


- This session is recorded.
- Mute your mic.
- You can leave messages in the chat panel to ask questions.
- When I check your progress, please click  if you are done.
yes
- We will start at 6:35.

CS5242 Neural Networks and Deep Learning

Lecture 01: Introduction and Linear Regression

Wei WANG

TA: FU Di, LIANG Yuxuan, FuYujian

cs5242@comp.nus.edu.sg

change log:
V1: slide 41, the superscript should be n



Agenda

- Introduction of this module
 - Logistics
 - History
- Linear regression
 - Univariate
 - Multivariate

Instructor

- WANG Wei
 - COM2-04-09
 - wangwei@comp.nus.edu.sg
- Research
 - Machine learning system optimization
 - Multi-modal data analysis/applications

Teaching Assistants



LIANG Yuxuan

e0427783@u.nus.edu

Quiz



FU Di

e0409760@u.nus.edu

Assignment



FU Yujian

e0427770@u.nus.edu

Final project

Your background?

LumiNUS Poll

Pre-requisite

Course

- Machine Learning (CS3244)
 - Or <https://www.coursera.org/learn/machine-learning>
- Linear Algebra (MA1101R)
 - Or <https://ocw.mit.edu/courses/mathematics/18-06-linear-algebra-spring-2010/>
- Calculus (MA1521)
- Probability (ST2334)
- [Brief summary](#)

Coding

- Python (ONLY, version 3.x)
 - [Numpy](#)
 - Keras, TensorFlow, PyTorch, MxNet, Caffe, SINGA
- Jupyter notebook or Google colab

Grading policy

Weightage:

- Assignment 30 points
 - Two assignments
 - 15% off per day late (17:01 is the start of one day); 0 score if you submit it 7 days after the deadline
- Online quiz 30 points
 - 3 quizzes in total
 - You need a camera (from your laptop, mobile phone or external camera)
- Project 40 points
 - Kaggle competition
 - Group size ≤ 4

Collaboration

- Every assignment is an individual task
- The project is a group-based task
- Avoid academic offence (cheating, plagiarism including copying code from the internet, e.g., github)

Contacts

- LumiNUS forum
 - For all issues
- Email: cs5242@comp.nus.edu.sg
 - For private/personal issues
- Consultation (make appointments on LumiNUS)
 - Wei WANG, wangwei@comp.nus.edu.sg
 - Wed, 11:00-12:00

GPU machines

- SoC GPU machines
 - <https://dochub.comp.nus.edu.sg/cf/guides/compute-cluster/hardware>
- Google Cloud Platform
 - Some free credit for new register
 - GPU is expensive
- Amazon EC2 (g2.8xlarge)
 - Some free credit for students
 - GPU is expensive
- NSCC (National Super-Computing Center)
 - Free for NUS students
 - The libraries are sometimes outdated
 - Jobs will be submitted into a queue for executing
- **Important notes if you use cloud platforms:**
 - STOP/TERMINATE the instance immediately after your program terminates
 - Check the usage status frequently to avoid high outstanding bills
 - Amazon/Google may charge you for additional storage volume

Syllabus & Schedule

- Learn various neural network models
 - From shallow to deep neural networks
 - Including the model structure, training algorithms, tricks and applications
 - Covering the principles, coding practices and a bit theory
- Check LumiNUS for the latest plan

Intended learning outcomes (my expectation)

1

Explain the principles of the operations of different layers and training algorithms

2

Compare different neural network architectures

3

Implement popular neural networks

4

Solve problems using neural networks and deep learning techniques

Your expectation?

*I want to get an easy boost to my
GPA.*

This course may not be easy

*I want to earn the big bucks as a
data scientist / ML engineer
at Google / Amazon / Facebook.*



This course provides the basics; but is not enough.

check out:

CS5260 – NN & Deep Learning II

CS4243 – Computer vision & Pattern Recognition

CS4248 – Natural Language Processing

indeed Find jobs Company reviews Find salaries

What
Job title, keywords, or company

Where
Country, Town or MRT Station

data scientist Singapore **Find jobs** Advanced Job Search

Salary Estimate Job Type Location Company

data scientist jobs in Singapore

Sort by: **relevance** - date Page 1 of 613 jobs ?

Data Scientist
Kelly Services Singapore 3.9 ★
Singapore

- A Bachelor's degree in Data Scientist or Computer Science (or equivalent experience).
- A Bachelor in Data Scientist/Computer Course.

Sponsored · 30+ days ago · [Save job](#)

Senior Data Scientist, Credit Scoring
Singapore
\$180,000 a year

- This is a great opportunity for an experienced data scientist to develop innovative products and services using large data sets and advanced algorithms.

Sponsored · 24 days ago · [Save job](#)

Computer Vision Researcher
Trakomatic pte ltd
Singapore
\$6,500 - \$7,500 a month

- Trakomatic has deployed hundreds of cameras in strategic locations across Asia-Pacific region for data collection and real-time analytics.

Sponsored · 30+ days ago · [Save job](#)

Get new jobs for this search by email

My email:

[Activate](#)

By creating a job alert or receiving recommended jobs, you agree to our [Terms](#). You can change your consent settings at any time by unsubscribing or as detailed in our terms.

