# Deep Learning for Data Science DS 542

https://dl4ds.github.io/sp2026/

**Gradient Descent**

# Announcements

- Re: discussion deadlines are moving to 11:59pm on the day of discussion.
  - Why? The practice is more important than the timing.
  - Still targeting $\leq$ 30 minutes to do, but more time if you need/want it.

- Shared Compute Cluster (SCC) Tutorial next Monday.
  - Please bring your laptop to class.
  - No graded exercise, but will be walking through account setup.

# Plan for Today

- Loss functions for multiclass classification (spillover)
- Example of gradient descent
- Basics of gradient descent
- Gradient descent as a statistical process
- Challenges with gradient descent

# Loss Function for Regression

If you recast regression as

1. Predicting the mean of a normal distribution with a fixed variance and
2. Optimize output for maximum likelihood,

Then the optimization is equivalent to optimizing with least squared errors ($L_2$) as your loss function.

# Loss Function for Binary Classification

If you are modeling a binary classification problem,

1. The sigmoid function is handy to map arbitrary "scores" into probabilities, so

2. Your loss function is equivalent to

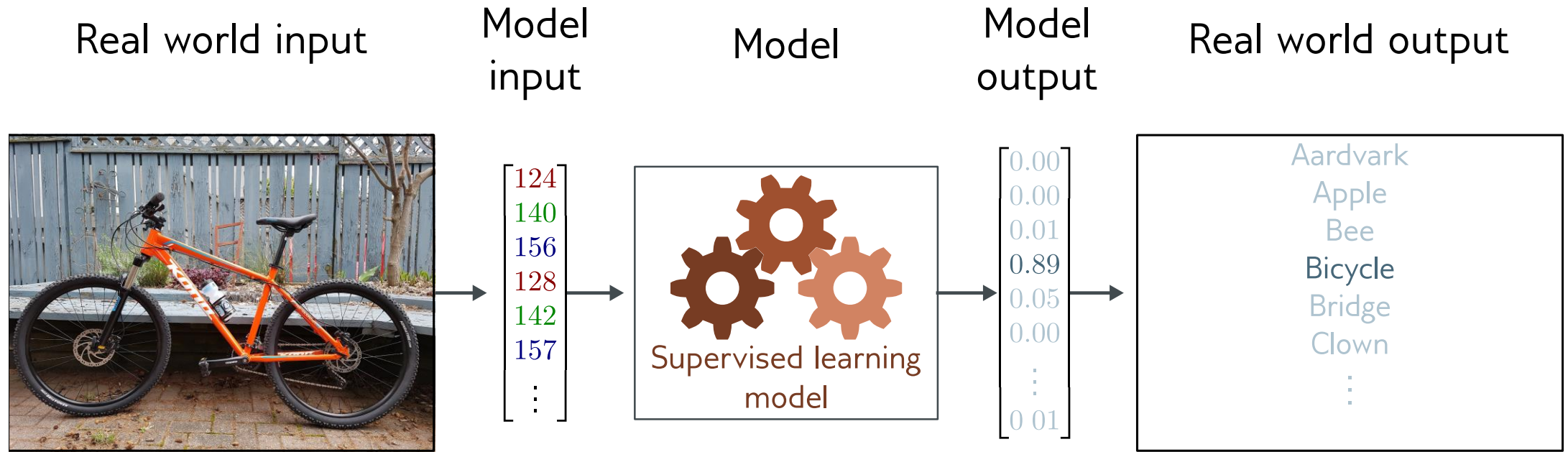$$L[\phi] = \sum_i -(1 - y_i)\log[1 - \text{sig}[f[x_i|\phi]] - y_i \log[\text{sig}[f[x_i|\phi]]$$

# Conceptualizing the Binary Loss Function

$$L[\phi] = \sum_i -(1 - y_i)\log[1 - \text{sig}[f[x_i|\phi]] - y_i\log[\text{sig}[f[x_i|\phi]]$$

