \dots & X_n^{(m)} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

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

Minimum when Derivative is 0

- Lets look at the error vector $e = \vec{x} \cdot \vec{\theta} - y$
- Minimizing the squared length of this error vector is equivalent to minimizing the sum of squared distance function:

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

- Forming the derivatives yields (see notes):

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

Pseudo-inverse & the Normal Equation

- Setting the gradient to zero yields the necessary condition for minimum (see notes):

$$\mathbf{X}^T \mathbf{X} \cdot \vec{\theta} = \mathbf{X}^T \vec{\mathbf{Y}}$$

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

$$\vec{\theta} = \text{pinv}(\mathbf{X}) \vec{\mathbf{Y}} \quad \text{where} \quad \text{pinv}(\mathbf{X}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

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

Gradient Descent and Pseudo Inverse

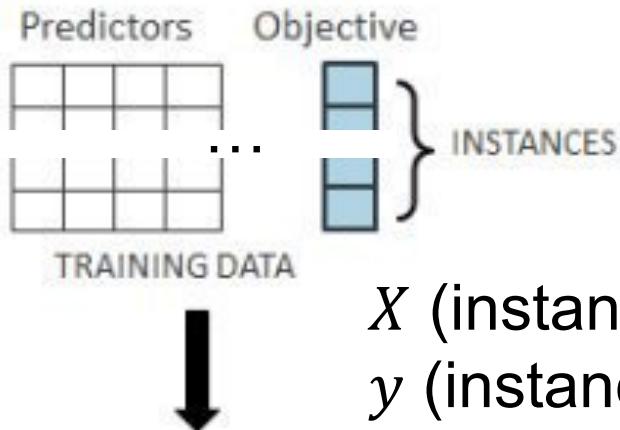
- $\mathbf{X}^T \mathbf{X}$ can still be singular (or not full rank) and not have an inverse. This can be resolved with some more algebra.
- This pinv technique doesn't work for all error functions J. Gradient descent is more general.
- Gradienet descent allows parallelization.

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



The learning
algorithm
(regression)