My recent searches

- [Computer Vision - Singapore](#)
- [copmputer vision - Singapore](#)
- [computer vision developer - Singapore](#)

[» clear searches](#)

Page 1 of 613

\$180,000 a year

*My PhD advisor is making me take
this course so that I can use deep
learning in our research.*

Neural Networks and Deep Learning:

- Andrew Ng
- <https://www.coursera.org/learn/neural-networks-deep-learning>

CSC 321: Intro to Neural Networks and Machine Learning

- Roger Grosse
- http://www.cs.toronto.edu/~rgrosse/courses/csc321_2017/

Neural Networks for Machine Learning

- Geoffrey Hinton
- <https://www.coursera.org/learn/neural-networks>

CS231n: Convolutional Neural Networks for Visual Recognition

- Fei-Fei Li, Justin Johnson, Serena Yeung
- <http://cs231n.stanford.edu/>

CS224d: Deep Learning for Natural Language Processing

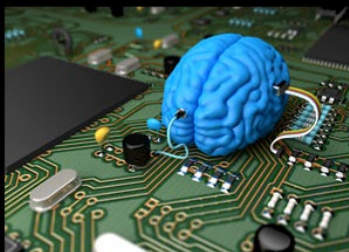
- Richard Socher
- <http://cs224d.stanford.edu/>

AI is cool and I want to build the next Skynet.

Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



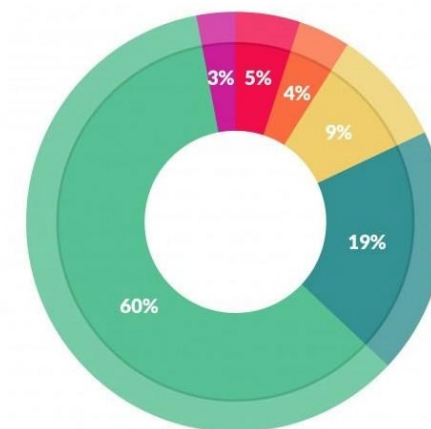
What I think I do

```
In [1]:  
import keras  
Using TensorFlow backend.
```

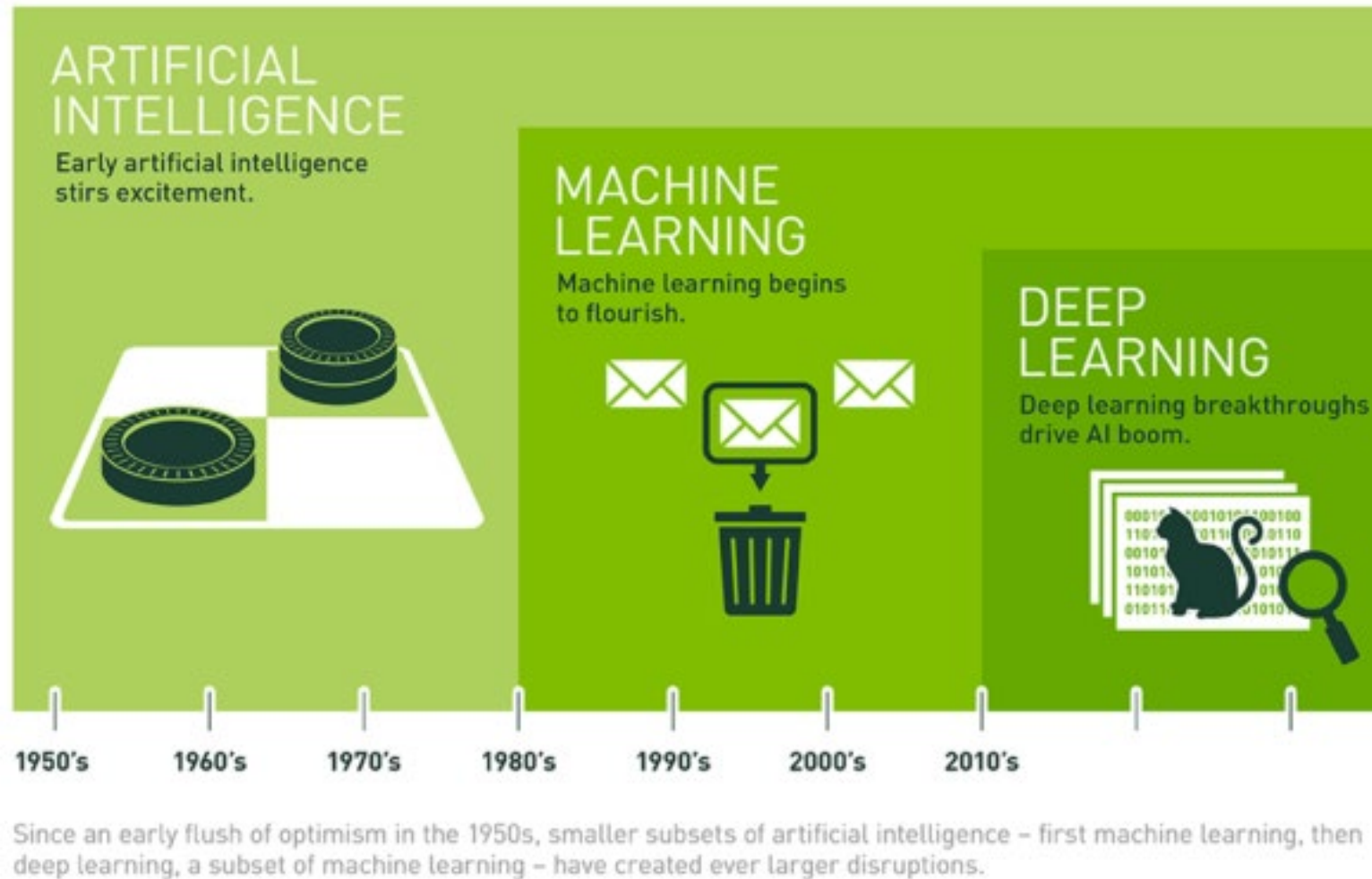
What I actually do

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



Terminology Overload: AI vs. ML vs. DL



AI

- 1950's
- “Human intelligence exhibited by machines”
 - Expert systems; Rules
 - Machine learning algorithms
- Narrow AI: image recognition, machine translation

