

# INTRODUCTION TO DATA SCIENCE

**JOHN P DICKERSON**

Lecture #17 – 10/26/2021

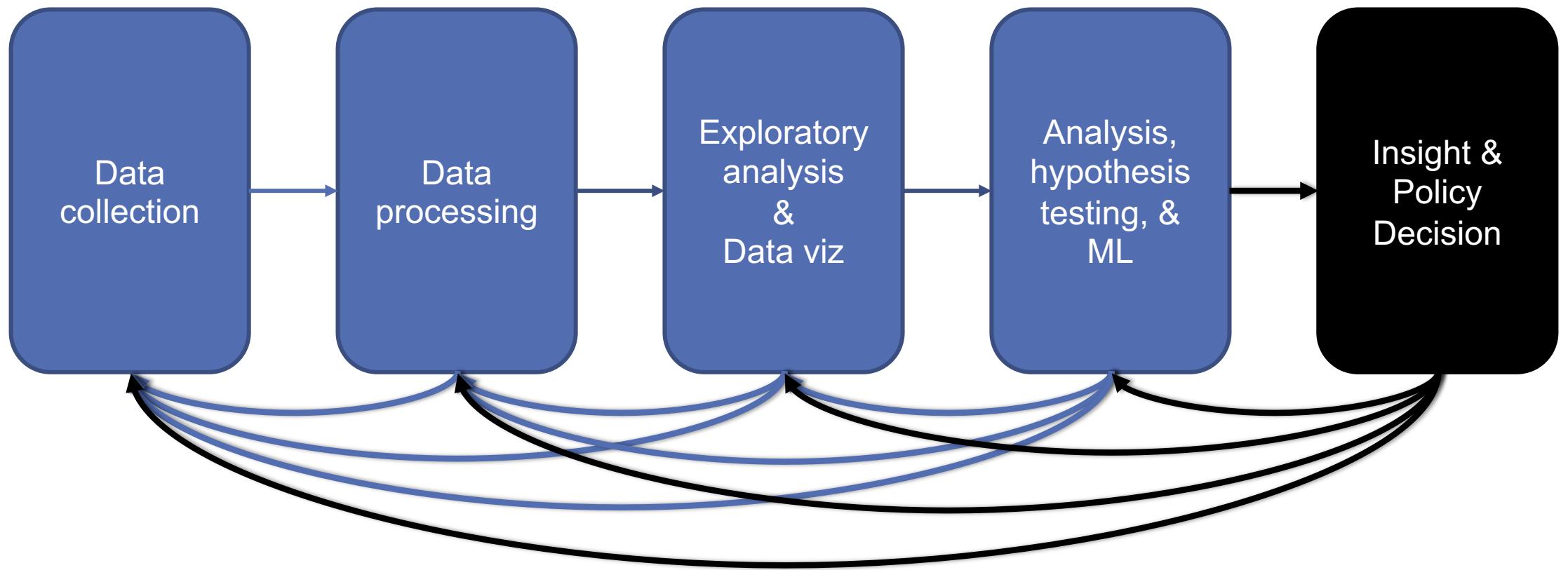
**CMSC320**  
Tuesdays & Thursdays  
5:00pm – 6:15pm

<https://cmsc320.github.io/>



**COMPUTER SCIENCE**  
UNIVERSITY OF MARYLAND

# THE DATA LIFECYCLE

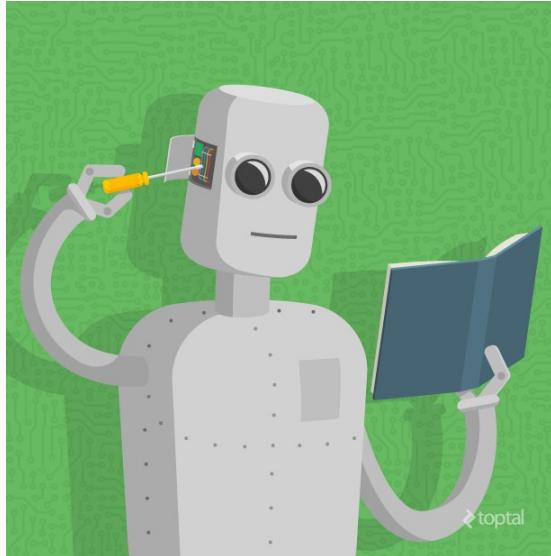


# TODAY'S LECTURE

## Introduction to machine learning

- How did we actually come up with that linear model from last class?
- Basic setup and terminology; linear regression & classification

Thanks to: Zico Kolter (CMU), David Kauchak (Pomona), Nick Mattei (Tulane)



*First GIS result for “machine learning”*

# RECALL: EXPLICIT EXAMPLE FROM THE NLP LECTURES

Score  $\psi$  of an instance  $x$  and class  $y$  is the sum of the weights for the features in that class:

$$\begin{aligned}\psi_{xy} &= \sum \theta_n f_n(x, y) \\ &= \boldsymbol{\theta}^T \mathbf{f}(x, y)\end{aligned}$$

Let's compute  $\psi_{x_1, y=hates\_cats} \dots$

- $\psi_{x_1, y=hates\_cats} = \boldsymbol{\theta}^T \mathbf{f}(x_1, y = hates\_cats = 0)$
- $= 0*1 + -1*1 + 1*0 + -0.1*1 + 0*0 + 1*0 + -1*0 + 0.5*0 + 1*1$
- $= -1 - 0.1 + 1 = -0.1$

$$\boldsymbol{\theta}^T = \boxed{0 \quad -1 \quad 1 \quad -0.1 \quad 0 \quad 1 \quad -1 \quad 0.5 \quad 1}$$



<i>hates_cats</i>	<i>likes_cats</i>	(1)
1	I	
1	like	
0	hate	
1	cats	
0	I	
0	like	
0	hate	
0	cats	
1	-	

$$\mathbf{f}(x_1, y = 0)$$

# RECALL: EXPLICIT EXAMPLE FROM THE NLP LECTURES

Saving the boring stuff:

- $\psi_{x_1,y=hates\_cats} = -0.1$ ;  $\psi_{x_1,y=likes\_cats} = +2.5$  Document 1: I like cats
- $\psi_{x_2,y=hates\_cats} = +1.9$ ;  $\psi_{x_2,y=likes\_cats} = +0.5$  Document 2: I hate cats

We want to predict the class of each document:

$$\hat{y} = \arg \max_y \theta^T f(\mathbf{x}, y)$$

Document 1:  $\text{argmax}\{ \psi_{x_1,y=hates\_cats}, \psi_{x_1,y=likes\_cats} \}$  ????????

Document 2:  $\text{argmax}\{ \psi_{x_2,y=hates\_cats}, \psi_{x_2,y=likes\_cats} \}$  ????????



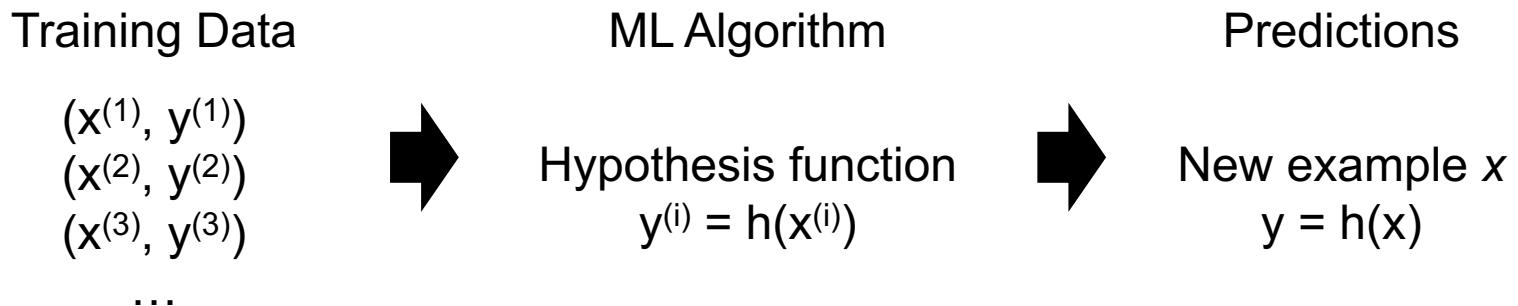
# MACHINE LEARNING

We used a **linear model** to classify input documents

The model parameters  $\theta$  were given to us a priori

- (John created them by hand.)
- Typically, we cannot specify a model by hand.

**Supervised machine learning provides a way to automatically infer the predictive model from labeled data.**



# MACHINE LEARNING

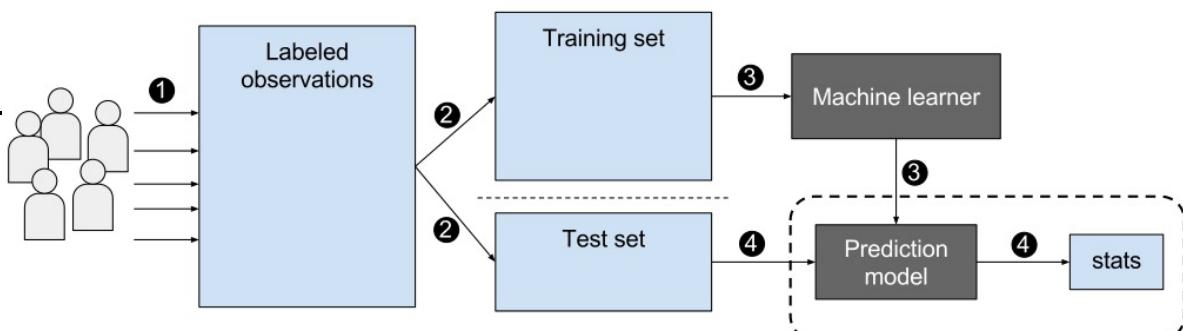
Arthur Samuel 1959, “... give[s] computers the ability to learn without being explicitly programmed.”



Tom M. Mitchell, “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$  if its performance at tasks in  $T$  as measured by  $P$ , improves with experience  $E$ .”

## Major Machine Learning Paradigms:

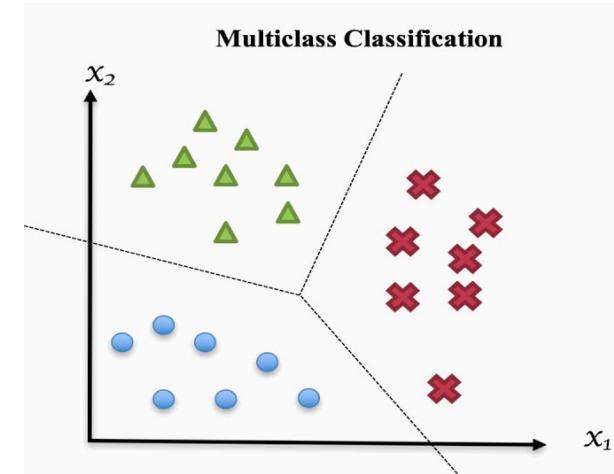
- **Supervised Learning:** provide labeled examples – machine learns to identify these (classification).
- **Unsupervised Learning:** unlabeled examples – machine learns to differentiate these (clustering).
- **Reinforcement Learning:** provide reward signal – given as feedback to the machine (RL).



# LEARNING: TYPES OF FEEDBACK

## Supervised Learning.

- Learn a **function** from examples of its inputs and outputs.
- E.g., An agent is presented with many camera images and is told to learn which ones contain busses.
- Agent learns to map from images to Boolean output 0/1 of bus not present/present.
- Learning decision trees is a form of supervised learning.



## Unsupervised Learning.

- Learn patterns in the input with no output values supplied.
- E.g.: Identify communities on the Internet.