$$L[\phi] = \sum_i -(1 - y_i)\log\Pr[y_i = 0|x_i] - y_i\log\Pr[y_i = 1|x_i]$$

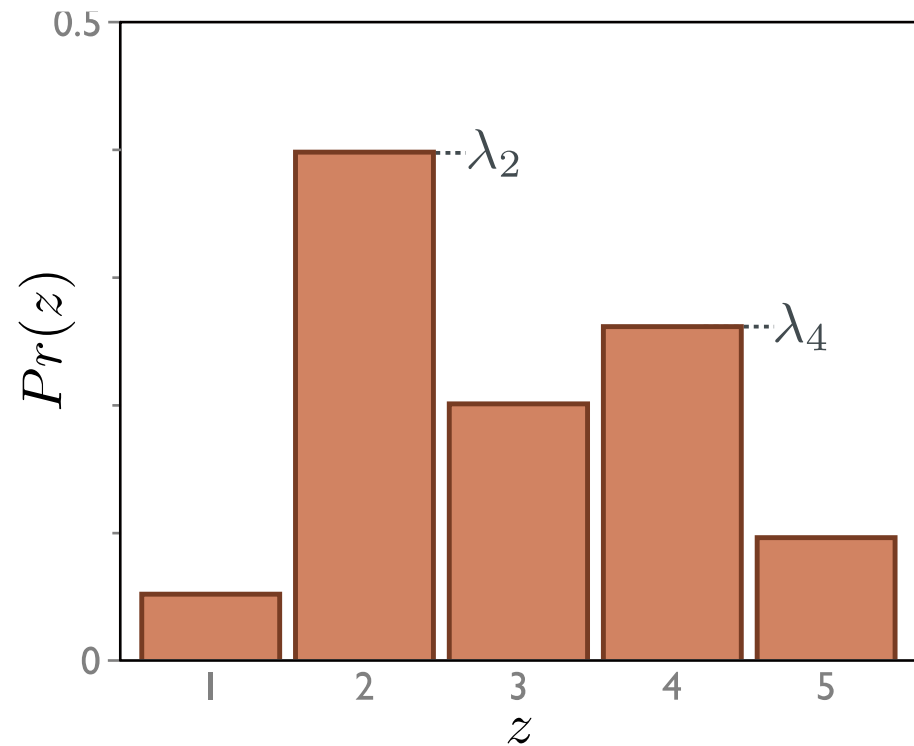$$L[\phi] = \sum_i -\log\Pr[y_i|x_i]$$

# Example 3: multiclass classification

Real world input

Model input

Model

Model output

Real world output



Goal: predict which of K classes $y \in \{1, 2, \ldots, K\}$ the input $x$ belongs to.

# Example 3: multiclass classification

1. Choose a suitable probability distribution $Pr(\mathbf{y}|\boldsymbol{\theta})$ that is defined over the domain of the predictions $\mathbf{y}$ and has distribution parameters $\boldsymbol{\theta}$.

- Domain: $y \in \{1, 2, \ldots, K\}$

- **Categorical distribution**

- $K$ parameters $\lambda_k \in [0, 1]$

- $\sum_k \lambda_k = 1$

$$Pr(y = k) = \lambda_k$$

# Example 3: multiclass classification

2. Set the machine learning model $\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ to predict one or more of these parameters so $\boldsymbol{\theta} = \mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]$ and $Pr(\mathbf{y}|\boldsymbol{\theta}) = Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])$.

Problem:

- Output of neural network can be anything
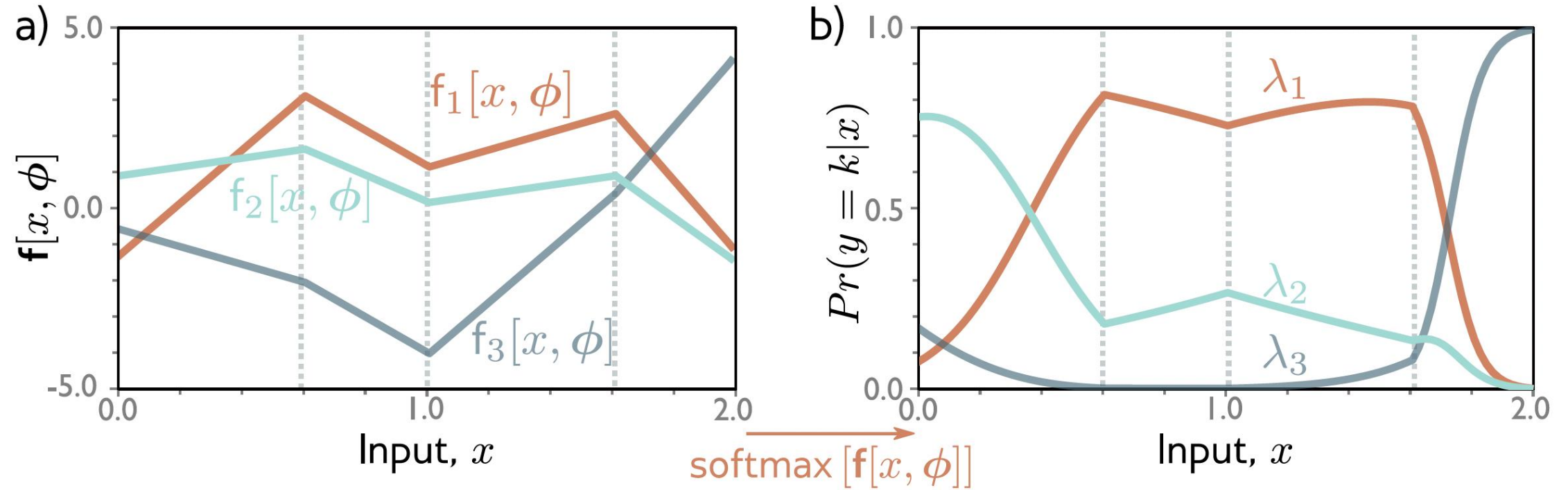- Parameters $\lambda_k \in [0,1]$, sum to one