<https://www.youtube.com/watch?v=nASDYRkbQIY>

Machine Learning

- 1980's
- “An approach to achieve AI through systems that can learn from experience to find patterns in that data”
 - Can we code a program with rules (like bubble sorting) to do
 - Image recognition



[This Photo](#) by Unknown Author is licensed under [CC BY-NC-ND](#)



- Speech recognition

Neural networks and Deep Learning

- 1950's
- A class of machine learning algorithm that use a cascade of multiple layers of nonlinear processing units for feature extraction and transformation---Wikipedia

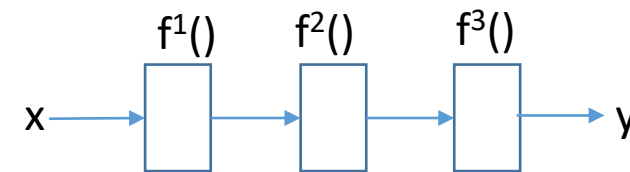
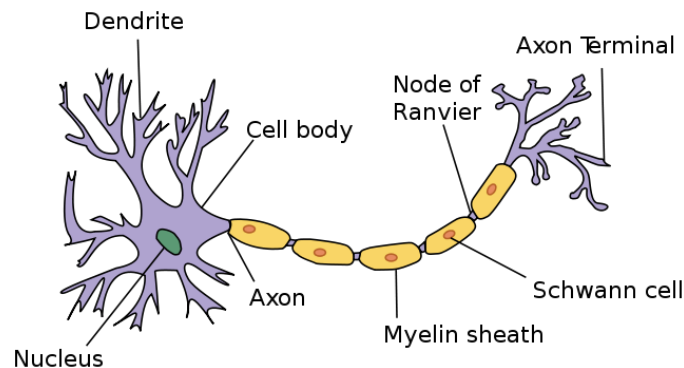
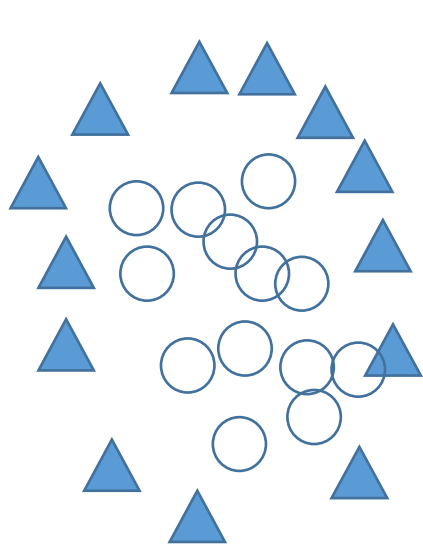


Image source: <http://pediaa.com/difference-between-nerve-and-neuron/>

Neural nets/perceptrons are *loosely* inspired by biology. But they certainly are *not* a model of how the brain works, or even how neurons work.

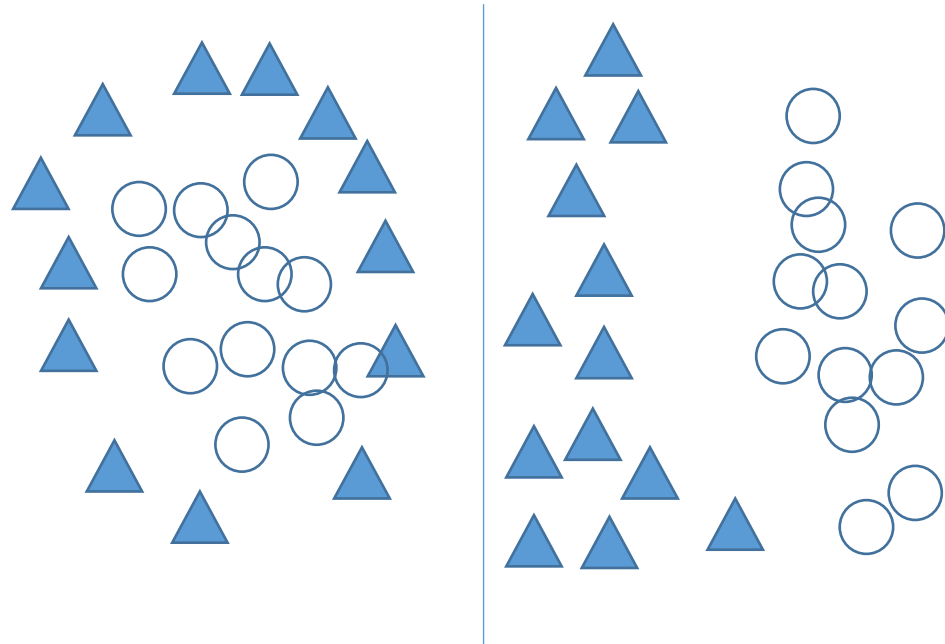
Neural networks and Deep Learning

- Feature learning (transformation)

