# CS-GY 6923: Lecture 1
# Introduction to Machine Learning

NYU Tandon School of Engineering, Prof. Christopher Musco

Who has tried ChatGPT? DALLE? Imagen?



Edit the detailed description  Surprise me  Upload  →|

Paint Washington Square Park in a snow storm in the style of Vincent Van Gogh.  **Generate**

## Who has tried ChatGPT? DALLE? Imagen?

Can you write a bash script to resize all images in a directory to be less than 100kb?

Yes, here is an example bash script that uses the `mogrify` command from the `imagemagick` package to resize all images in a directory to be less than 100KB:

```bash
#!/bin/bash

# specify the directory containing the images
dir="path/to/images"

# resize all images in the directory
mogrify -path $dir -resize x800 -quality 50 -define jpeg:extent=100kb
-format jpg $dir/*
```

This script takes a directory specified by the variable `dir` and uses the `mogrify` command to resize all images in the directory with the following options:

May look back on 2015-now as a the dawn of true artificial intelligence.



*The New York Times*

Showing 77 results for:

## chatgpt

Sort by Relevance

Date Range ∨          Section ∨          Type ∨

---

Dec. 5, 2022

TECHNOLOGY

### The Brilliance and Weirdness of ChatGPT

A new chatbot from OpenAI is inspiring awe, fear, stunts and attempts to circumvent its guardrails.

By Kevin Roose

PRINT EDITION  The Brilliance And Weirdness Of ChatGPT | December 9, 2022, Page B1

---

Dec. 9, 2022

PODCASTS

### Can ChatGPT Make This Podcast?

OpenAI's new chatbot is a coder, a teacher, a potential Google killer, and … a friend?

By Kevin Roose, Casey Newton, Davis Land, Paula Szuchman, Dan Powell, Elisheba Ittoop, Rowan Niemisto and Alyssa Moxley

---

Jan. 15

OPINION

### How ChatGPT Hijacks Democracy

Other developments in recent years:

- Human-level image classification and understanding.
- Near perfect machine translation.
- Human level game play in very complex games (Go, Starcraft).
- Machine learning as a central tool in science.

<span style="color:orange">What technologies have caught people's eye?</span>

Give you a <u>foundation</u> to understand the main ideas in modern machine learning.

We will do so through a combination of:

- Hands on implementation.
    - Demos and take-home labs using `Python` and `Jupyter notebooks`. 20% of grade
    - We will use `Google Colab` as the primary programming environment.
- Theoretical exploration.
    - Written problem sets. 20%
    - Midterm and final exam. 25% of grade each.

## COURSE OBJECTIVES

### Goals of theoretical component:

1. Build experience with the most important mathematical tools used in machine learning, including probability, statistics, and linear algebra. This experience will prepare you for more advanced coursework in ML, or research.

2. Be able to understand contemporary research in machine learning, including papers from NeurIPS, ICML, ICLR, and other major machine learning venues.

3. Learn how theoretical analysis can help explain the performance of machine learning algorithms and lead to the design of entirely new methods.

Goals of hands-on component:

1. Reinforce theory learned in class, and make sure you understand algorithms described by implementing them.
2. Learn how to view and formulate real world problems in the language of machine learning.
3. Gain experience applying the most popular and successful machine learning algorithms to these problems.

- CS-GY 6953: **Deep Learning** (Prof. Chinmay Hegde)
- ECE-GY 7143: **Advanced Machine Learning** (Prof. Anna Chromanska)
- CS-GY 6763: **Algorithmic Machine Learning and Data Science** (me)
- Keep your eyes out for special topics courses.

Given recent progress, we will run this semester of the course as a special semester focused on generative machine learning.

- Still cover all the basics.
- Add material in each section on background for topics like text and image generation, style transfer, etc.

All class information can be found at:

www.chrismusco.com/machinelearning2023_grad

1. Make sure you are signed into and follow EdStem, which will be used for all classroom communication (no email). Now integrated into Brightspace.
2. Don't hesitate to ask me or the TAs for help.[1]



Aarshvi Gajjar



Atsushi Shimizu



Teal Witter

---

[1]Fill out office hours poll on Ed!

Class participation accounts for 10% of your grade. It's easy to get a perfect score:

- Ask and answer questions in lecture.
- Post questions or responses to other students on Ed. Or other things you find interesting.
- Participate in professor or TA office hours.