## Reinforcement Learning.

- Learn from reinforcement (occasional rewards).
- E.g., An agent learns how to play backgammon or go or chess against itself.



# TERMINOLOGY

**Input features:**  $x^{(i)} \in \mathbb{R}^n, i = 1, \dots, m$

	-	like	hate	cats
$x^{(1)T} =$	1	1	0	1
$x^{(2)T} =$	1	0	1	1

**Outputs:**  $y^{(i)} \in \mathcal{Y}, i = 1, \dots, m$

$$y^{(i)} \in \{0, 1\} = \{\text{hates\_cats}, \text{likes\_cats}\}$$

**Model parameters:**  $\theta \in \mathbb{R}^n$

$$\theta^T = [0 \quad -1 \quad 1 \quad -0.1 \quad 0 \quad 1 \quad -1 \quad 0.5 \quad 1]$$

# TERMINOLOGY

**Hypothesis function:**  $h_\theta: \mathbb{R}^n \rightarrow \mathcal{Y}$

E.g., linear classifiers predict outputs using:

$$h_\theta(x) = \theta^T x = \sum_{j=1}^n \theta_j \cdot x_j$$

**Loss function:**  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

- Measures difference between a prediction and the true output
- E.g., squared loss:  $\ell(\hat{y}, y) = (\hat{y} - y)^2$
- E.g., hinge loss:  $\ell(y) = \max(0, 1 - t \cdot y)$

Output  $t = \{-1, +1\}$  based  
on -1 or +1 class label

Classifier score  $y$

# THE CANONICAL MACHINE LEARNING PROBLEM

At the end of the day, we want to learn a hypothesis function that predicts the actual outputs well.

Choose the parameterization  
that minimizes loss!

$\text{minimize}_{\theta}$

Over all possible  
parameterizations

$$\sum_{i=1}^m \ell(h_{\theta}(x^{(i)}), y^{(i)})$$

Given an hypothesis  
function and loss function

\*Not actually what we want – want it over the world of inputs – will discuss later ...

# HOW DO I MACHINE LEARN?

## 1. What is the hypothesis function?

- Domain knowledge and EDA can help here.

## 2. What is the loss function?

- We've discussed two already: squared and absolute.

## 3. How do we solve the optimization problem?

- (We'll cover gradient descent and stochastic gradient descent in class, but if you are interested, take CMSC422!)



*First GIS result for “optimization”*

## ASIDE: LOSS FUNCTIONS



# QUICK ASIDE ABOUT LOSS FUNCTIONS

Say we're back to classifying documents into:

- hates\_cats, translated to label  $y = -1$
- likes\_cats, translated to label  $y = +1$

We want some parameter vector  $\theta$  such that:

- $\psi_{xy} > 0$  if the feature vector  $x$  is of class likes\_cat; ( $y = +1$ )
- $\psi_{xy} < 0$  if  $x$ 's label is  $y = -1$

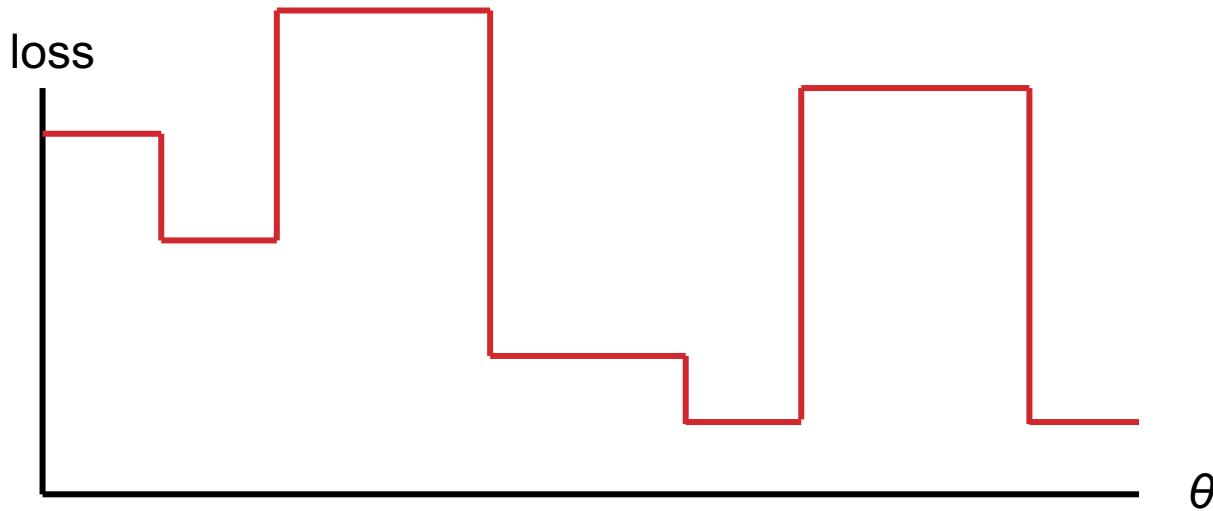
We want a hyperplane that separates positive examples from negative examples.

Why not use 0/1 loss; that is, the number of wrong answers?

$$\arg \min_{\theta} \sum_{i=1}^n \mathbf{1} [y^{(i)} \cdot \langle \theta, x^{(i)} \rangle \leq 0]$$

# MINIMIZING 0/1 LOSS IN A SINGLE DIMENSION

$$\sum_{i=1}^n \mathbf{1} [y^{(i)} \cdot \langle \theta, x^{(i)} \rangle \leq 0]$$



Each time we change  $\theta$  such that the example is right (wrong) the loss will increase (decrease)

# MINIMIZING 0/1 LOSS OVER ALL $\theta$

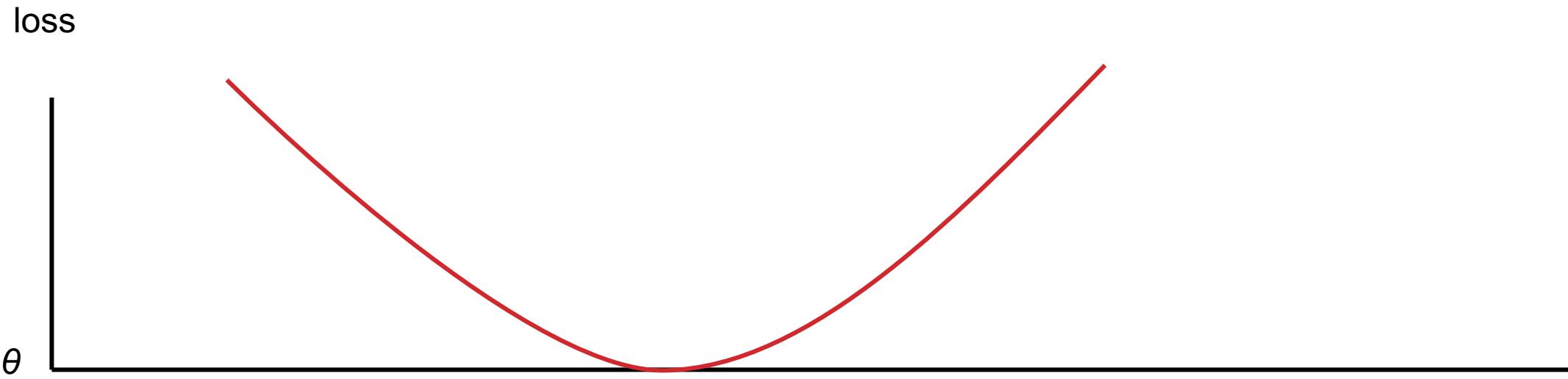
This is NP-hard.

$$\arg \min_{\theta} \sum_{i=1}^n \mathbf{1} \left[ y^{(i)} \cdot \langle \theta, x^{(i)} \rangle \leq 0 \right]$$

- Small changes in any  $\theta$  can have large changes in the loss (the change isn't continuous)
- There can be many local minima
- At any give point, we don't have much information to direct us towards any minima

**Maybe we should consider other loss functions.**

# DESIRABLE PROPERTIES



What are some desirable properties of a loss function???????

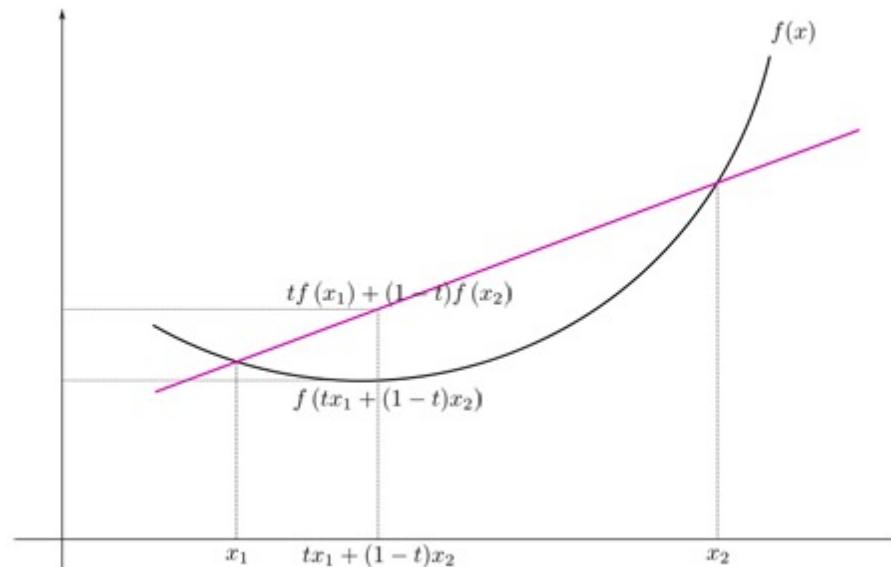
- **Continuous** so we get a local indication of the direction of minimization
- **Only one (i.e., global) minimum**

# CONVEX FUNCTIONS

**“A function is convex if the line segment between any two points on its graph lies above it.”**

**Formally, given function  $f$  and two points  $x, y$ :**

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad \forall \lambda \in [0, 1]$$



# SURROGATE LOSS FUNCTIONS

For many applications, we really would like to minimize the 0/1 loss

A **surrogate loss function** is a loss function that provides an upper bound on the actual loss function (in this case, 0/1)

We'd like to identify **convex** surrogate loss functions to make them easier to minimize

Key to a loss function is how it scores the difference between the actual label  $y$  and the predicted label  $y'$

# SURROGATE LOSS FUNCTIONS

0/1 loss:  $\ell(\hat{y}, y) = \mathbf{1} [y\hat{y} \leq 0]$

Any ideas for surrogate loss functions ??????????

Want: a function that is continuous and convex and is **consistent with the 0/1 loss**  
– that is, minimizing the surrogate loss also minimizes the 0/1 loss

- **Hinge:**  $\ell(\hat{y}, y) = \max(0, 1 - y\hat{y})$
- **Exponential:**  $\ell(\hat{y}, y) = e^{-y\hat{y}}$
- **Squared:**  $\ell(\hat{y}, y) = (y - \hat{y})^2$

What do each of these penalize?????????