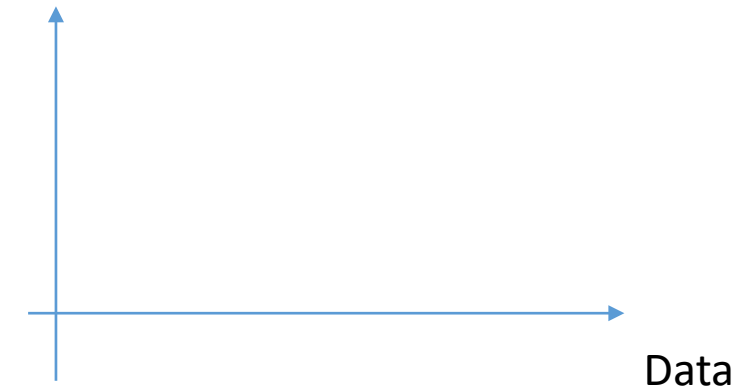


Neural networks and Deep Learning

- Feature learning (transformation)



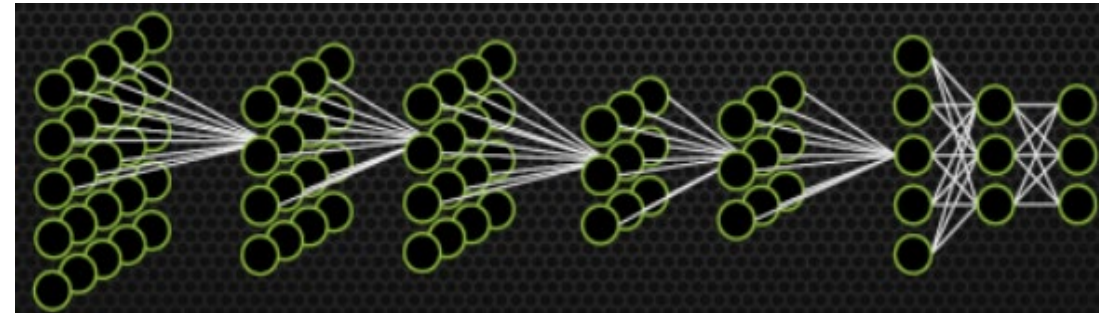
Performance



Deep learning is not a killer for everything

Neural networks (NN)

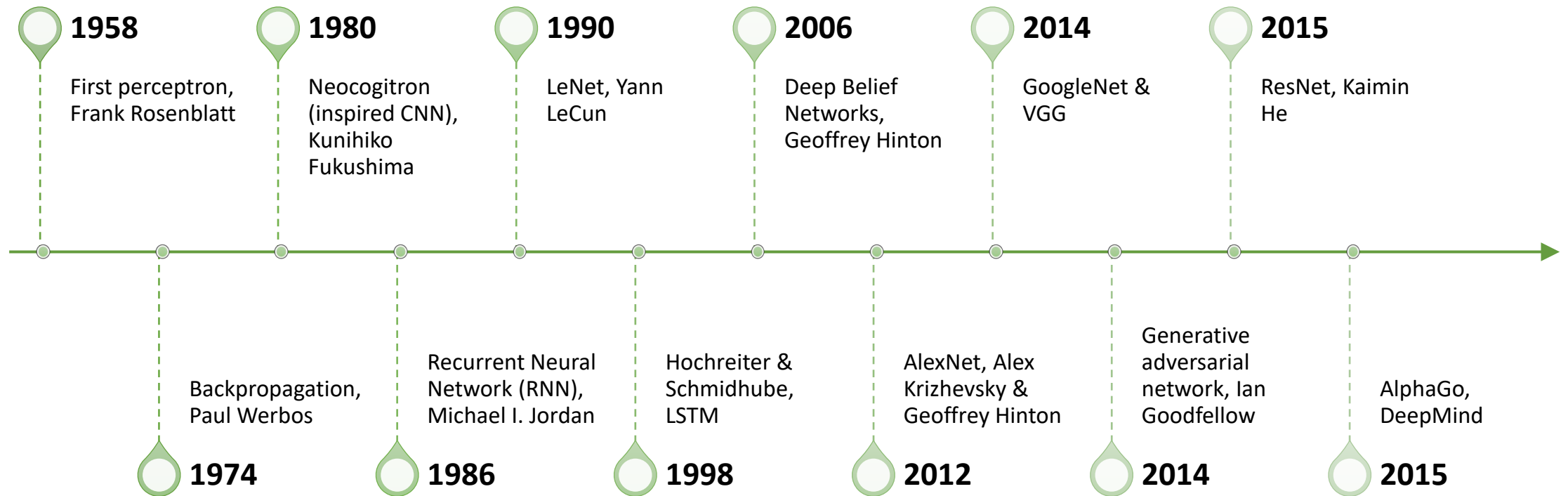
- Linear, polynomial, logistic, multinomial regression
- Perceptron and multi-layer perceptron (MLP)
- Convolutional neural network (CNN)
- Recurrent neural networks (RNN)
- Generative adversarial networks (GAN)
- (Restricted) Boltzmann machine
- Deep belief network
- Spike neural network
- Radial basis function neural network
- Hopfield networks



Source from [3]

History [2]