THE PREDICTION PROBLEM

**Goal:** Develop algorithms to make predictions based on data.

- **Input:** A single piece of data (an image, audio file, patient healthcare record, MRI scan).



- **Output:** A prediction (this image is a stop sign, this stock will go up 10% next quarter, this song is in French).

**Optical character recognition (OCR)**: Decide if a handwritten character is an $a, b, \ldots, z, 0, 1, \ldots, 9, \ldots$.

Optical character recognition (OCR): Decide if a handwritten character is an $a, b, \ldots, z, 0, 1, \ldots, 9, \ldots$.

Applications:

- Automatic mail sorting.
- Text search in handwritten documents.
- Digitizing scanned books.
- License plate detection for tolls.
- Etc.

How would you write an **code** to distinguish these digits?



Suppose you just want to distinguish <u>between a 1 and a 7</u>.

## 1S VS. 7S ALGORITHM

Reasonable approach: A number which contains one vertical line is a 1, if it contains one vertical and one horizontal line, it's a 7.

```python
1    def count_vert_lines(image):
2        ...
3
4    def count_horiz_lines(image):
5        ...
6
7    def classify(image):
8        ...
9        nv = count_vert_lines(image)
10       nh = count_vert_lines(image)
11
12       if (nv == 1) and (nh == 1):
13           return '7'
14       elif (nv == 1) and (nh == 0):
15           return '1'
16       elif ...
```

This rule breaks down in practice:



Even fixes/modifications of the rule tend to be brittle... Maybe you could get 80% accuracy, but not nearly good enough.

Rule based systems, also called <u>Expert Systems</u> were <u>the dominant approach</u> to artificial intelligence in the 1970s and 1980s.

Major limitation: While human's are very good at many tasks,

- It's often hard to encode <u>why</u> humans make decisions in simple programmable logic.
- We think in abstract concepts with no mathematical definitions (how exactly do you define a line? how do you define a curve? straight line?)

Focus on what humans do well: solving the task at hand!

**Step 1:** Collect and label many input/output pairs $(x_i, y_i)$. For our digit images, we have each $x_i \in \mathbb{R}^{28 \times 28}$ and $y_i \in \{0, 1, \ldots, 9\}$.



Input data: $(x_1)(x_2)$ ...

8 9 1 2 5 0 0 6 6 4

Output labels: $(y_1)(y_2)$ ..

This is called the training dataset.

Step 2: Learn from the examples we have.

- Have the computer <u>automatically</u> find some function $f(\mathbf{x})$ such that $f(\mathbf{x}_i) = y_i$ for most $(\mathbf{x}_i, y_i)$ in our training data set (by searching over many possible functions).

Think of $f$ as any crazy equation, or an arbitrary program:

$$f(\mathbf{x}) = 10 \cdot x[1,1] - 6 \cdot x[3,45] \cdot x[9,99] + 5 \cdot \text{mean}(\mathbf{x}) + \ldots$$

This approach of learning a function from <u>labeled</u> data is called **supervised learning.**

23

National Institute for Standards and Technology collected a huge amount of handwritten digit data from census workers and high school students in the early 90s:



This is called the NIST dataset, and was used to create the famous
MNIST handwritten digit dataset.

24

Since the 1990s machine learning have overtaken expert systems as the dominant approach to artificial intelligence.

- Current methods achieve .17% error rate for OCR on benchmark datasets (MNIST).[2]
- Very successful on other problems as well. The big break through for supervised learning in the 2010s was image classification.

_____

[2]Not because of overfitting! See: *Cold Case: The Lost MNIST Digits* by Chhavi Yadav + Léon Bottou.

Once we have the basic supervised machine learning setup, many very difficult questions remain:

- How do we parameterize a class of functions $f$ to search?
- How do we efficiently find a good function in the class?
- How do we ensure that an $f(x)$ which works well on our training data will generalize to perform well on future data?
- How do we deal with imperfect data (noise, outliers, incorrect training labels)?

Recall that in the supervised learning setup every input $x_i$ in our training dataset comes with a desired output $y_i$ (typically generated by a human, or some other process).

Types of supervised earning:

- **Classification** – predict a discrete class label.
- **Regression** – predict a continuous value.
  - Dependent variable, response variable, target variable, lots of different names for $y_i$.