# SURROGATE LOSS FUNCTIONS

0/1 loss:

$$\ell(\hat{y}, y) = \mathbf{1} [y\hat{y} \leq 0]$$

Hinge:

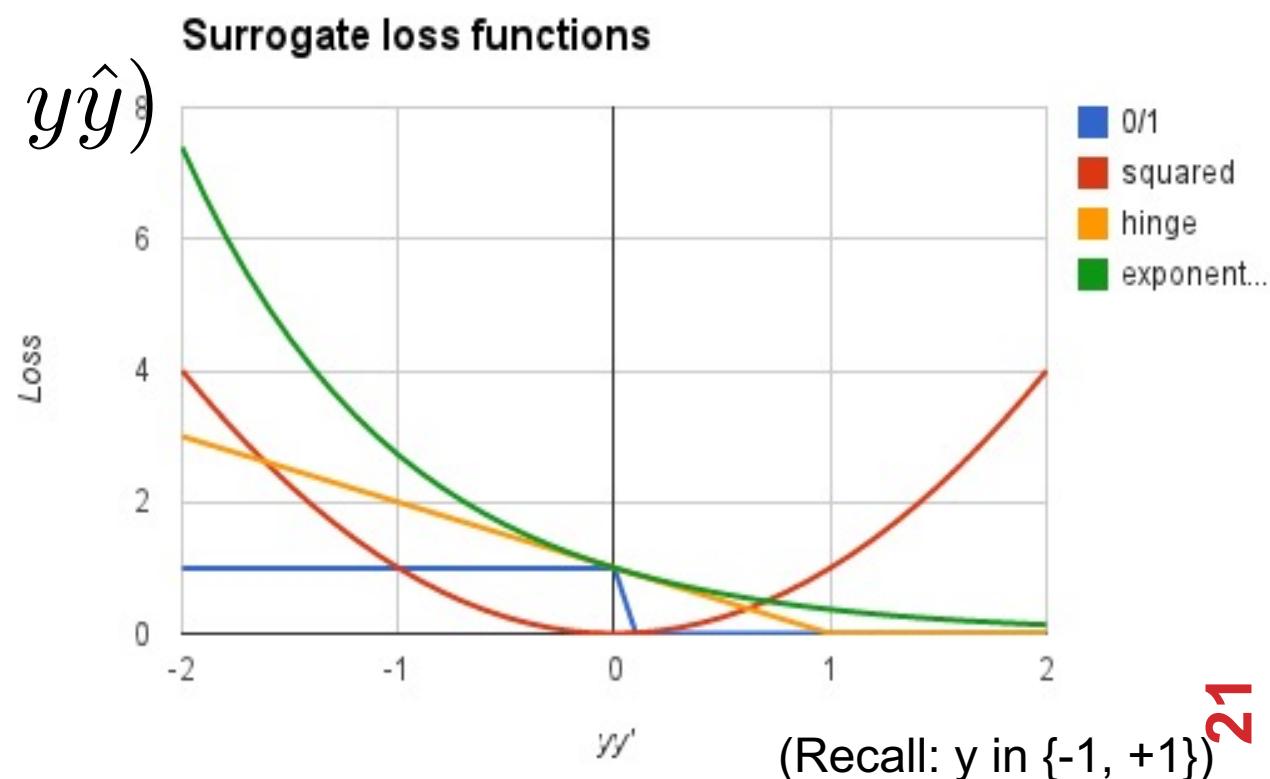
$$\ell(\hat{y}, y) = \max(0, 1 - y\hat{y})$$

Exponential:

$$\ell(\hat{y}, y) = e^{-y\hat{y}}$$

Squared loss:

$$\ell(\hat{y}, y) = (y - \hat{y})^2$$

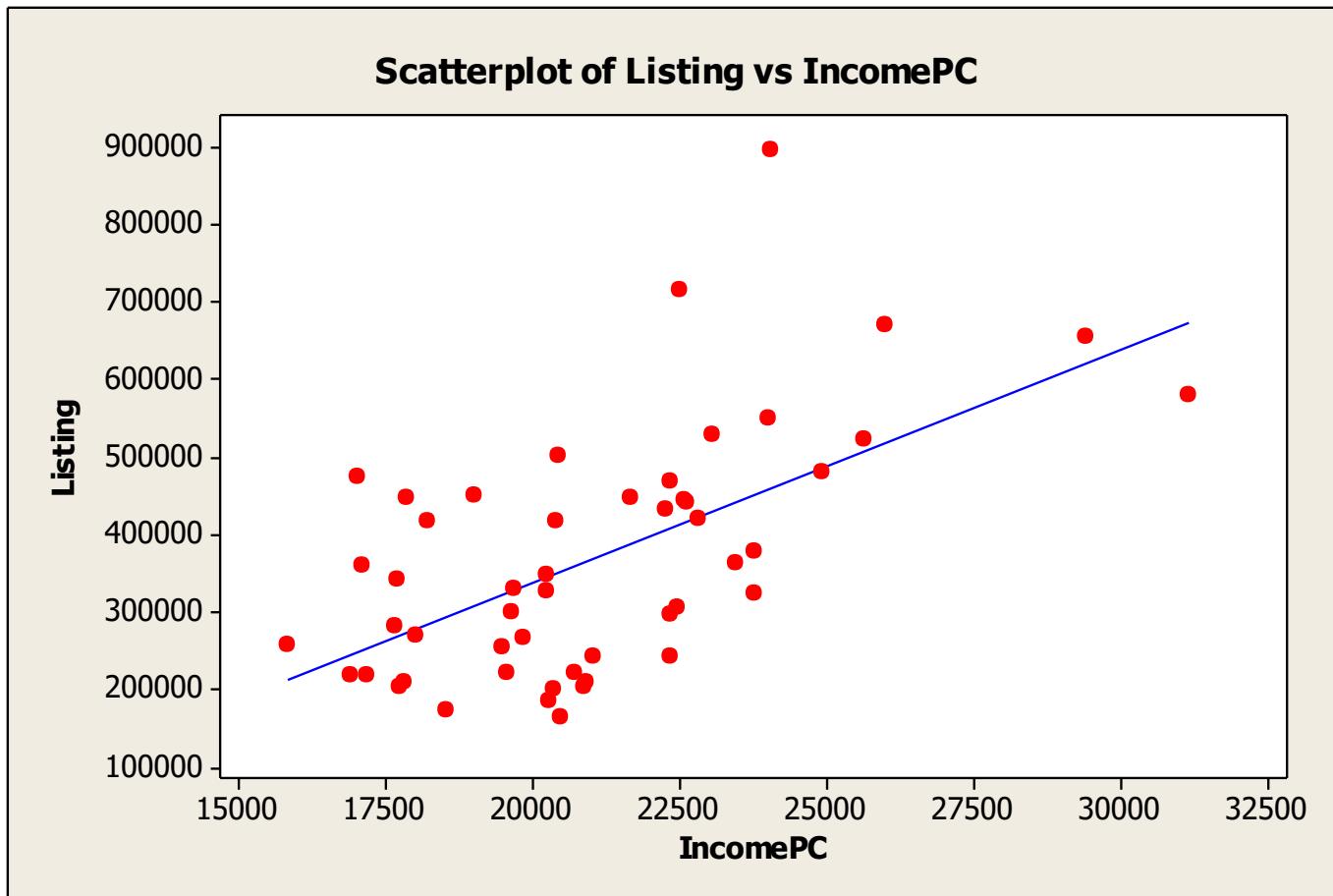


# SOME ML ALGORITHMS

Name	Hypothesis Function	Loss Function	Optimization Approach
Least squares	Linear	Squared	Analytical or GD
Linear regression	Linear	Squared	Analytical or GD
Support Vector Machine (SVM)	Linear, Kernel	Hinge	Analytical or GD
Perceptron	Linear	Perceptron criterion (~Hinge)	Perceptron algorithm, others
Neural Networks	Composed nonlinear	Squared, Hinge, Cross Ent, ...	SGD
Decision Trees	Hierarchical halfplanes	Many	Greedy
Naïve Bayes	Linear	Joint probability	#SAT

**EXAMPLE**

# RECALL: LINEAR REGRESSION



# LINEAR REGRESSION AS MACHINE LEARNING

Let's consider linear regression that minimizes the sum of squared error, i.e., least squares ...

1. Hypothesis function: ????????

- Linear hypothesis function

$$h_{\theta}(x) = \theta^T x$$

2. Loss function: ????????

- Squared error loss

$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

3. Optimization problem: ????????

$$\text{minimize}_{\theta} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

# LINEAR REGRESSION AS MACHINE LEARNING

Rewrite inputs:

Each row is a feature vector paired with a label for a single input

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

*m labeled inputs*

*n features*

Rewrite optimization problem:

$$\text{minimize}_{\theta} \frac{1}{2} \|X\theta - y\|_2^2$$

\*Recall:  $\|x\|_2^2 = z^T z = \sum_i z_i^2$

# GRADIENTS

In Lecture 10, we showed that the mean is the point that minimizes the residual sum of squares:

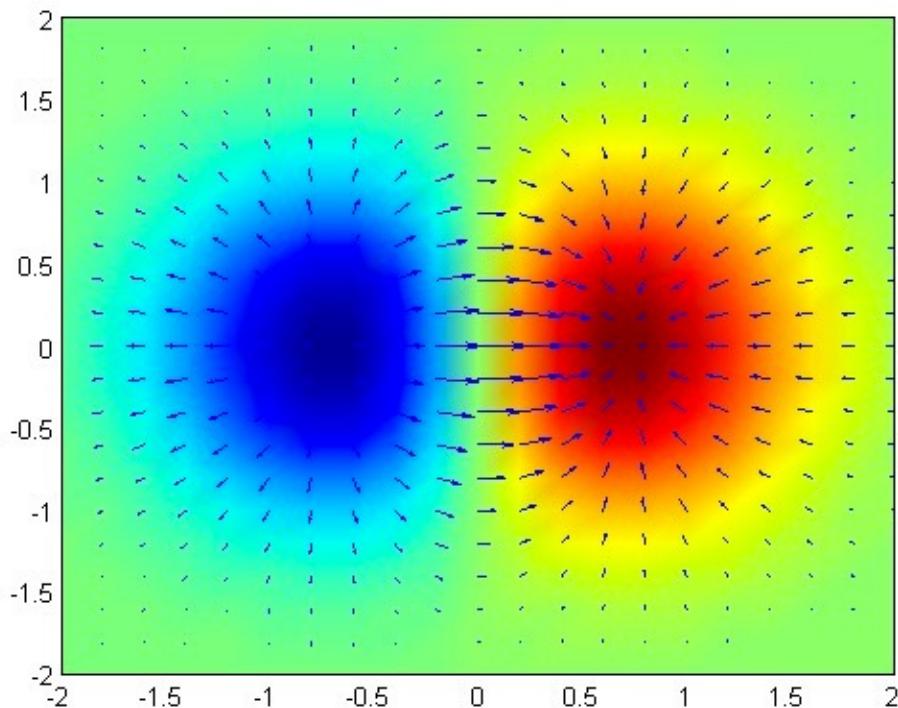
- Solved minimization by finding point where derivative is zero
- (Convex functions like RSS → single global minimum.)

The **gradient** is the multivariate generalization of a derivative.