<https://www.youtube.com/watch?v=yaL5ZMvRRqE>



Refer to <http://people.idsia.ch/~juergen/who-invented-backpropagation.html> for more papers

Applications

Computer Vision

- [Image classification](#)
- [Object detection](#), [demo](#)
- [Scene text recognition](#)
- [Neural style transferring](#)
- [Image generation](#)

Natural language processing

- Question answering
- Machine translation

Speech

- Speech recognition, e.g. [Amazon Echo](#)

Univariate Linear Regression

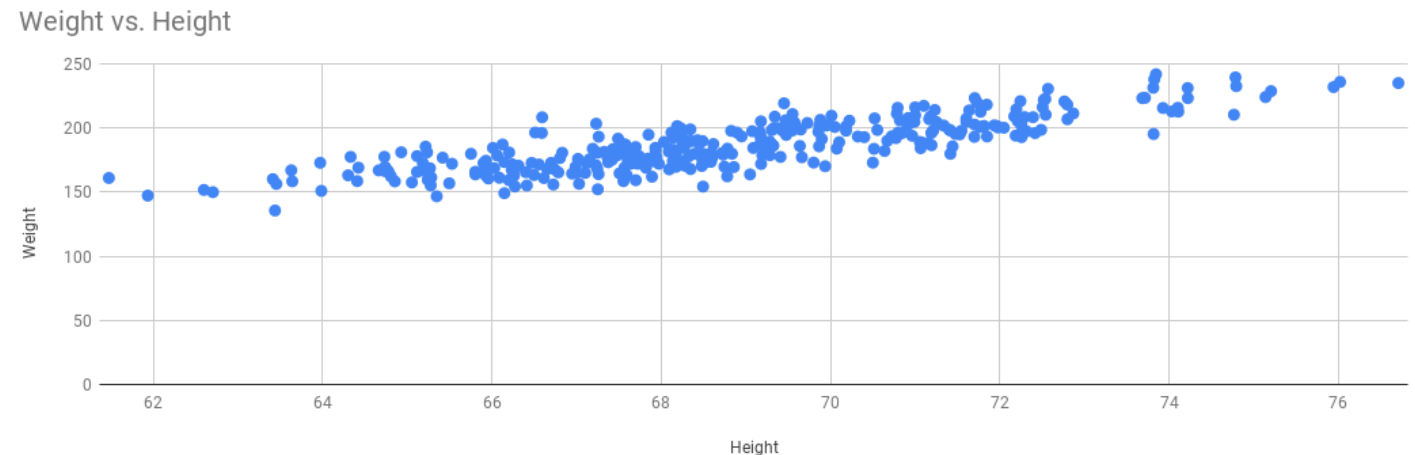
Weight vs Height

Problem definition:

- Predict the weight given the height of a person.

Data

- Input: Height (inches)
- Output: Weight (pounds)
- Split the data into
 - 80% for training
 - 20% for evaluation



Data source: <https://www.kaggle.com/mustafaali96/weight-height>

Modeling

- Notation
 - A person is called an example/instance
 - Height denoted as $x \in R$: input feature
 - Weight denoted as $y \in R$: the target or ground truth
- Map from input to output by **linear regression**
 - $\tilde{y} = xw + b, w \in R, b \in R$
 - w is the slope and b is the intercept
 - \tilde{y} is called the **prediction**

Training

- Learn w and b to fit the training data $S_{train} = \{(x^{(i)}, y^{(i)})\}, i = 1 \dots m$
 - That fit the data well
 - How to measure the quality of w and b ?
- Loss function
 - The **smaller** the loss value, the better the prediction (closer to the target)
 - $L(x, y|w, b) = |\tilde{y} - y|$
 - $L(x, y|w, b) = \frac{1}{2} ||\tilde{y} - y||^2 = \frac{1}{2} (\tilde{y} - y)^2$
 - Easier for optimization/training because it is differentiable
 - The coefficient $\frac{1}{2}$ is to make the gradient simple
- Define the training objective
 - $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)}|w, b)$

Optimization/training

Tune the model parameters over the training instances to minimize the average loss, i.e., the training objective

$$\min_{w,b} J(w, b)$$

$$\rightarrow \min_{w,b} \frac{1}{m} \sum_{i=1}^m L(x^{(i)}, y^{(i)} | w, b)$$

$$\rightarrow \min_{w,b} \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} + b - y^{(i)})^2$$

Optimization/training

1. Fix b and learn w

$$\min_w \frac{1}{2m} \sum_{i=1}^m (x^{(i)})^2 w^2 + \frac{1}{m} \sum_{i=1}^m x^{(i)} (b - y^{(i)}) w + \frac{1}{2m} \sum_{i=1}^m (b - y^{(i)})^2$$

$$\rightarrow \min_w c_1 w^2 + c_2 w + c_3$$

2. Fix w and learn b

Gradient for univariate simple functions

Gradient: recall high school calculus and how to take derivatives

- $\frac{\partial f(x)}{\partial x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

- $f(x) = 3, \frac{\partial f(x)}{\partial x} = 0$

- $f(x) = 3x, \frac{\partial f(x)}{\partial x} = 3$

- $f(x) = x^2, \frac{\partial f(x)}{\partial x} = 2x$

Gradient for univariate composite functions

- $f(x) = g(x) + h(x)$
 - $f(x) = 3x + x^2$

- $f'(x) = g'(x) + h'(x)$

We use $f'(x)$ and $\frac{\partial f(x)}{\partial x}$ interchangeably for the gradient of $f(x)$ with respect to x , where x and $f(x)$ are scalars

- $f(x) = g(x)h(x)$
 - $f(x) = xx^2$

- $f'(x) = g'(x)h(x) + g(x)h'(x)$

- $f(x) = \frac{g(x)}{h(x)}$
 - $f(x) = \frac{e^x}{x}$

- $f'(x) = \frac{g'(x)h(x) - g(x)h'(x)}{h(x)^2}$

- $f(x) = g(u), u = h(x)$
 - $f(x) = \ln(x + x^2)$

- $f'(x) = g'(u)h'(x)$

HW1: solve for the derivatives / gradients of these functions.