Another example of supervised classification: **Face Detection**.



Each input data example $x_i$ is an image. Each output $y_i$ is 1 if the image contains a face, 0 otherwise.

- Harder than digit recognition, but we now have essentially perfect methods (used in nearly all digital cameras, phones, etc.)

Other examples of supervised classification:

- Object detection (Input: image, Output: dog or cat)
- Spam detection (Input: email text, Output: spam or not)
- Medical diagnosis (Input: patient data, Output: disease condition or not)
- Credit decision making (Input: financial data, Output: offer loan or not)

Example of supervised regression: Stock Price Prediction.



Each input x is a vector of metrics about a company (sales volume, PE ratio, earning reports, historical price data).

Each output $y_i$ is the **price of the stock** 3 months in the future.

Other examples of supervised regression:

- Home price prediction (Inputs: square footage, zip code, number of bathrooms, Output: Price)
- Car price prediction (Inputs: make, model, year, miles driven, Output: Price)
- Weather prediction (Inputs: weather data at nearby stations, Output: tomorrows temperature )
- Robotics/Control (Inputs: information about environment and current position at time $t$, Output: estimate of position at time $t + 1$)

Later in the class we will talk about other frameworks:

- **Unsupervised learning** (no labels or response variable)
  - Important for representation learning and generative ML.
- **Semi-supervised learning, self-supervised learning**.

Focus less in this class on:

- **Reinforcement learning**
  - Game playing
- **Active-learning.**
  - The learning algorithms can request labels.

Types of supervised learning:

- **Classification** – predict a <u>discrete</u> class label.
- **Regression** – predict a <u>continuous</u> value.
  - Dependent variable, response variable, target variable, lots of different names for $y_i$.

Motivating example: Predict the highway miles per gallon (MPG) of a car given quantitative information about its engine. Demo in `demo_auto_mpg.ipynb`.

What factors might matter?

Data set available from the UCI Machine Learning Repository: https://archive.ics.uci.edu/.

Datasets from UCI (and many other places) comes as tab, space, or comma delimited files.

Check dataset description to know what each column means.



'mpg', 'cylinders','displacement', 'horsepower', 'weight',
'acceleration', 'model year', 'origin', 'car name'

- Use `pandas` for reading data from delimited files. Stores data in a type of table called a "data frame" but this is just a wrapper around a `numpy` array.
- Use `matplotlib` for initial exploration.



Displacement　　　　Horsepower　　　　Acceleration

# SIMPLE LINEAR REGRESSION

Linear regression from a Machine Learning (not a Statistics) perspective. Our first supervised machine learning model.



Only focus on one predictive variable at a time (e.g. horsepower). This is why it's called simple linear regression.

Dataset:

- $x_1, \ldots, x_n \in \mathbb{R}$ (horsepowers of $n$ cars – this is the predictor/independent variable)
- $y_1, \ldots, y_n \in \mathbb{R}$ (MPG – this is the response/dependent variable)

- **Model** $f_{\boldsymbol{\theta}}(x)$: Class of equations or programs which map input $x$ to predicted output. We want $f_{\boldsymbol{\theta}}(x_i) \approx y_i$ for training inputs.

- **Model Parameters** $\boldsymbol{\theta}$: Vector of numbers. These are numerical knobs which parameterize our class of models.

- **Loss Function** $L(\boldsymbol{\theta})$: Measure of how well a model fits our data. Often some function of $|f_{\boldsymbol{\theta}}(x_1) - y_1|, \ldots, |f_{\boldsymbol{\theta}}(x_n) - y_n|$

**Common Goal:** Choose parameters $\boldsymbol{\theta}^*$ which minimize the Loss Function:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

$$a x^2 + b x + c$$
$$\theta = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

Choosing $\boldsymbol{\theta}^*$ based on minimizing the empirical error on our training data is called Empirical Risk Minimization. It is by far the most common approach to solving supervised learning problems.