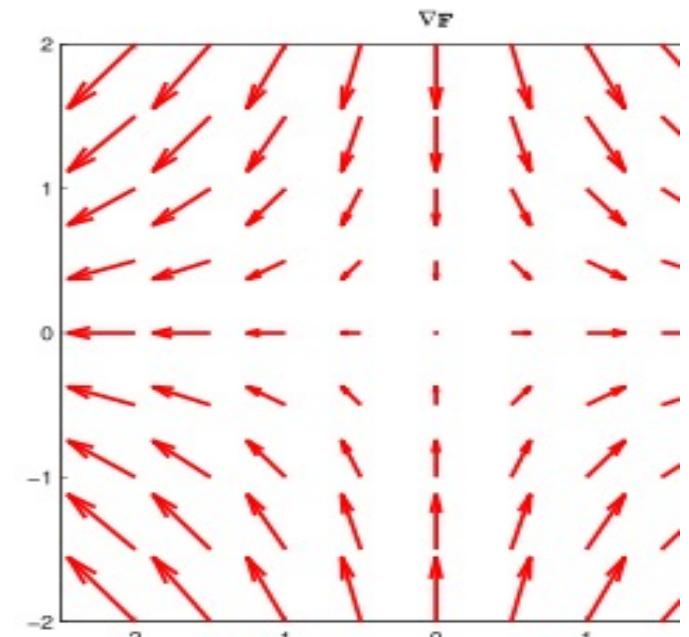
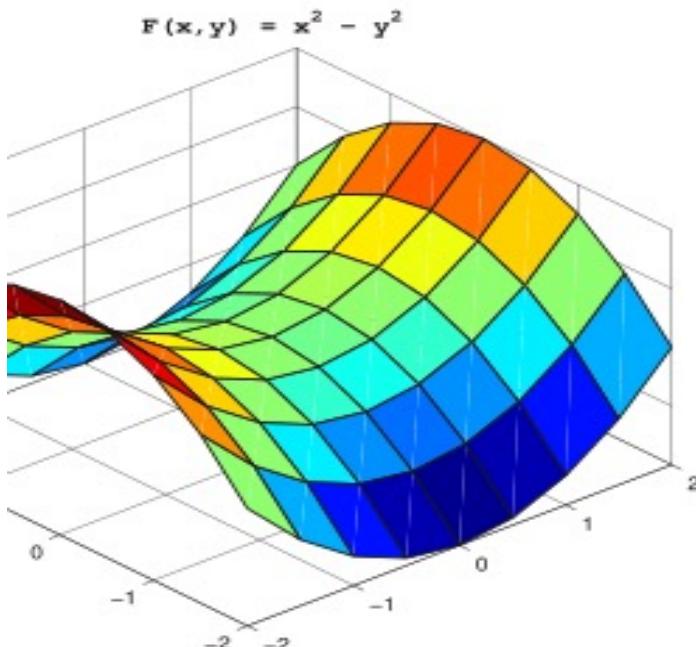
For a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , the gradient is a vector of all  $n$  partial derivatives:

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_n} \end{bmatrix} \in \mathbb{R}^n$$

# GRADIENTS



Gradient of  $f(x,y) = xe^{-(x^2 + y^2)}$



# GRADIENTS

Minimizing a multivariate function involves finding a point where the gradient is zero:

$$\nabla_{\theta} f(\theta) = 0 \text{ (the vector of zeros)}$$

Points where the gradient is zero are **local** minima

- If the function is convex, also a **global** minimum

Let's solve the least squares problem!

We'll use the multivariate generalizations of some concepts from MATH141/142 ...

- Chain rule:  $\nabla_{\theta} f(X\theta) = X^T \nabla_{X\theta} f(X\theta)$
- Gradient of squared  $\ell^2$  norm:  $\nabla_{\theta} \|\theta - z\|_2^2 = 2(\theta - z)$

# LEAST SQUARES

Recall the least squares optimization problem:

$$\text{minimize}_{\theta} \frac{1}{2} \|X\theta - y\|_2^2$$

What is the gradient of the optimization objective ?????

$$\nabla_{\theta} \frac{1}{2} \|X\theta - y\|_2^2 =$$

Chain rule:  
 $\nabla_{\theta} f(X\theta) = X^T \nabla_{X\theta} f(X\theta)$

$$X^T \nabla_{X\theta} \frac{1}{2} \|X\theta - y\|_2^2 =$$

Gradient of norm:  
 $\nabla_{\theta} \|\theta - z\|_2^2 = 2(\theta - z)$

$$\nabla_{\theta} \frac{1}{2} \|X\theta - y\|_2^2 = X^T(X\theta - y)$$

# LEAST SQUARES

Recall: points where the gradient **equals zero** are minima.

$$\nabla_{\theta} \frac{1}{2} \|X\theta - y\|_2^2 = X^T(X\theta - y)$$

So where do we go from here?????????

$$X^T(X\theta - y) = 0 \qquad \text{Solve for model parameters } \theta$$

$$X^T X\theta - X^T y = 0 \rightarrow X^T X\theta = X^T y$$

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

$$\boxed{\theta = (X^T X)^{-1} X^T y}$$

# ML IN PYTHON



**Python has tons of hooks into a variety of machine learning libraries. (Part of why this course is taught in Python!)**

**Scikit-learn is the most well-known (non-deep-learning) library:**

- Classification (SVN, K-NN, Random Forests, ...)
- Regression (SVR, Ridge, Lasso, ...)
- Clustering (k-Means, spectral, mean-shift, ...)
- Dimensionality reduction (PCA, matrix factorization, ...)
- Model selection (grid search, cross validation, ...)
- Preprocessing (cleaning, EDA, ...)
- Neural nets and some deep learning, but not the right library for this

**Built on the NumPy stack; plays well with Matplotlib.**

# LEAST SQUARES IN PYTHON

You don't need Scikit-learn for OLS ...

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

```
# Analytic solution to OLS using Numpy  
params = np.linalg.solve(X.T.dot(X), X.T.dot(y))
```

But let's say you did want to use it.

```
from sklearn import linear_model  
  
X = [[0,0], [1,1], [2,2]]  
Y = [0, 1, 2]  
  
# Solve OLS using Scikit-Learn  
reg = linear_model.LinearRegression()  
reg.fit(X, Y)  
reg.coef_
```

```
array([ 0.5,  0.5])
```

*NEXT UP:*

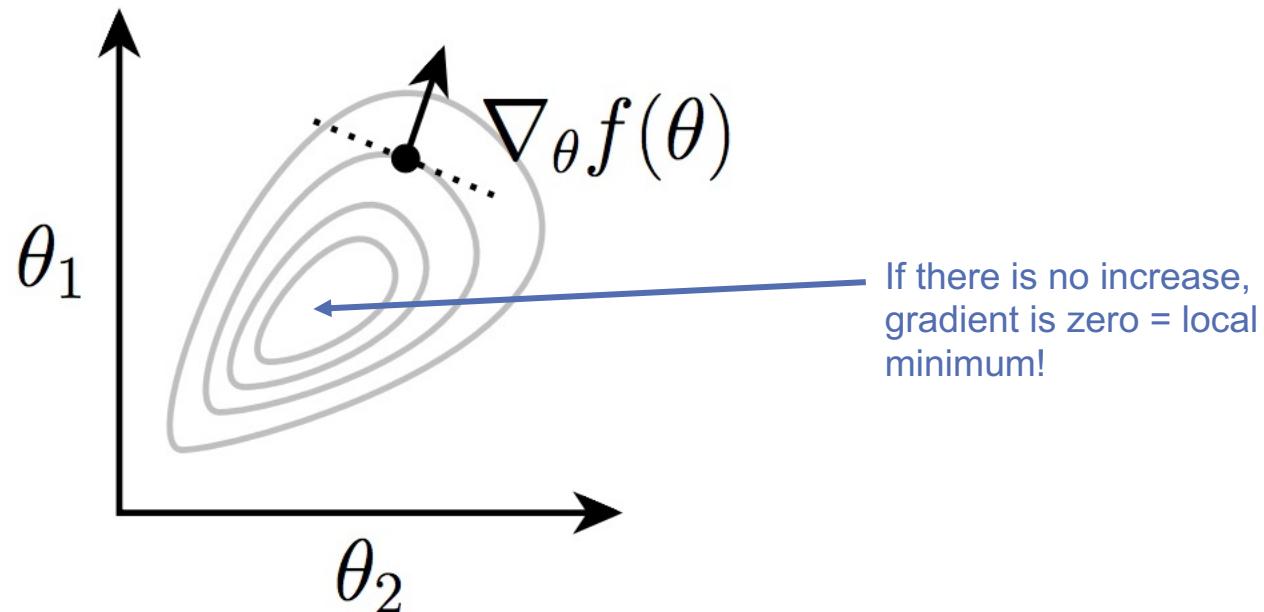
**(STOCHASTIC)  
GRADIENT DESCENT**



# GRADIENT DESCENT

We used the gradient as a condition for optimality

It also gives the local **direction of steepest increase** for a function:



Intuitive idea: take small steps **against** the gradient.

# GRADIENT DESCENT

Algorithm for any\* hypothesis function  $h_\theta: \mathbb{R}^n \rightarrow \mathcal{Y}$ , loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ , step size  $\alpha$ :

Initialize the parameter vector:

- $\theta \leftarrow 0$

Repeat until satisfied (e.g., exact or approximate convergence):

- Compute gradient:  $g \leftarrow \sum_{i=1}^m \nabla_\theta \ell(h_\theta(x^{(i)}), y^{(i)})$
- Update parameters:  $\theta \leftarrow \theta - \alpha \cdot g$

\*must be reasonably well behaved

# GRADIENT DESCENT

**Step-size ( $\alpha$ ) is an important parameter**

- Too large → might oscillate around the minima
- Too small → can take a long time to converge

**If there are no local minima, then the algorithm eventually converges to the optimal solution**

**Very widely used in Machine Learning**

# EXAMPLE

Function:  $f(x,y) = x^2 + 2y^2$

Gradient: ??????????

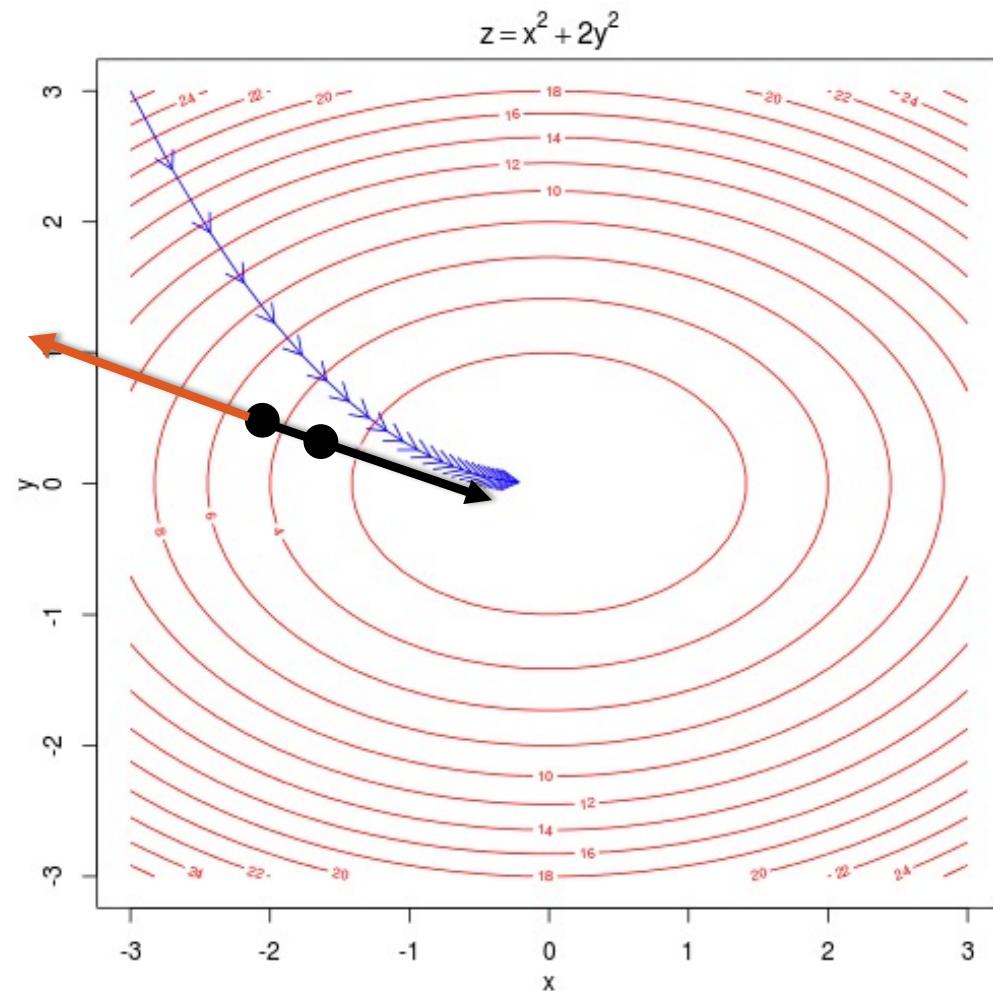
$$\nabla f(x, y) = \begin{bmatrix} 2x \\ 4y \end{bmatrix}$$

Let's take a gradient step from  $(-2, +1/2)$ :

$$\nabla f(-2, 1) = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$

Step in the direction  $(+4, -2)$ , scaled by step size

Repeat until no movement



# A simple example: predicting electricity use

What will peak power consumption be in Pittsburgh tomorrow?