$$\text{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^{K} \exp[z_{k'}]}$$

Solution:

- Pass through function that maps "anything" to [0,1] and sums to one

$$Pr(y = k|\mathbf{x}) = \text{softmax}_k[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]]$$

# Example 3: multiclass classification



$$Pr(y = k|\mathbf{x}) = \mathrm{softmax}_k[\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}]]$$

# Example 3: multiclass classification

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L[\phi]] = \underset{\phi}{\operatorname{argmin}} \left[ -\sum_{i=1}^{I} \log \left[ Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \qquad (5.7)$$

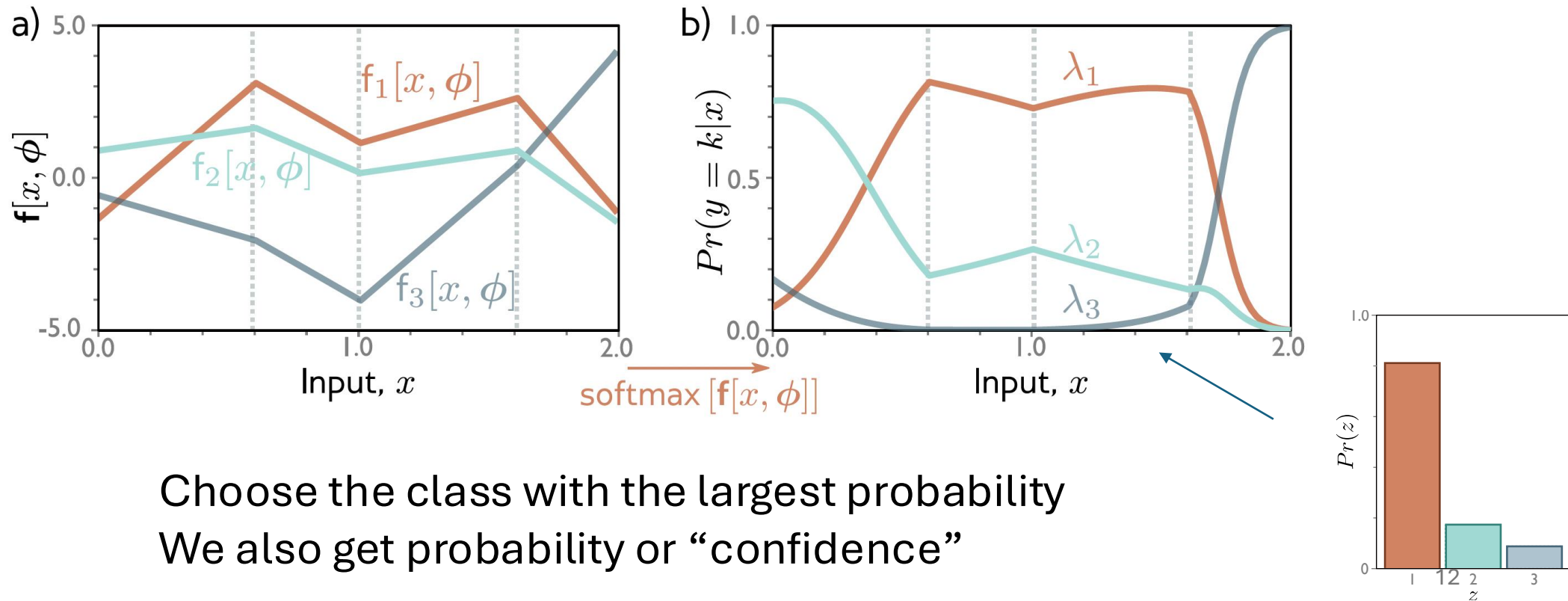$$L[\phi] = -\sum_{i=1}^{I} \log \left[ \operatorname{softmax}_{y_i} \left[ \mathbf{f}\left[ \mathbf{x}_i, \phi \right] \right] \right]$$

$$\operatorname{softmax}_k[\mathbf{z}] = \frac{\exp[z_k]}{\sum_{k'=1}^{K} \exp[z_{k'}]}$$

$$= -\sum_{i=1}^{I} \mathrm{f}_{y_i} \left[ \mathbf{x}_i, \phi \right] - \log \left[ \sum_{k=1}^{K} \exp \left[ \mathrm{f}_k \left[ \mathbf{x}_i, \phi \right] \right] \right]$$

