



# Surviving the Singularity: Linear Regression

From [Twitch.tv/1bit2far](https://www.twitch.tv/1bit2far)

# Background

Last week we got into basic mathematics and ML fundamentals

This week we will be looking at actual models

Class based off [https://d2l.ai/chapter\\_linear-regression/index.html](https://d2l.ai/chapter_linear-regression/index.html)

You can find old lessons on youtube: @1bit2far

You can find slides and materials at [github/krisciu/SurvivingTheSingularity](https://github.com/krisciu/SurvivingTheSingularity)

Please ask questions





# Welcome to 2024

RING HAS ERADICATED 99.9999% OF HUMANITY

THE UN-CULLED HUMANS LIVE OFF NITREO CUBES

ALL HUMANS CAUGHT BY THE NEO-AUTOMATA ARE PROCESSED  
AT THE CLIPPY-REHABILITATION-UNIT



# Linear Regression

A brief review

Used for numerical analysis

Relates an output:  $Y$  to features:  $X$

Simplest form:  $Y = mX + b$

Output = features  $\times$  weight + bias

Our goal is to find the best fit

---

# Linear Regression: Let's break it down!

$$\hat{y} = w_1x_1 + \dots + w_dx_d + b.$$

**Predictions are based on weighted features + a bias**

$$\hat{y} = \mathbf{w}^\top \mathbf{x} + b.$$

**Predictions are based on the dot product of a vector of features on a vector of weights + bias**

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b,$$

**For a whole Dataset**

Prediction is based on a matrix-vector product of a matrix where rows are examples and columns are features with a vector of weights + bias

# Loss Functions

For linear regression

For regression we usually want to use Squared Error

$$l^{(i)}(\mathbf{w}, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2.$$

We can calculate the quality of a model on a dataset by averaging our losses

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)})^2$$

---

# Gradient Descent and Optimization

Gradient Descent: A method of updating an ML models parameters by minimizing the gradient of the loss function with respect to said parameters

Learning Rate: A [hyperparameter](#) number that represents the step size of the Gradient Descent

Batch Gradient Descent: Most basic implementation, done on the losses of every single datapoint in our dataset

Stochastic Gradient Descent: Gradient Descent done on only one datapoint

MiniBatch Stochastic Gradient Descent: Gradient Descent done on a small batch of samples



# Implement Linear Regression from Scratch



# Neural Networks



## Biological

Dendrite

Nucleus

Axon

Synaptic Weights

Regression can  
be thought of  
as a single layer  
Neural Network

## Technological

Input layer

Model/**Activation Function**

Output layer

Weights/Bias



# Implement Linear Regression EZ Mode

# Analyzing Faults in Results

## Underfitting

Training/validation error substantial

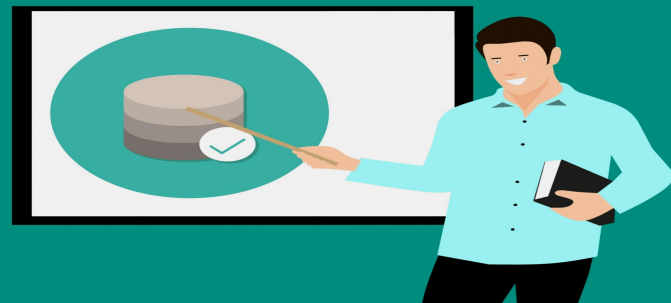
Training/validation very close



## Overfitting

Training data significantly lower than validation error

Limited samples can make this more likely



# Weight Decay

A regularization technique

Regularize our weights

To do this: Add norm as a penalty to minimizing loss

$$L(\mathbf{w}, b) + \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

L2 regularization = ridge regression, L1 regularization = lasso regression

L2 norms distribute weight more evenly

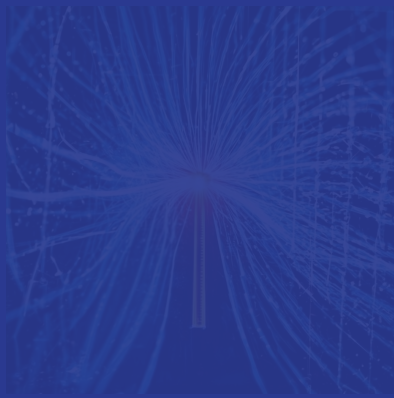
L1 norms concentrate weight on a small set of features

$$L(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(\mathbf{w}, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)})^2$$

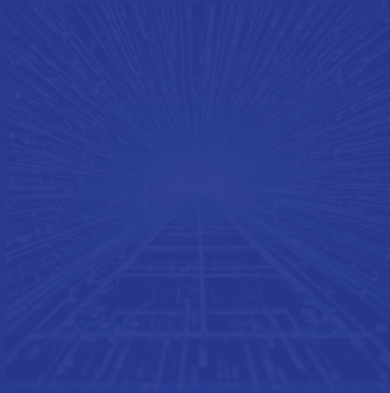
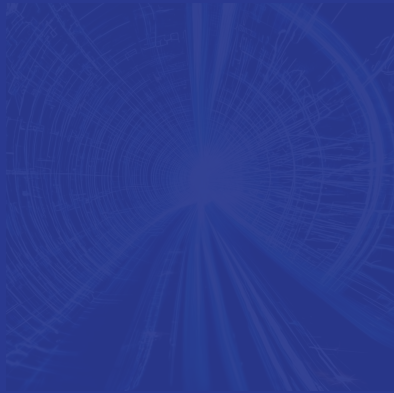
# Weight Decay from Scratch

# Weight Decay EZ mode





# Questions?





# Big Project Brainstorming

# Next up

More Neural Nets

How to softmax

How to classify

---

# Shilling