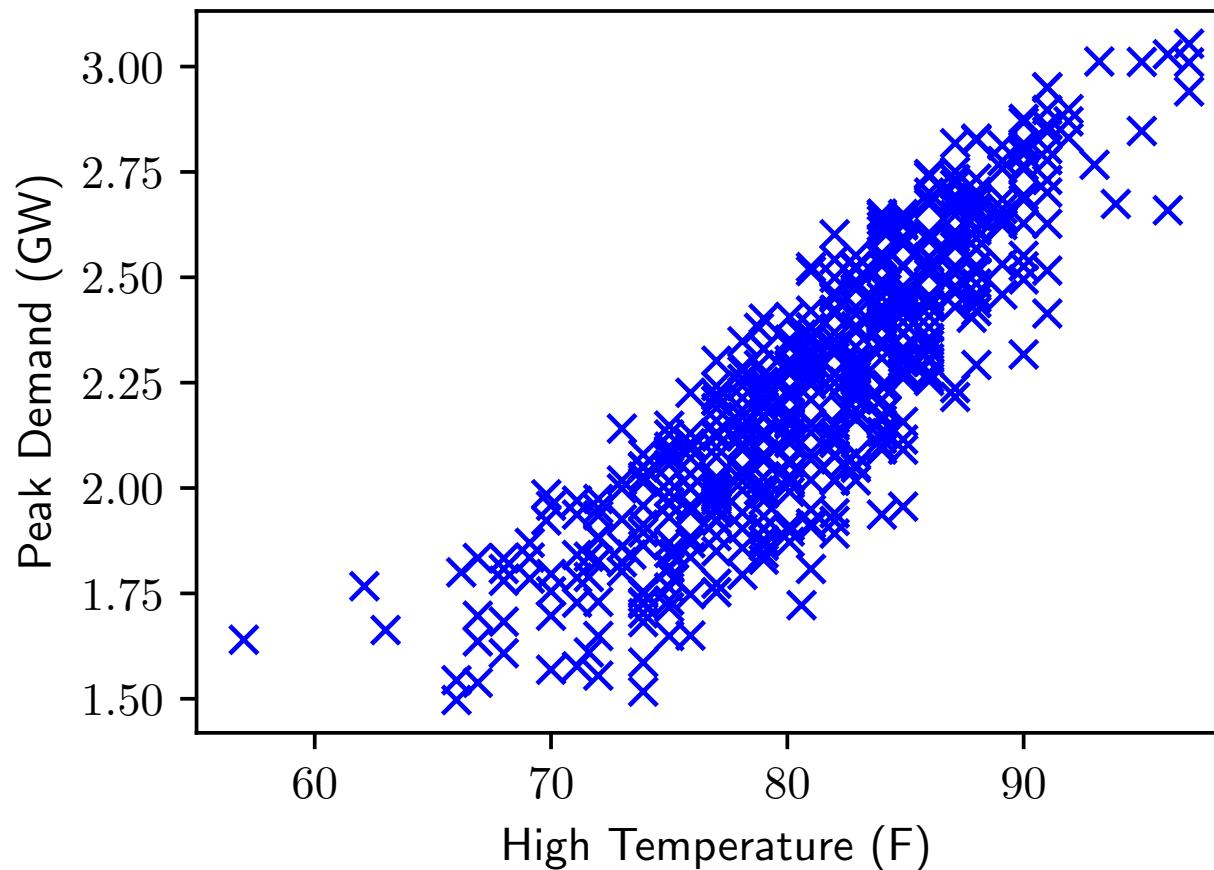
Difficult to build an “a priori” model from first principles to answer this question

But, relatively easy to record past days of consumption, plus additional features that affect consumption (i.e., weather)

Date	High Temperature (F)	Peak Demand (GW)
2011-06-01	84.0	2.651
2011-06-02	73.0	2.081
2011-06-03	75.2	1.844
2011-06-04	84.9	1.959
...	...	...

# Plot of consumption vs. temperature

Plot of high temperature vs. peak demand for summer months (June – August) for past six years



# Hypothesis: linear model

Let's suppose that the peak demand approximately fits a *linear model*

$$\text{Peak_Demand} \approx \theta_1 \cdot \text{High_Temperature} + \theta_2$$

Here  $\theta_1$  is the “slope” of the line, and  $\theta_2$  is the intercept

How do we find a “good” fit to the data?

Many possibilities, but natural objective is to minimize some difference between this line and the observed data, e.g. squared loss

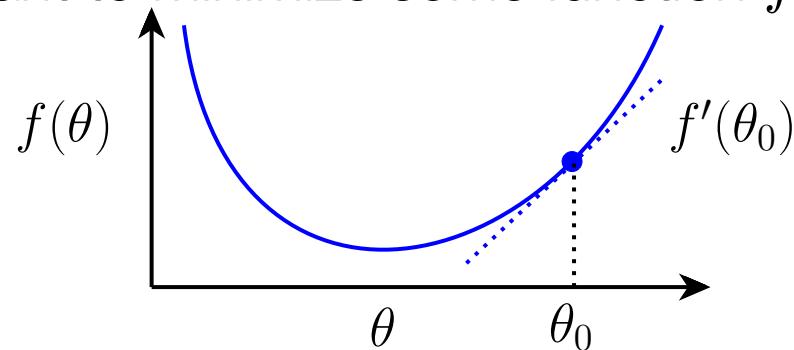
$$E(\theta) = \sum_{i \in \text{days}} (\theta_1 \cdot \text{High_Temperature}^{(i)} + \theta_2 - \text{Peak_Demand}^{(i)})^2$$

# How do we find parameters?

How do we find the parameters  $\theta_1, \theta_2$  that minimize the function

$$\begin{aligned} E(\theta) &= \sum_{i \in \text{days}} (\theta_1 \cdot \text{High_Temperature}^{(i)} + \theta_2 - \text{Peak_Demand}^{(i)})^2 \\ &\equiv \sum_{i \in \text{days}} (\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)})^2 \end{aligned}$$

General idea: suppose we want to minimize some function  $f(\theta)$



Derivative is slope of the function, so negative derivative points “downhill”

# Computing the derivatives

What are the derivatives of the error function with respect to each parameter  $\theta_1$  and  $\theta_2$ ?

$$\begin{aligned}\frac{\partial E(\theta)}{\partial \theta_1} &= \frac{\partial}{\partial \theta_1} \sum_{i \in \text{days}} (\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)})^2 \\ &= \sum_{i \in \text{days}} \frac{\partial}{\partial \theta_1} (\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)})^2 \\ &= \sum_{i \in \text{days}} 2(\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)}) \cdot \frac{\partial}{\partial \theta_1} \theta_1 \cdot x^{(i)} \\ &= \sum_{i \in \text{days}} 2(\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)}) \cdot x^{(i)} \\ \frac{\partial E(\theta)}{\partial \theta_2} &= \sum_{i \in \text{days}} 2(\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)})\end{aligned}$$

# Finding the best $\theta$

To find a good value of  $\theta$ , we can repeatedly take steps in the direction of the negative derivatives for each value

Repeat:

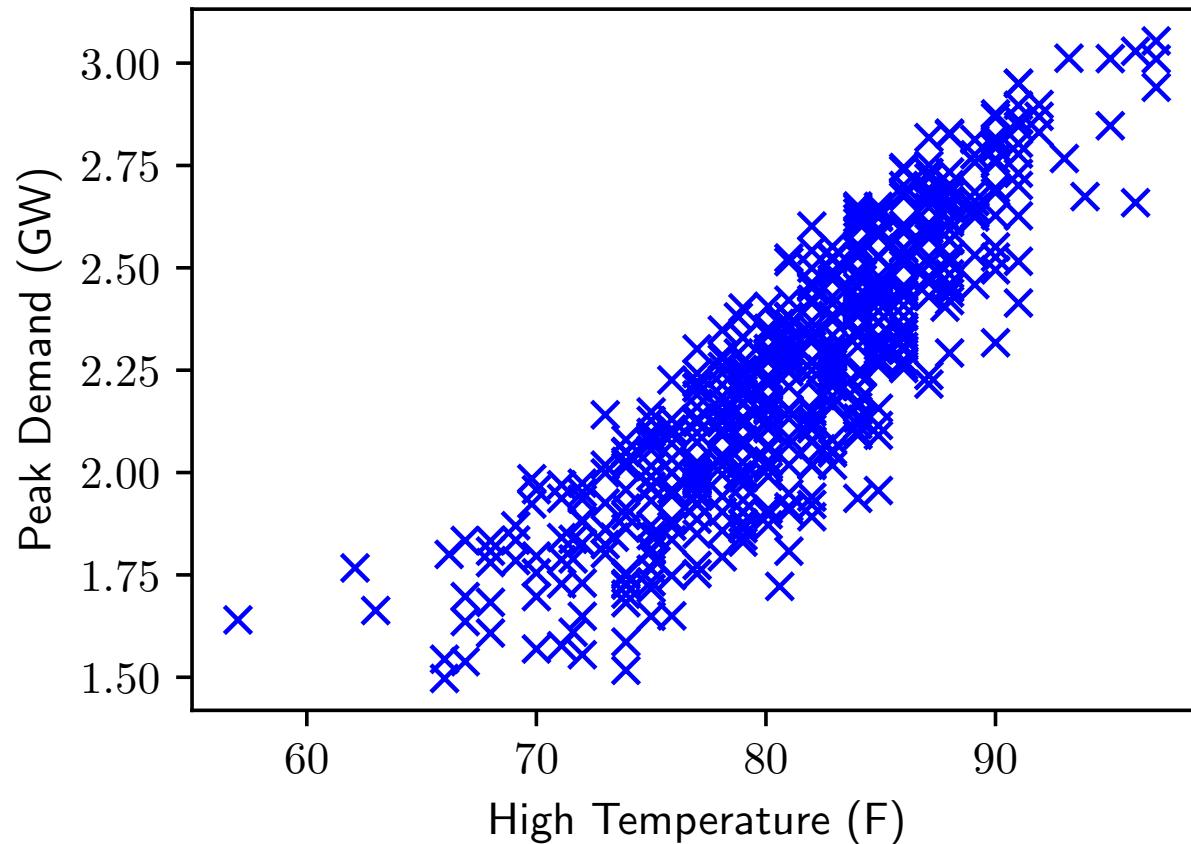
$$\theta_1 := \theta_1 - \alpha \sum_{i \in \text{days}} 2(\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)}) \cdot x^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \sum_{i \in \text{days}} 2(\theta_1 \cdot x^{(i)} + \theta_2 - y^{(i)})$$