*Multiclass cross-entropy loss*

# Example 3: multiclass classification

4. To perform inference for a new test example $\mathbf{x}$, return either the full distribution $Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \hat{\boldsymbol{\phi}}])$ or the maximum of this distribution.



softmax $[\mathbf{f}[x, \boldsymbol{\phi}]]$

Choose the class with the largest probability
We also get probability or "confidence"

# Any questions?

# Multiple outputs

- Treat each output $y_d$ as *independent*:

$$Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \phi]) = \prod_d Pr(y_d|\mathbf{f}_d[\mathbf{x}_i, \phi])$$

where $\mathbf{f}_d[\mathbf{x}, \phi]$ is the $d^{th}$ set of network outputs

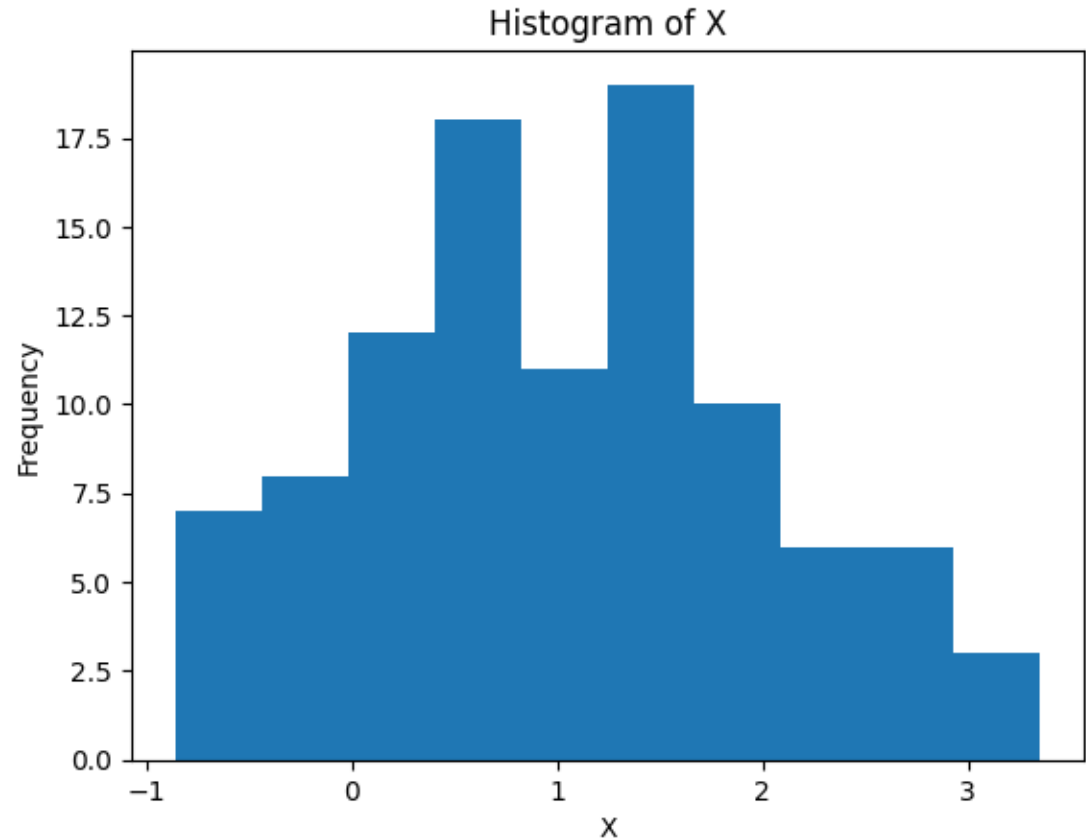- Negative log likelihood becomes sum of terms:

$$L[\phi] = -\sum_{i=1}^{I} \log\Big[Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}_i, \phi])\Big] = -\sum_{i=1}^{I}\sum_{d} \log\Big[Pr(y_{id}|\mathbf{f}_d[\mathbf{x}_i, \phi])\Big]$$