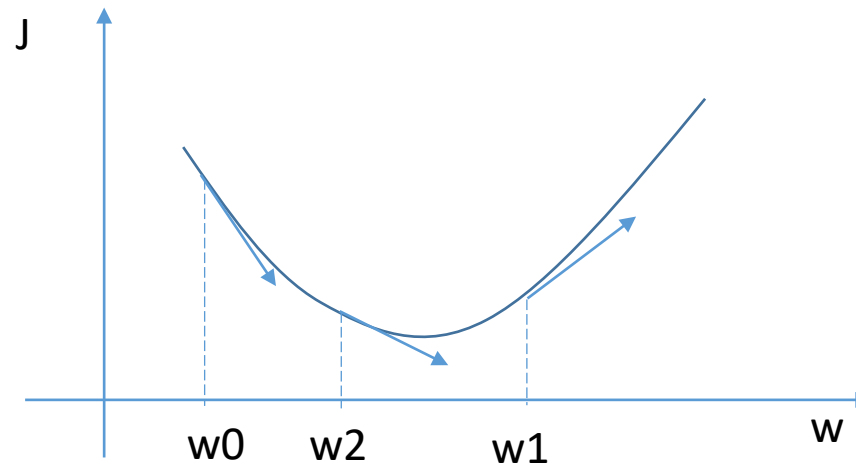
Training by gradient descent

$$J = c_1 w^2 + c_2 w + c_3$$

Find the w for which we have the lowest* J .

For example:

$$c_1 = 1, c_2 = -2, c_3 = 1$$



α is the **learning rate**, which controls the moving step length. It's important for convergence. If it is large, w oscillates around the optimal position. If it is small, it takes many iterations to reach the optimum.

Initialize w as w_0

Compute $\frac{\partial J}{\partial w} \big|_{w=w_0}$, negative;

Move w from w_0 to the right by

$$w_1 = w_0 - \alpha \frac{\partial J}{\partial w} \big|_{w=w_0}$$



Compute $\frac{\partial J}{\partial w} \big|_{w=w_1}$, positive;

Move w from w_1 to the left by

$$w_2 = w_1 - \alpha \frac{\partial J}{\partial w} \big|_{w=w_1}$$



Compute $\frac{\partial J}{\partial w} \big|_{w=w_2}$, negative

Move w from w_2 to the right by

$$w_3 = w_2 - \alpha \frac{\partial J}{\partial w} \big|_{w=w_2}$$

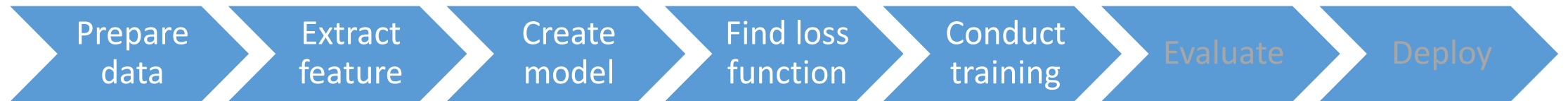
Training by gradient descent

- Gradient descent algorithm for optimization
- Set $w = 0.1$ or a random number
- Repeat
 - For each data sample, compute $\tilde{y} = xw + b$
 - Compute the average loss, $\sum_{\langle x, y \rangle \in S_{train}} L(x, y | w, b) / |S_{train}|$
 - Compute $\frac{\partial J}{\partial w}$
 - Update $w = w - \alpha \frac{\partial J}{\partial w}$

Update w, b repeatedly...

What if there are multiple parameters, e.g., multiple variables?

Machine learning pipeline



Multivariate Linear Regression

Multivariate linear regression

- Consider the problem of [house price prediction](#)
- Each instance in the training dataset
 - Denote the features using a column vector
 - $\mathbf{x} \in R^{n \times 1}$: x_i is the i-th feature
 - size, floor, location, age, lease, etc.
 - Target $y \in R$: price

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix}$$

Model

- Map from input to output

$$\tilde{y} = \mathbf{w}^T \mathbf{x} + b, \mathbf{w} \in R^{n \times 1}, b \in R, w_i \text{ is the } i\text{-th element of } \mathbf{w} \text{ (importance of } i\text{-th feature)}$$

$$= \sum_{i=1}^n w_i x_i + b \quad \text{The summation is over } n \text{ elements, where } n \text{ is the number of features per instance.}$$

$$= (w_1, w_2, \dots, w_n) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + b$$

$$= (w_1, w_2, \dots, w_n, b) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \\ 1 \end{pmatrix}$$

$$\rightarrow \tilde{y} = \bar{\mathbf{w}}^T \bar{\mathbf{x}}$$

For the rest of this module, we use \mathbf{w} and \mathbf{x} for $\bar{\mathbf{w}}$ and $\bar{\mathbf{x}}$ respectively unless there is a special definition of \mathbf{w} and \mathbf{x} .