41

### General Supervised Learning

- Model: $f_\theta(\boxed{x})$

$$\underline{\beta_1 x} + \underline{\beta_0} = f_{(\beta_0, \beta_1)}(x)$$

- Model Parameters: $\theta$

$$\text{intercept} \swarrow \beta_0, \quad \beta_1 \longrightarrow \text{slope}$$

- Loss Function: $L(\theta)$

$$\sum_{i=1}^{u} \left( y_i - f_\theta(x_i) \right)^2$$

### Linear Regression

- Model:

- Model Parameters:

- Loss Function:

42

$$\sum_{i=1}^{y} | f_\theta(x_i) - y_i |$$

What is a natural **loss function** for linear regression?



$$\max_i | f_\theta(x_i) - y_i |$$

$f_{\beta 0,\ \beta 1}(x)$

Typical choices are a function of $y_1 - f_{\beta_0,\beta_1}(x_1), \ldots, y_n - f_{\beta_0,\beta_1}(x_n)$



$y_i - f_{\beta_0,\beta_1}(x)$

$f_{\beta_0,\beta_1}(x)$

- $\ell_2$/Squared Loss: $L(\beta_0, \beta_1) = \sum_{i=1}^{n} \left(y_i - f_{\beta_0,\beta_1}(x_i)\right)^2$.
- $\ell_1$/Least absolute deviations: $L(\beta_0, \beta_1) = \sum_{i=1}^{n} |y_i - f_{\beta_0,\beta_1}(x_i)|$.
- $\ell_\infty$ Loss $L(\beta_0, \beta_1) = \max_{i \in 1,\ldots,n} |y_i - f_{\beta_0,\beta_1}(x_i)|$.

44

We're going to start with the Squared Loss/Sum-of-Squares Loss.
Also called "Residual Sum-of-Squares (RSS)"



$f_{\beta_0, \beta_1}(x)$

- Relatively <u>robust</u> to outliers.
- Simple to define, leads to simple algorithms for finding $\beta_0, \beta_1$
- Theoretically justified from <u>classical statistics</u> related to assumptions about Gaussian noise. Will discuss later in the course.

45

## General Supervised Learning

- Model: $f_{\boldsymbol{\theta}}(x)$

- Model Parameters: $\boldsymbol{\theta}$

- Loss Function: $L(\boldsymbol{\theta})$

## Linear Regression

- Model:
  $f_{\beta_0, \beta_1}(x) = \underbrace{\beta_0 + \beta_1 \cdot x}$

- Model Parameters: $\underbrace{\beta_0, \beta_1}$

- Loss Function: $L(\beta_0, \beta_1) = \sum_{i=1}^{n} (\underbrace{y_i} - f_{\beta_0, \beta_1}(\underbrace{x_i}))^2$

**Goal:** Choose $\beta_0, \beta_1$ to minimize
$L(\underbrace{\beta_0, \beta_1}) = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$.

This is the entire job of any Supervised Learning Algorithm.

Univariate function:



$$x^3 + 3 \cdot x^2 - 5 \cdot x + 1$$

· Find all places where <u>derivative</u> $f'(x) = 0$ and check which
  has the smallest value.

Multivariate function: $L(\beta_0, \beta_1)$

- Find values of $\beta_0, \beta_1$ where all partial derivatives equal 0.
- $\frac{\partial L}{\partial \beta_0} = 0$ and $\frac{\partial L}{\partial \beta_1} = 0$.

Multivariate function: $\underline{L(\beta_0, \beta_1)} = \sum_{i=1}^{n}\underline{(y_i - \beta_0 - \beta_1 x_i)^2}$

· Find values of $\beta_0, \beta_1$ where <u>all</u> partial derivatives equal 0.

· $\frac{\partial L}{\partial \beta_0} = 0$ and $\frac{\partial L}{\partial \beta_1} = 0$.

Some definitions:

· Let $\bar{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$.                $\underline{\bar{y}}$ is the <u>mean</u> of $y$.
· Let $\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i$.                $\bar{x}$ is the <u>mean</u> of $x$.
· Let $\sigma_y^2 = \frac{1}{n}\sum_{i=1}^{n}(y_i - \bar{y})^2$.      $\sigma_y^2$ is the <u>variance</u> of $y$.
· Let $\sigma_x^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2$.      $\sigma_x^2$ is the <u>variance</u> of $x$.
· Let $\sigma_{xy} = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})$.    $\sigma_{xy}$ is the <u>covariance</u>.