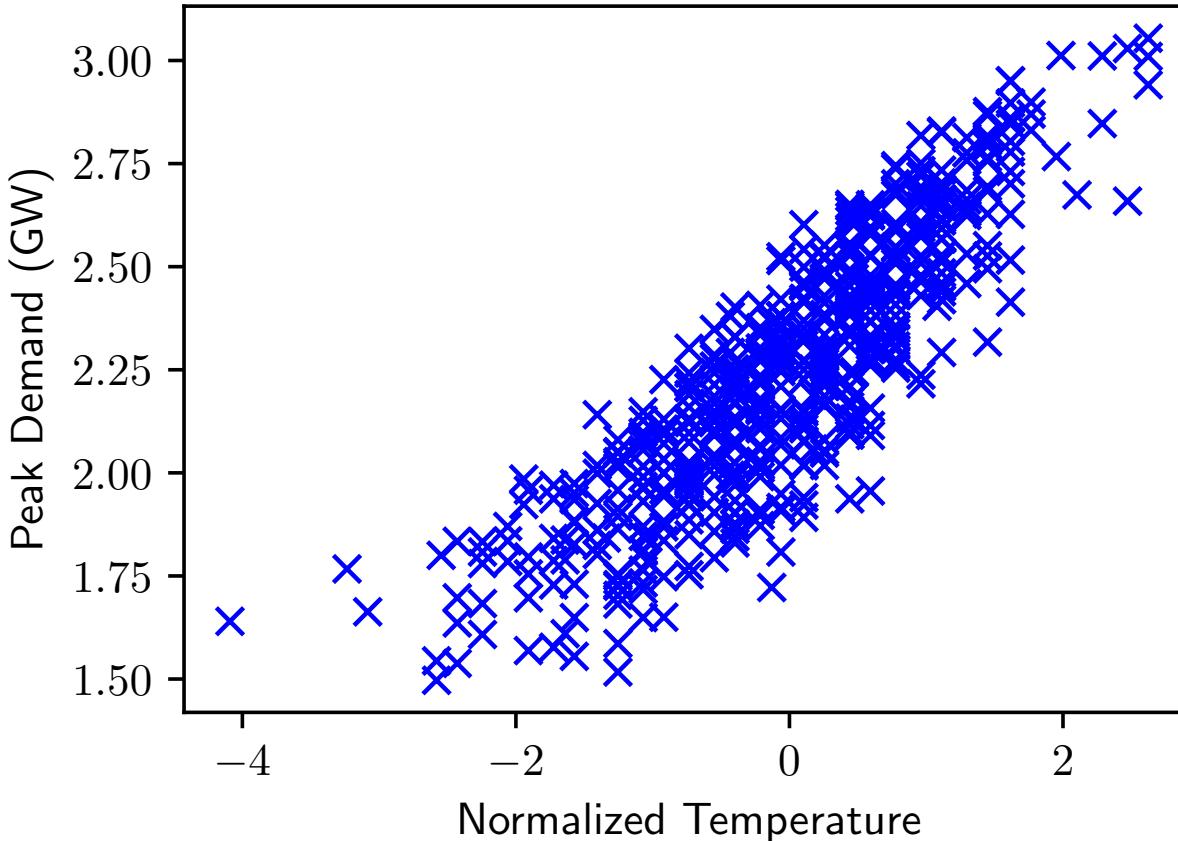
where  $\alpha$  is some small positive number called the step size

This is the *gradient decent algorithm*, the workhorse of modern machine learning

# Gradient descent

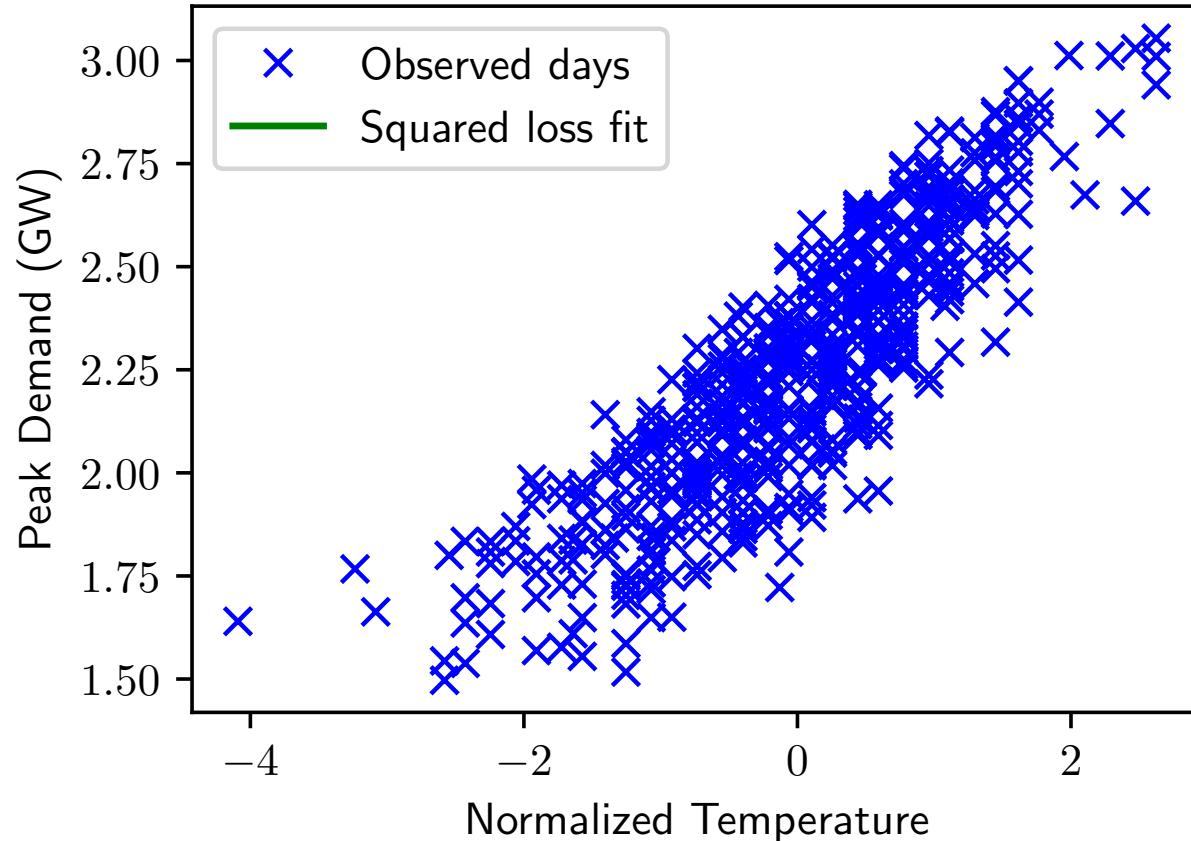


# Gradient descent



**Normalize** input by subtracting the mean and dividing by the standard deviation

# Gradient descent – Iteration 1

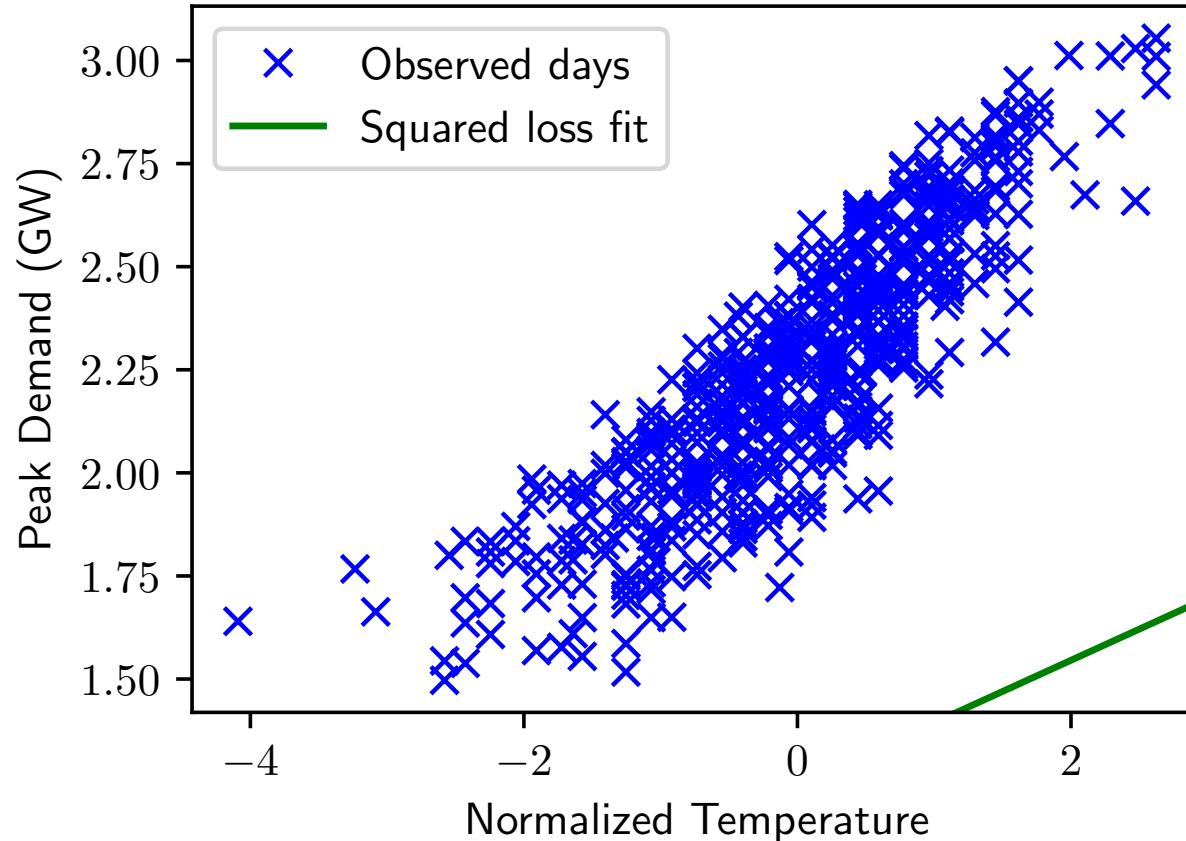


$$\theta = (0.00, 0.00)$$

$$E(\theta) = 1427.53$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-151.20, -1243.10)$$