Optimization

- Compute the gradient of the loss with respect to **(w.r.t) \mathbf{w}**
 - Consider a single instance

$$J(\mathbf{w}) = L(\mathbf{x}, y | \mathbf{w}) = \frac{1}{2} \|\tilde{y} - y\|^2 = \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^2$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = ?$$

Gradient of vector and matrix (denominator layout)

- Vectors
 - By default, is a column vector
 - Denoted as \mathbf{x} , \mathbf{y} , \mathbf{z}

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}.$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial \mathbf{x}} & \frac{\partial y_2}{\partial \mathbf{x}} & \cdots & \frac{\partial y_m}{\partial \mathbf{x}} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \frac{\partial y_2}{\partial x_n} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}.$$

Denominator layout: the result shape is (n, m)

Gradient of vector and matrix (denominator layout)

- Matrix
 - Denoted as **X**, **Y**, **Z**

$$\frac{\partial \mathbf{Y}}{\partial x} = \begin{bmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{21}}{\partial x} & \dots & \frac{\partial y_{m1}}{\partial x} \\ \frac{\partial y_{12}}{\partial x} & \frac{\partial y_{22}}{\partial x} & \dots & \frac{\partial y_{m2}}{\partial x} \\ \frac{\partial y_{1n}}{\partial x} & \frac{\partial y_{2n}}{\partial x} & \dots & \frac{\partial y_{mn}}{\partial x} \end{bmatrix}$$

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{12}} & \dots & \frac{\partial y}{\partial x_{1q}} \\ \frac{\partial y}{\partial x_{21}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{2q}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{p1}} & \frac{\partial y}{\partial x_{p2}} & \dots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

Gradient of matrix with respect to (w.r.t) matrix,
matrix w.r.t vector, vector w.r.t matrix?
Too complex. Not commonly used.

Gradient table (denominator layout)

Shape of $\frac{\partial y}{\partial x}$	Scalar y	Vector \mathbf{y} (m ,1)	Matrix \mathbf{Y} (m, n)
Scalar x	?	?	?
Vector \mathbf{x} (n ,1)	(n, 1)	?	
Matrix \mathbf{X} (p, q)	?		

HW2: fill in the table with the matrix dimensions of the resulting gradients.

Gradient table: vector by vector (denominator layout)

HW3: fill in the gradient table with the resulting gradients

$y =$	a	x	Ax	$x^T A$
$\frac{\partial y}{\partial x} =$?	?	A^T	?

$y =$	au $u = u(x)$	vu $v = v(x), u = u(x)$	$v + u$ $v = v(x), u = u(x)$	Au $u = u(x)$	$g(u)$ $u = u(x)$
$\frac{\partial y}{\partial x} =$?	$v \frac{\partial u}{\partial x} + \frac{\partial v}{\partial x} u^T$?	?	?

$y =$	a	$u^T v$ $v = v(x), u = u(x)$	$g(u)$ $u = u(x)$	$x^T Ax$
$\frac{\partial y}{\partial x} =$?	?	?	?

Optimization

- Compute the gradient of the loss w.r.t \mathbf{w}
 - Consider a single instance

$$J(\mathbf{w}) = L(\mathbf{x}, y | \mathbf{w}) = \frac{1}{2} \|\tilde{y} - y\|^2 = \frac{1}{2} (\mathbf{w}^T \mathbf{x} - y)^2$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = ?$$

$$z = \mathbf{w}^T \mathbf{x} - y$$

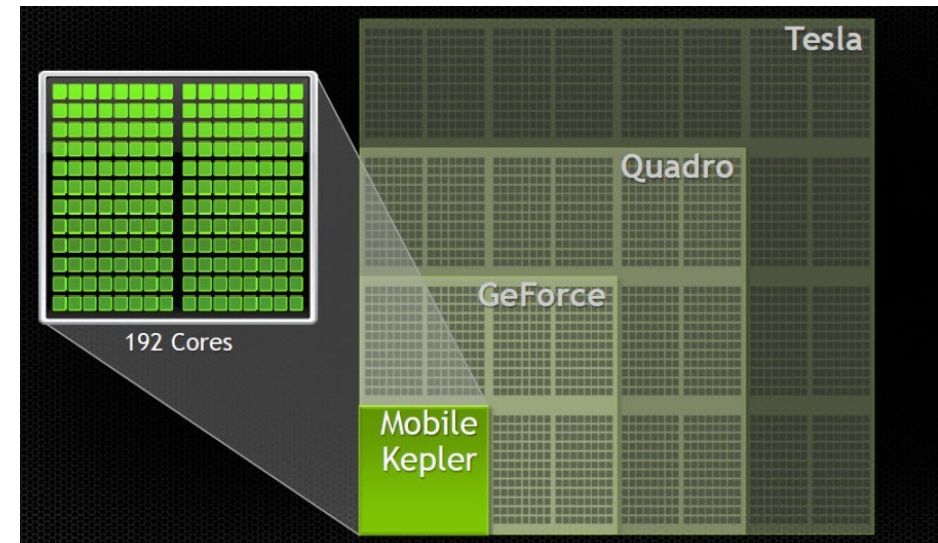
$$J(\mathbf{w}) = \frac{1}{2} z^2$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial \mathbf{w}} = z \mathbf{x} = (\mathbf{w}^T \mathbf{x} - y) \mathbf{x}$$

Shape Check!