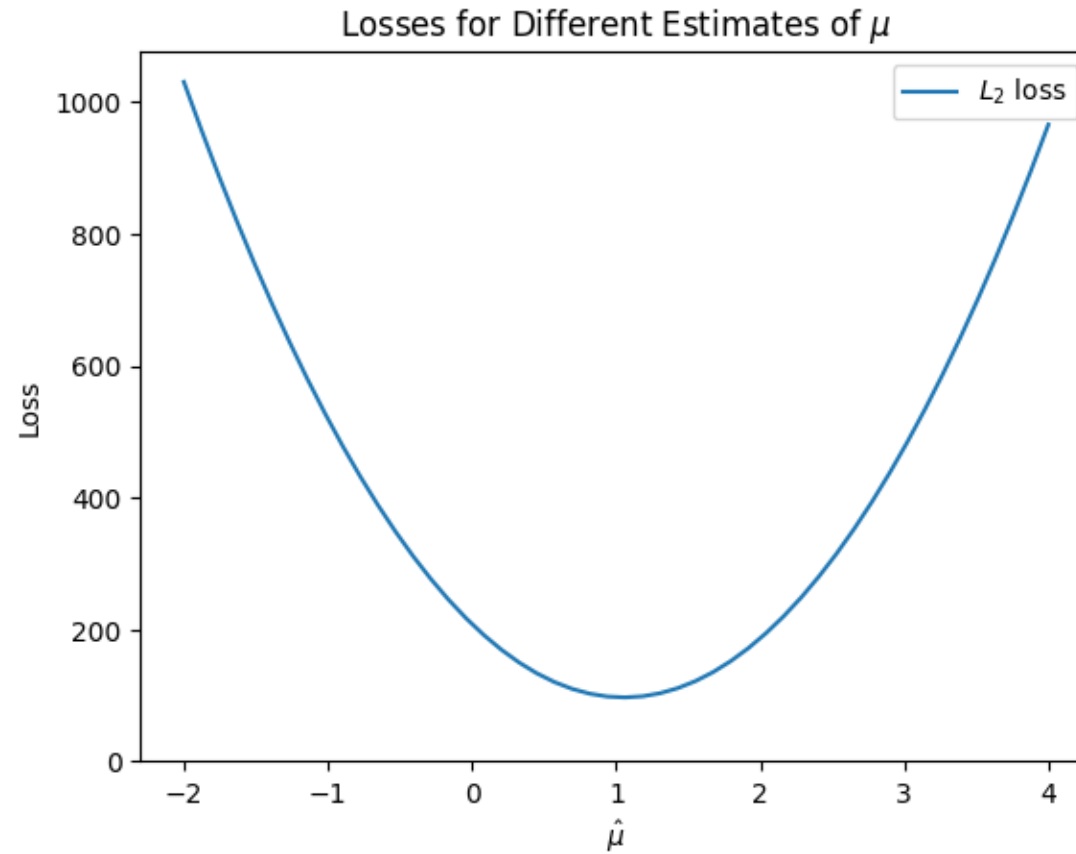$d^{th}$ output of the $i^{th}$ training example

14

# Any questions?

# An Example of Gradient Descent

- X is 100 samples from a normal distribution.
  - What were the parameters of that normal distribution?
  - What was the mean of that normal distribution?