Claim: $L(\beta_0, \beta_1)$ is minimized when:

· $\beta_1 = \sigma_{xy}/\sigma_x^2$
· $\beta_0 = \bar{y} - \beta_1 \bar{x}$

Loss function: $L(\beta_0, \beta_1) = \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)^2$

$\dfrac{\partial L}{\partial \beta_0} = 0 \qquad \dfrac{\partial L}{\partial \beta_1} = 0$

$\dfrac{\partial}{\partial \beta_0} \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^{n} -2(y_i - \beta_0 - \beta_1 x_i) = 0$

$= \sum_{i=1}^{n} \dfrac{\partial}{\partial \beta_0}(y_i - \beta_0 - \beta_1 x_i)^2 = \sum_{i=1}^{n} 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-1)$

$\sum_{i=1}^{n} -y_i + \sum_{i=1}^{n} \beta_0 + \sum_{i=1}^{n} \beta_1 x_i = 0$

$n \cdot \beta_0 + \beta_1 \sum_{i=1}^{n} x_i = \sum_{i=1}^{n} y_i$

$\beta_0 + \beta_1 \dfrac{1}{n} \sum_{i=1}^{n} x_i = \dfrac{1}{n} \sum_{i=1}^{n} y_i$

$\beta_0 + \beta_1 \bar{x} = \bar{y}$

$\beta_0 = \bar{y} - \beta_1 \bar{x}$

50

$$L(\beta_0, \beta_1)$$

Loss function after substitution:
$$\tilde{L}(\beta_1) = \sum_{i=1}^{n}(y_i - \bar{y} + \beta_1\bar{x} - \beta_1 x_i)^2$$

$$\frac{\partial \tilde{L}}{\partial \beta_1} = 0 \qquad \frac{\partial}{\partial \beta_1} \sum_{i=1}^{n}\left(y_i - \bar{y} + \beta_1\bar{x} - \beta_1 x_i\right)^2 = 0$$

$$= \sum_{i=1}^{n} 2\left(y_i - \bar{y} + \beta_1\bar{x} - \beta_1 x_i\right)\left(\bar{x} - x_i\right) = 0$$

$$\sum_{i=1}^{n}\left(y_i - \bar{y}\right)\left(\bar{x} - x_i\right) + \beta_1 \sum_{i=1}^{n}\underbrace{\left[\bar{x} - x_i\right]\left(\bar{x} - x_i\right)}_{} = 0$$

$$= -6_{xy} \cdot n \qquad \beta_1 \cdot 6_x^2 \cdot n = 0$$

$$\beta_1 = \frac{6_{xy} \cdot n}{6_x^2 \cdot n} = 6_{xy}/6_x^2$$

**Takeaways:**

- Minimizing functions exactly is sometimes easy with calculus, but not always! We will learn much more general tools (like gradient descent).
- Simple closed form formula for optimal parameters $\beta_0^*$ and $\beta_1^*$ for squared-loss!



$$B_0^* + B_1^* x$$

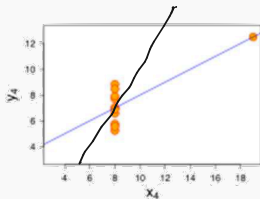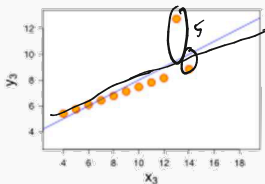Let $L(\beta_0, \beta_1) = \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)^2$.
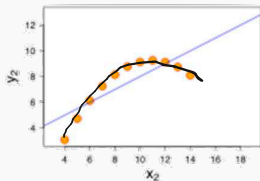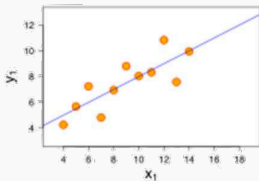
$$R^2 = 1 - \frac{L(\beta_0, \beta_1)}{n\sigma_y^2}$$

is exactly the $R^2$ value ("coefficient of determination") you may remember from statistics.

The smaller the loss, the closer $R^2$ is to 1, which means we have a better regression fit.
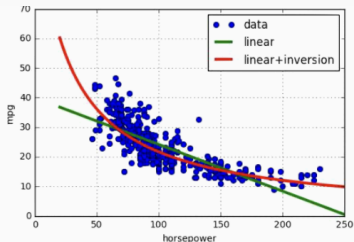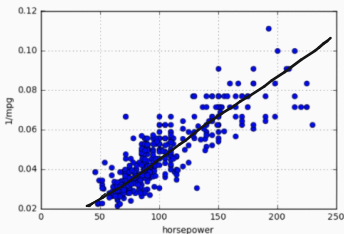
53

Many reasons you might get a poor regression fit:

Some of these are fixable!