Vectorization

- Consider multiple examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}$, $i=1, 2, \dots, m$
- Naïve approach of computing the gradient
 - for each instance $(\mathbf{x}^{(i)}, y^{(i)})$
 - Accumulate the loss $L(\mathbf{x}^{(i)}, y^{(i)})$ into J
 - Average J over m
 - for each instance $(\mathbf{x}^{(i)}, y^{(i)})$
 - Accumulate the gradient of $\frac{\partial L(\mathbf{x}^{(i)}, y^{(i)})}{\partial \mathbf{w}}$
- Not as fast as matrix operations



Vectorization

- Consider multiple examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}$, $i=1, 2, \dots, m$
- Put each instance (the feature vector) into one row of a matrix
 - For fast memory access as Numpy stores data in [row-major format](#)
- Put each target into one element of a column vector

$$X = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \dots \\ \mathbf{x}^{(m)T} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \dots \\ y^{(m)} \end{pmatrix} \quad \tilde{\mathbf{y}} = X\mathbf{w}$$

$J(\mathbf{w}) = ?$

Vectorization

- Consider multiple examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}$, $i=1, 2, \dots, m$
- Put each instance (the feature vector) into one row of a matrix
 - For fast memory access as Numpy stores data in [row-major format](#)
- Put each target into one element of a column vector

$$X = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \vdots \\ \mathbf{x}^{(m)T} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} \quad \tilde{\mathbf{y}} = X\mathbf{w}$$

$$J(\mathbf{w}) = \frac{1}{2m} \|\tilde{\mathbf{y}} - \mathbf{y}\|^2 = \frac{1}{2m} (\tilde{\mathbf{y}} - \mathbf{y})^T (\tilde{\mathbf{y}} - \mathbf{y}) \quad \mathbf{u} = \tilde{\mathbf{y}} - \mathbf{y}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = ?$$

Vectorization

- Consider multiple examples $\{(\mathbf{x}^{(i)}, y^{(i)})\}$, $i=1, 2, \dots, m$
- Put each instance (the feature vector) into one row of a matrix
 - For fast memory access as numpy stores data in [row-major format](#)
- Put each target into one element of a column vector

$$X = \begin{pmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \vdots \\ \mathbf{x}^{(m)T} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} \quad \tilde{\mathbf{y}} = X\mathbf{w}$$

$$J(\mathbf{w}) = \frac{1}{2m} \|\tilde{\mathbf{y}} - \mathbf{y}\|^2 = \frac{1}{2m} (\tilde{\mathbf{y}} - \mathbf{y})^T (\tilde{\mathbf{y}} - \mathbf{y}) \quad \mathbf{u} = \tilde{\mathbf{y}} - \mathbf{y}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2m} \left(\frac{\partial \mathbf{u}}{\partial \mathbf{w}} \mathbf{u} + \frac{\partial \mathbf{u}}{\partial \mathbf{w}} \mathbf{u} \right) = \frac{1}{m} X^T \mathbf{u} = \frac{1}{m} X^T (\tilde{\mathbf{y}} - \mathbf{y}) \quad \text{Shape Check!}$$

Training by gradient descent

- Gradient descent algorithm for optimization
- initialize \mathbf{w} randomly
- Repeat
 - Compute the objective $J(\mathbf{w})$ over all training instances
 - Compute $\frac{\partial J}{\partial \mathbf{w}}$
 - Update $\mathbf{w} = \mathbf{w} - \alpha \frac{\partial J}{\partial \mathbf{w}}$

Notation summary

- Scalar x
- Vector $\mathbf{x} \in \mathbb{R}^{n \times 1}$
 - i-th element of a vector, x_i
 - i-th example $\mathbf{x}^{(i)}$
- Matrix \mathbf{X}
 - i-th row and j-th column X_{ij}
- If we choose denominator layout for $\frac{\partial y}{\partial \mathbf{x}}$ we should lay out the gradient $\frac{\partial y}{\partial \mathbf{x}}$ as a column vector, and $\frac{\partial \mathbf{y}}{\partial x}$ as a row vector.

Summary

- AI vs. ML vs. Deep Learning
- Terminologies & Notations
 - Feature, label, loss
- Univariate & multivariate linear regression
 - Solving parameters via gradient descent
 - Taking gradients of vectors & matrices

References

- [1] Goodfellow Ian, Bengio Yoshua, Courville Aaron. Deep learning. MIT Press. <http://www.deeplearningbook.org>
- [2] Haohan Wang, Bhiksha Raj. On the Origin of Deep Learning. 2017 <https://arxiv.org/abs/1702.07800>
- [3] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng. “Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations.” In ICML 2009
- <http://neuralnetworksanddeeplearning.com/chap1.html>
- <https://www.analyticsvidhya.com/blog/2017/06/a-comprehensive-guide-for-linear-ridge-and-lasso-regression/>
- <https://medium.com/meta-design-ideas/math-stats-and-nlp-for-machine-learning-as-fast-as-possible-915ef47ced5f>

Homework

Theory

1. Solve for the gradients of the functions. (slide 34)
2. Solve for the dimensionality of the gradient vectors (slide 44)
3. Fill out the gradient tables (slide 45).

Practical

1. <https://tinyurl.com/y26qj2ts>

Next Lecture

(Shallow) Neural Networks