Histogram of X

# Visualizing Gradient Descent
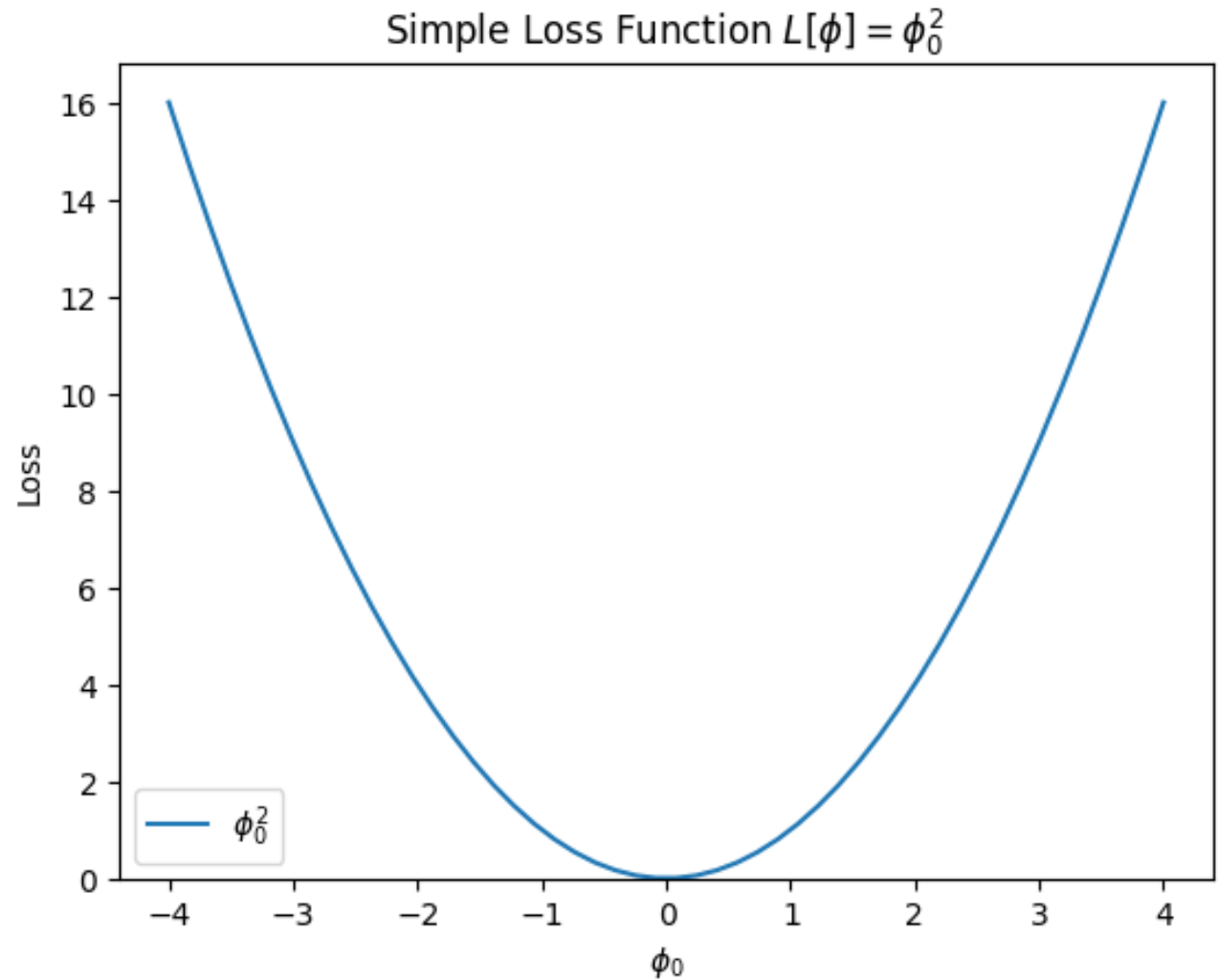


Losses for Different Estimates of $\mu$

# Basics of Gradient Descent

- Given current set of parameters $\phi_t$,

  - Calculate all partial derivatives $\frac{\partial L[\phi]}{\partial \phi_i}$ based on current parameters $\phi_t$.

  - The vector of these $\frac{\partial L[\phi]}{\partial \phi_i}$ is the gradient of the loss function $\nabla L[\phi]$.

  - Update $\phi_{t+1} = \phi_t - \alpha \nabla L[\phi]$

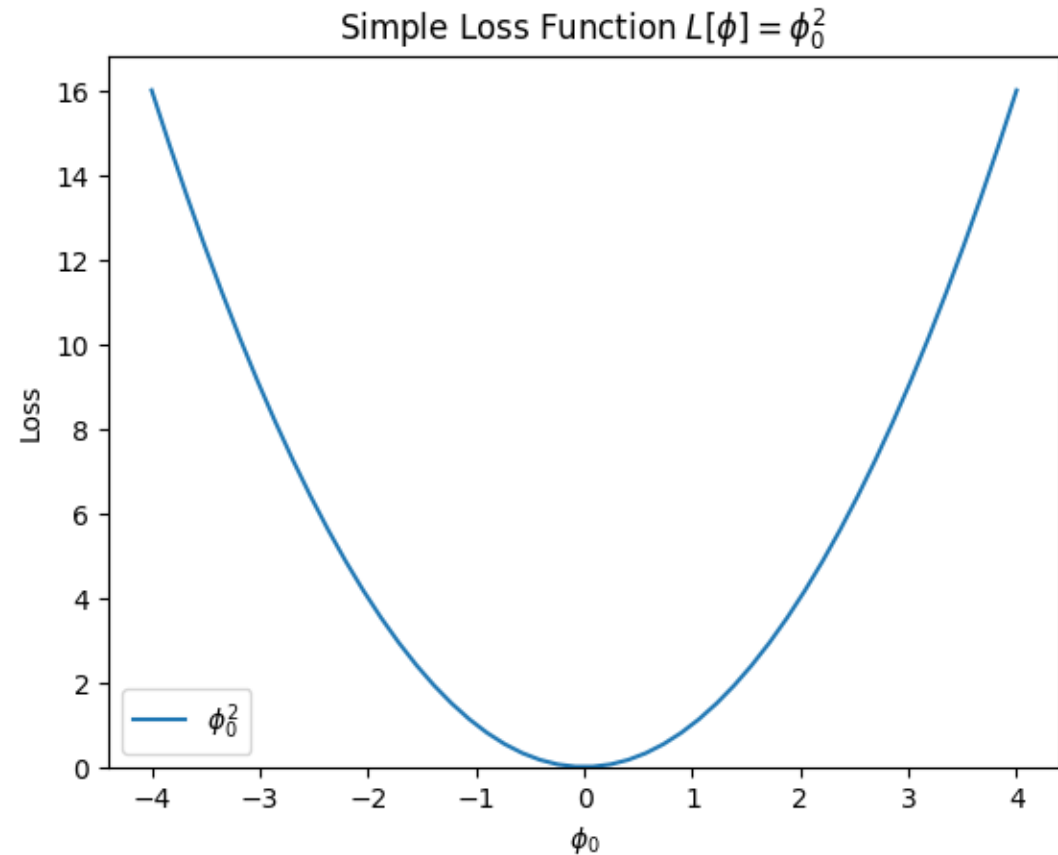    where $\alpha$ is the learning rate.