- Remove outliers, use more robust loss function.
- Non-linear model transformation.

Fit the model $\frac{1}{\text{mpg}} \approx \beta_0 + \beta_1 \cdot \text{horsepower}$.



Much better fit, same exact learning algorithm!

# MULTIPLE LINEAR REGRESSION

Predict target *y* using multiple features, simultaneously.

**Motivating example:** Predict diabetes progression in patients after 1 year based on health metrics. (Measured via numerical score.)

**Features:** Age, sex, average blood pressure, six blood serum measurements (e.g. cholesterol, lipid levels, iron, etc.)

Demo in `demo_diabetes.ipynb`.

## Introducing **Scikit Learn**.

Pros:

- One of the most popular "traditional" ML libraries.

- Many built in models for regression, classification, dimensionality reduction, etc.

- Easy to use, works with 'numpy', 'scipy', other libraries we use.

- Great for rapid prototyping, testing models.

Cons:

- Everything is very "black-box": difficult to debug, understand why models aren't working, speed up code, etc.

Modules used:

- `datasets` module contains a number of pre-loaded datasets. Saves time over downloading and importing with `pandas`.
- `linear_model` can be used to solve Multiple Linear Regression. A bit overkill for this simple model, but gives you an idea of `sklearn`'s general structure.

Target variable:

- Scalars $y_1, \ldots, y_n$ for $n$ data examples (a.k.a. samples).

Predictor variables:

- $d$ dimensional vectors $\mathbf{x}_1, \ldots, \mathbf{x}_n$ for $n$ data examples and $d$ features

Now it the time to review your linear algebra!

Notation:

- Let $\underline{\underline{X}}$ be an $n \times d$ matrix. Written $X \in \mathbb{R}^{n \times d}$.
- $\underline{x_i}$ is the $i^{\text{th}}$ row of the matrix.
- $\underline{x}^{(j)}$ is the $j^{\text{th}}$ column.

$X^T \in \mathbb{R}^{d \times n}$

- $\underline{x_{ij}}$ is the $i, j$ entry.
- For a vector $\underline{y}$, $y_i$ is the $i^{\text{th}}$ entry.
- $X^T$ is the matrix transpose.
- $y^T$ is a vector transpose.

$$1 \times u \quad u \times 1 \quad \rightarrow \quad 1 \times 1$$

Things to remember:

- Matrix multiplication. If I multiply $X \in \mathbb{R}^{n \times d}$ by $Y \in \mathbb{R}^{d \times k}$ I get $XY = Z \in \mathbb{R}^{n \times k}$.
- Inner product/dot product. $\langle y, z \rangle = \sum_{i=1}^{n} y_i z_i$.
- $\langle y, z \rangle = y^T z = z^T y$.
- Euclidean norm: $\|y\|_2 = \sqrt{y^T y}$.
- $(XY)^T = Y^T X^T$.

$$\|y\|_2 = \left( \sum_{i=1}^{u} y_i^2 \right)^{1/2} = \left( y^T y \right)^{1/2} = \left( \langle y, y \rangle \right)^{1/u}$$

$$\begin{pmatrix} 1 & 6 & 6 & 5 \\ 0 & 1 & 6 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Things to remember:

- Identity matrix is denoted as $I$.
- "Most" square matrices have an inverse: i.e. if $Z \in \mathbb{R}^{n \times n}$, there is a matrix $Z^{-1}$ such that $Z^{-1}Z = ZZ^{-1} = I$.
- Let $D = \text{diag}(d)$ be a diagonal matrix containing the entries in $d$.
- $XD$ scales the columns of $X$. $DX$ scales the rows.

You also need to be comfortable working with matrices in
numpy . Go through the demo_numpy_matrices.ipynb
slowly.
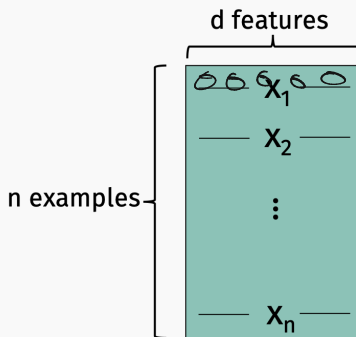
Target variable:

- Scalars $y_1, \ldots, y_n$ for $n$ data examples (a.k.a. samples).