# Gradient descent – Iteration 2

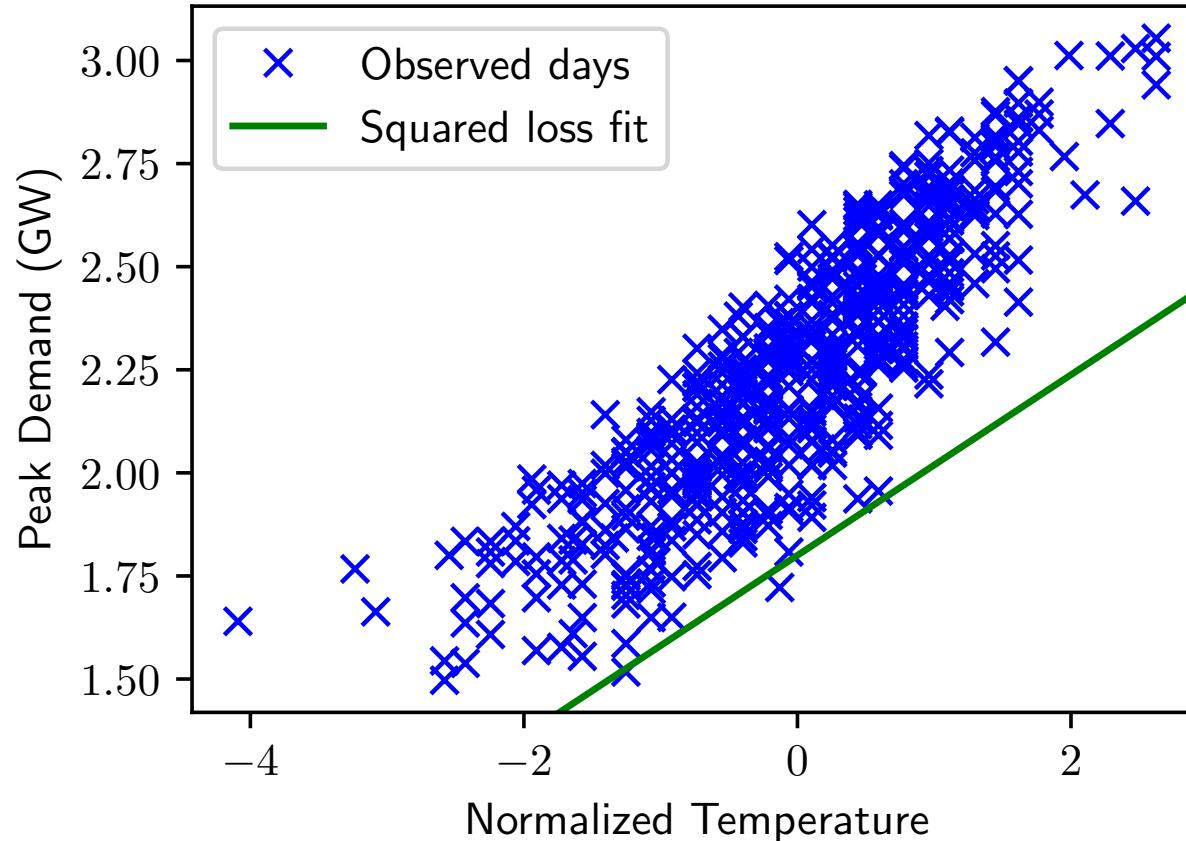


$$\theta = (0.15, 1.24)$$

$$E(\theta) = 292.18$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-67.74, -556.91)$$

# Gradient descent – Iteration 3

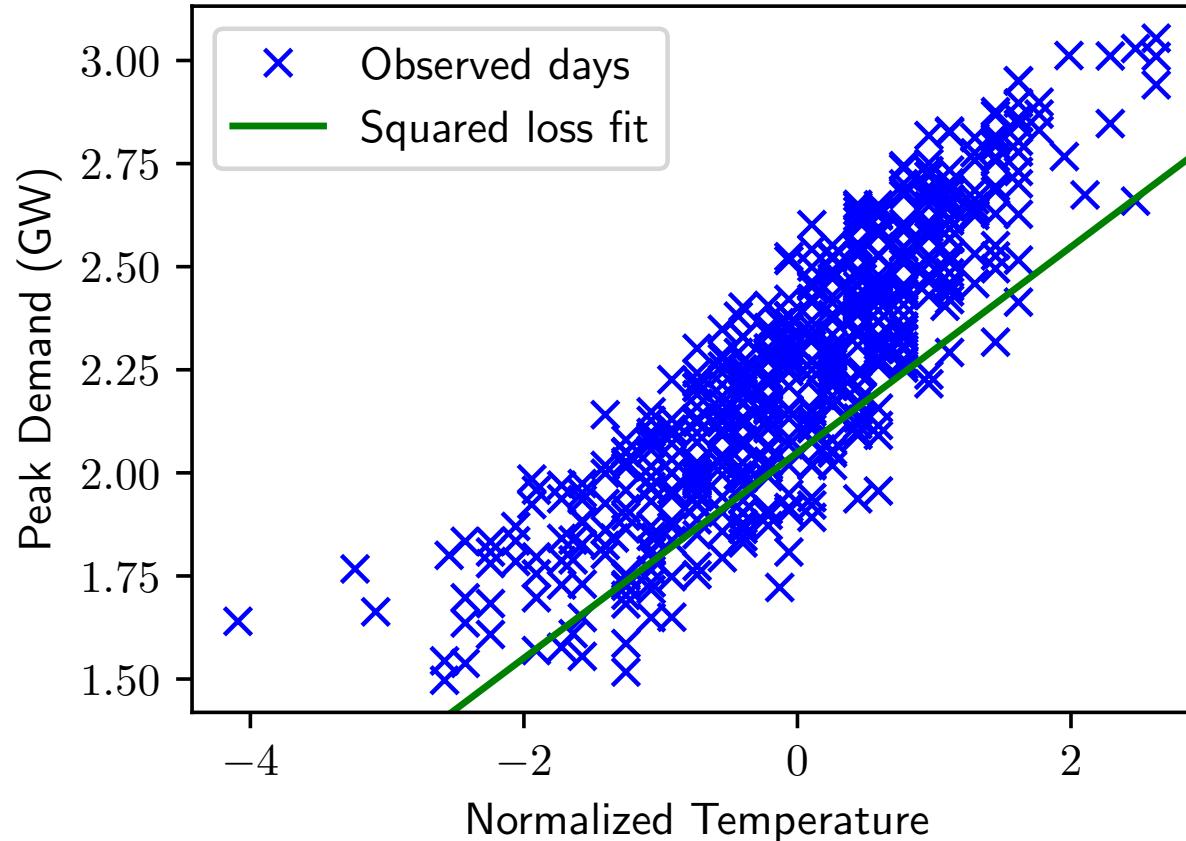


$$\theta = (0.22, 1.80)$$

$$E(\theta) = 64.31$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-30.35, -249.50)$$

# Gradient descent – Iteration 4

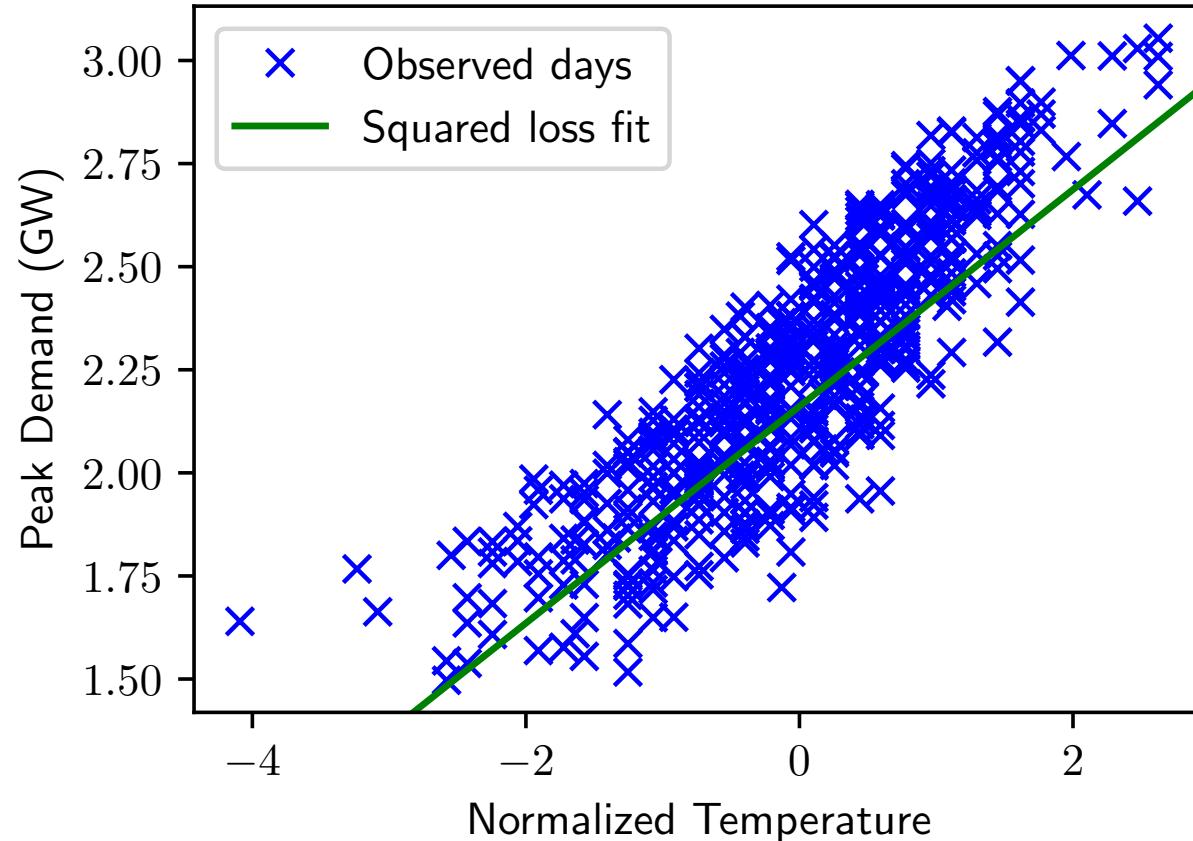


$$\theta = (0.25, 2.05)$$

$$E(\theta) = 18.58$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-13.60, -111.77)$$

# Gradient descent – Iteration 5

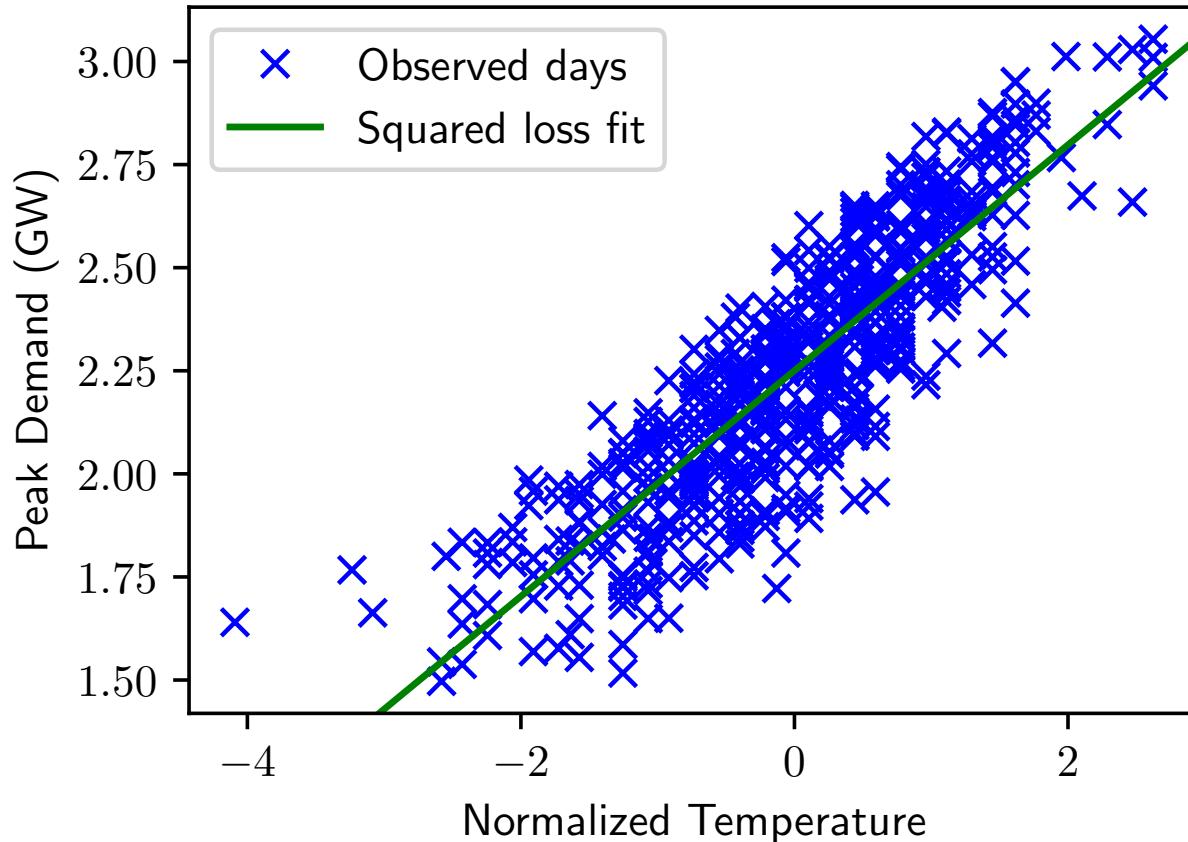


$$\theta = (0.26, 2.16)$$

$$E(\theta) = 9.40$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-6.09, -50.07)$$