Predictor variables:

- $d$ dimensional vectors $\mathbf{x}_1, \ldots, \mathbf{x}_n$ for $n$ data examples and $d$ features



$$(n \times d)(d \times 1) = n \times 1$$

d features

n examples

age | sex | BMI | ...

**X**

65

Data matrix indexing:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \ldots & x_{1d} \\ x_{21} & x_{22} & \ldots & x_{2d} \\ x_{31} & x_{32} & \ldots & x_{3d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \ldots & x_{nd} \end{bmatrix}$$

Multiple Linear Regression Model:

Predict $\quad y_i \approx \beta_1 x_{i1} + \beta_2 x_{i2} + \ldots + \beta_d x_{id}$

The rate at which diabetes progresses depends on many factors, with each factor having a different magnitude effect.

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ x_{31} & x_{32} & \dots & x_{3d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1d} \\ 1 & x_{22} & \dots & x_{2d} \\ 1 & x_{32} & \dots & x_{3d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

Multiple Linear Regression Model:

$\beta_0$

Predict $\qquad y_i \approx \beta_1 + \beta_2 x_{i2} + \dots + \beta_d x_{id}$

In this case, $\beta_1$ serves as the "intercept" parameter.

Use as much linear algebra notation as possible!

- Model Parameters: $\beta_1, \ldots \beta_d$

- Model: $y_i = \beta_1 x_{i_1} + \beta_2 x_{i_2} + \ldots + \beta_d x_{id}$

$\widehat{y}$   $\underline{X}$   $f_\beta(X) = X \beta$

- Loss Function:

$$L(\beta) = \| y - X\beta \|_2^2$$

## MULTIPLE LINEAR REGRESSION

### Linear <u>Least-Squares</u> Regression.

- Model Parameters:

$$\boldsymbol{\beta} = [\beta_1, \beta_2, \ldots, \beta_d]$$

- Model:

$$f_{\boldsymbol{\beta}}(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle$$

- Loss Function:

$$L(\boldsymbol{\beta}) = \sum_{i=1}^{n} |y_i - \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle|^2$$
$$= \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

**Machine learning goal:** minimize the loss function
$\underline{L(\boldsymbol{\beta})} : \mathbb{R}^d \to \mathbb{R}.$

$$\frac{\partial L}{\partial \beta_1} = 0 \quad \cdots \quad \frac{\partial L}{\partial \beta_d} = 0$$

Find optimum by determining for which $\boldsymbol{\beta} = [\beta_1, \dots, \beta_d]$ all
partial derivatives are 0. I.e. when do we have:

$$\begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$L(\boldsymbol{\beta})$

$L(\tilde{\boldsymbol{\beta}})$

$\nabla \underline{L(\beta)}$

For any function $L(\beta) : \underline{\mathbb{R}^d} \to \underline{\mathbb{R}}$, the gradient $\nabla L(\beta)$ is a function from $\underline{\underline{\mathbb{R}^d}} \to \underline{\underline{\mathbb{R}^d}}$ defined:

$$\nabla L(\beta) = \begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix}$$

The gradient of the loss function is a central tool in machine learning. We will use it again and again.

Loss function:

$$L(\beta) = \|y - X\beta\|_2^2$$

$X \in \mathbb{R}^{n \times d}$

$X^T \in \mathbb{R}^{d \times n}$

Gradient:

$(d \times n)(n \times 1) \rightarrow (n \times d)(d \times 1) = n \times 1$

$$-2 \cdot X^T (\underbrace{y - X\beta}_{n \times 1}) \qquad (d \times 1)$$

Find optimum by determining for which $\beta = [\beta_1, \ldots, \beta_d]$ the gradient is 0. I.e. when do we have:

$$\nabla L(\beta) = \begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

73

**Goal:** minimize the loss function $L(\beta) = \|y - X\beta\|_2^2.$

$$= -2X^T y + 2 X^T X \beta$$

$$\nabla L(\beta) = -2 \cdot X^T(y - X\beta) = 0$$

$$= 2X^TX\beta - 2X^Ty = 0$$

$(d \times u)(u \times d)$
$= (d \times d)$

**Solve for optimal $\beta^*$:**

$$2 X^T X \beta = 2 X^T y$$

$O(d^3)$

$$(X^TX)^{-1}(X^TX)\beta^* = (X^TX)^{-1}X^Ty$$

$$\boxed{\beta^* = (X^TX)^{-1}X^Ty}$$

$$O(udd) = O(ud^2)$$

## MULTIPLE LINEAR REGRESSION SOLUTION

Need to compute $\beta^* = \arg\min_{\beta} \|y - X\beta\|_2^2 = (X^TX)^{-1}X^Ty$.

- Main cost is computing $(X^TX)^{-1}$ which takes $O(nd^2)$ time.
- Can solve slightly faster using the method
  numpy.linalg.lstsq, which is running an algorithm
  based on QR decomposition.
- For larger problems, can solve much faster using an
  *iterative methods* like scipy.sparse.linalg.lsqr.

Will learn more about iterative methods when we study
Gradient Descent.

Function: $\mathbb{R}^d \to \mathbb{R}$

$$\underline{f(\mathbf{z})} = \underline{\mathbf{a}^T \mathbf{z}} \text{ for some fixed vector } \underline{\mathbf{a}} \in \mathbb{R}^d$$

$$q_1 z_1$$

Gradient:

$$f(z) = \langle a, z \rangle = \sum_{i=1}^{d} q_i z_i$$

$$\nabla f(z) = \begin{bmatrix} \partial f/z_1 \\ \vdots \\ \partial f z_d \end{bmatrix} = \begin{bmatrix} q_1 \\ a_2 \\ \vdots \\ a_d \end{bmatrix} = \boxed{a}$$

Function:

$$\frac{d}{z_1}(z_1)^2 = 2z_1$$

$$\underline{f(\mathbf{z}) = \underline{\|\mathbf{z}\|_2^2}}$$

$$\sum_{i=1}^{d}(z_i)^2$$

Gradient:

$$\nabla f(z) = \begin{bmatrix} \partial f/z_1 \\ \partial f/z_d \end{bmatrix} = \begin{matrix} 2z_1 \\ \vdots \\ 2z_d \end{matrix} = 2z$$

76

Loss function:

$$L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

Loss function:

$$L(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

Model: $f_{\boldsymbol{\beta}}(x) = \sum_{i=1}^{d} \beta_1 x_i$

Example from book: What is the sign of $\beta_1$ when we run a simple linear regression using the following predictors for number of sales in a particular market as a function of:

- Amount of TV advertising in that market:
- Amount of print advertising in that market:

What is the sign of the corresponding $\beta$'s when we run a
<u>multiple</u> linear regression using the following predictors
together:

- <u>Amount of TV advertising in that market</u>: Positive
- <u>Amount of print advertising in that market</u>: Negative, close
  to zero

Can you explain this? Try to think of your own example of a
regression problem where this phenomenon might show up.

## DEALING WITH CATEGORICAL VARIABLES

The sex variable in the diabetes problem was binary. We encoded it as 2 numbers – e.g. (0,1), (-1,1), (1,2).

Suppose we go back to the MPG prediction problem. What if we had a categorical predictor variable for car make with more than 2 options: e.g. Ford, BMW, Honda. How would you encode as a numerical column?

$$
\begin{bmatrix} \text{ford} \\ \text{ford} \\ \text{honda} \\ \text{bmw} \\ \text{honda} \\ \text{ford} \end{bmatrix} \rightarrow \begin{bmatrix} \\ \\ \\ \\ \\ \end{bmatrix}
$$

Better approach: <u>One Hot Encoding.</u>

$$
\begin{bmatrix}
\texttt{ford} \\
\texttt{ford} \\
\texttt{honda} \\
\texttt{bmw} \\
\texttt{honda} \\
\texttt{ford}
\end{bmatrix}
\rightarrow
\begin{bmatrix}
1 & 0 & 0 \\
1 & 0 & 0 \\
0 & 1 & 0 \\
0 & 0 & 1 \\
0 & 1 & 0 \\
1 & 0 & 0
\end{bmatrix}
$$

- Create a separate feature for every category, which is 1 when the variable is in that category, zero otherwise.
- Not too hard to do by hand, but you can also use library functions like `sklearn.preprocessing.OneHotEncoder`.

Avoids adding inadvertent linear relationships.