# Gradient descent – Iteration 10

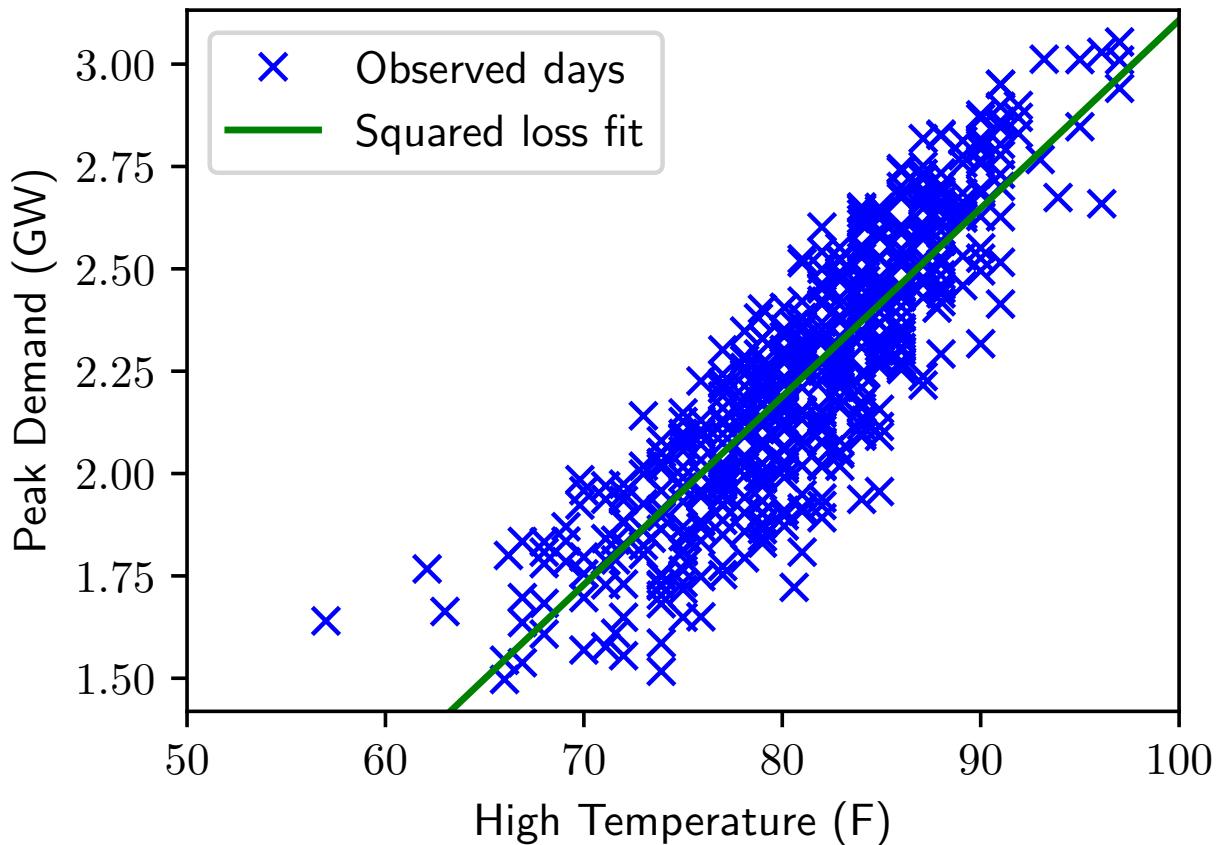


$$\theta = (0.27, 2.25)$$

$$E(\theta) = 7.09$$

$$\left( \frac{\partial E(\theta)}{\partial \theta_1}, \frac{\partial E(\theta)}{\partial \theta_2} \right) = (-0.11, -0.90)$$

# Fitted line in “original” coordinates



# Making predictions

Importantly, our model also lets us make *predictions* about new days

What will the peak demand be tomorrow?

If we know the high temperature will be 72 degrees (ignoring for now that this is also a prediction), then we can predict peak demand to be:

$$\text{Predicted\_demand} = \theta_1 \cdot 72 + \theta_2 = 1.821 \text{ GW}$$

(requires that we rescale  $\theta$  after solving to “normal” coordinates)

Equivalent to just “finding the point on the line”

# IN GENERAL: GRADIENT DESCENT FOR OLS

Algorithm for linear hypothesis function and squared error loss function (combined to  $1/2\|X\theta - y\|_2^2$ , like before):

Initialize the parameter vector:

- $\theta \leftarrow 0$

Repeat until satisfied:

- Compute gradient:  $g \leftarrow X^T(X\theta - y)$
- Update parameters:  $\theta \leftarrow \theta - \alpha \cdot g$

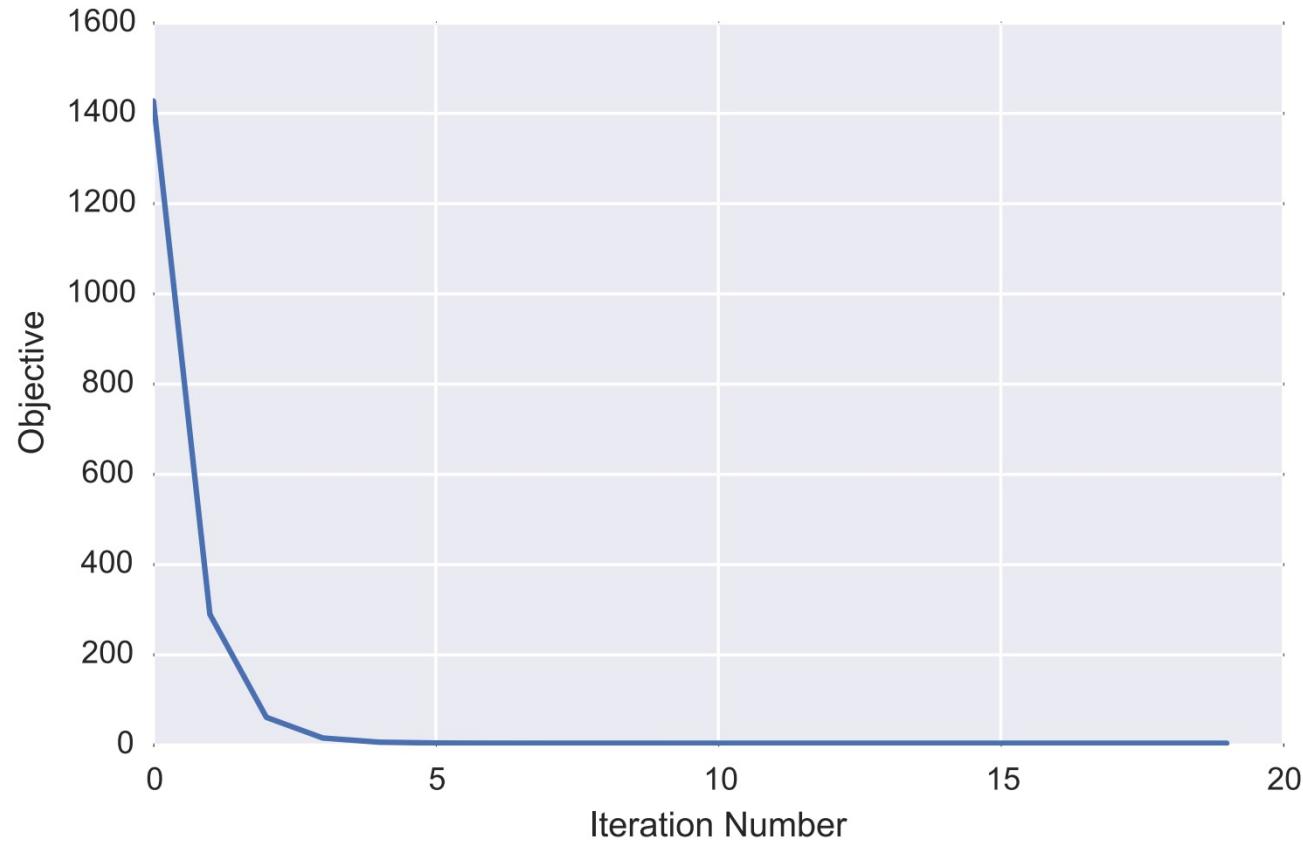
# GRADIENT DESCENT IN PURE(-ISH) PYTHON

```
# Training data (X, y), T time steps, alpha step
def grad_descent(X, y, T, alpha):
    m, n = X.shape           # m = #examples, n = #features
    theta = np.zeros(n)       # initialize parameters
    f = np.zeros(T)           # track loss over time

    for i in range(T):
        # loss for current parameter vector theta
        f[i] = 0.5*np.linalg.norm(X.dot(theta) - y)**2
        # compute steepest ascent at f(theta)
        g = X.T.dot(X.dot(theta) - y)
        # step down the gradient
        theta = theta - alpha*g
    return theta, f
```

Implicitly using squared loss and linear hypothesis function above; drop in your favorite gradient for kicks!

# PLOTTING LOSS OVER TIME



Why ????????

# ITERATIVE VS ANALYTIC SOLUTIONS

But we already had an analytic solution! What gives?

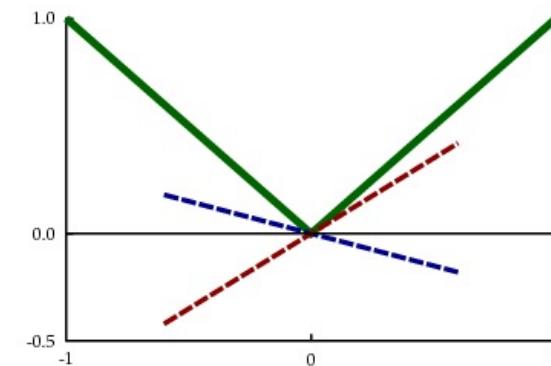
Recall: last class we discuss 0/1 loss, and using convex surrogate loss functions for tractability

One such function, the **absolute error loss function**, leads to:

$$\underset{\theta}{\text{minimize}} \sum_{i=1}^m |\theta^T x^{(i)} - y^{(i)}| \equiv \underset{\theta}{\text{minimize}} \|X\theta - y\|_1$$

Problems ????????

- Not differentiable! But subgradients?
- No closed form!
- So you **must** use iterative method



# LEAST ABSOLUTE DEVIATIONS

Can solve this using gradient descent and the gradient:

$$\nabla_{\theta} \|X\theta - y\|_1 = X^T \text{sign}(X\theta - y)$$

Simple to change in our Python code:

```
for i in range(T):
    # loss for current parameter vector theta
    f[i] = np.linalg.norm(X.dot(theta) - y, 1)
    # compute steepest ascent at f(theta)
    g = X.T.dot( np.sign(X.dot(theta) - y) )
    # step down the gradient
    theta = theta - alpha*g
return theta, f
```

# BATCH VS STOCHASTIC GRADIENT DESCENT

**Batch:** Compute a single gradient (vector) for the entire dataset (as we did so far)

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \quad (\text{for every } j).$$

**Incremental/Stochastic:**

- Do one training sample at a time, i.e., update parameters for every sample separately
- Much faster in general, with more pathological cases

Loop {

    for i=1 to m, {

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} \quad (\text{for every } j).$$

}

}