# What should the learning rate be?

- Try $\alpha = 1$.

- Too small, and it takes many steps to get close.

- Too big, and it overshoots.



Simple Loss Function $L[\phi] = \phi_0^2$

# Convex Loss Functions

- Generally, a lot easier to optimize…

- With gradient descent, main issue is not overshooting minimum too much.

Simple Loss Function $L[\phi] = \phi_0^2$
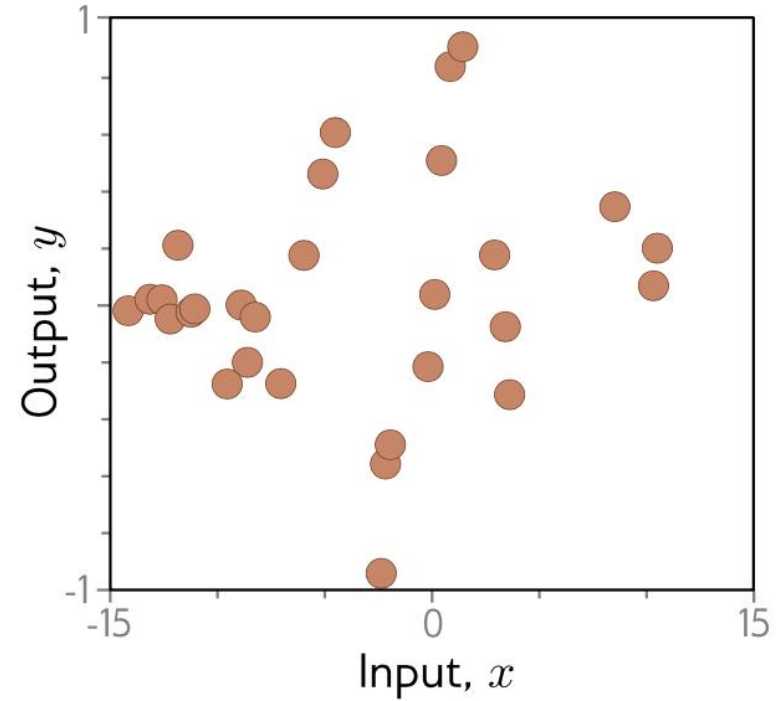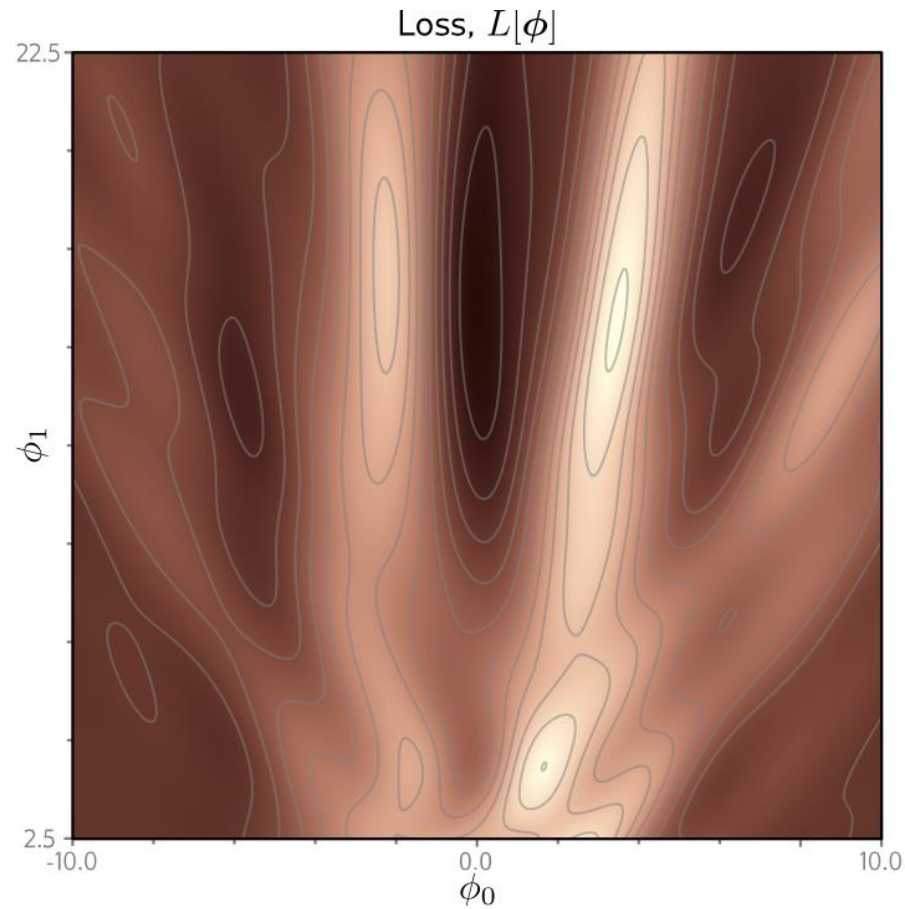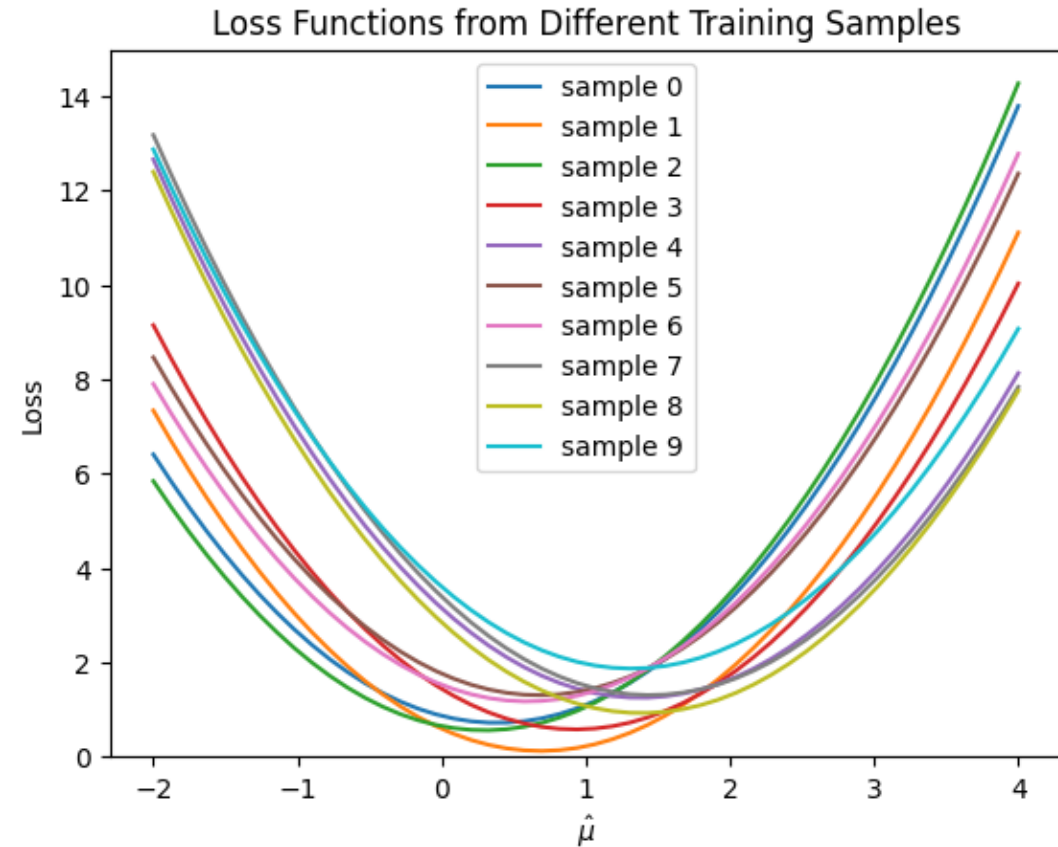
# Non-Convex Loss Functions



Image Source: Understanding Deep Learning, via https://udlbook.github.io/udlfigures/

# Any questions?

# Gradient Descent as a Statistical Process

- Our training data is a sample of the whole population.
  - Different training samples yield different training loss functions.



Loss Functions from Different Training Samples

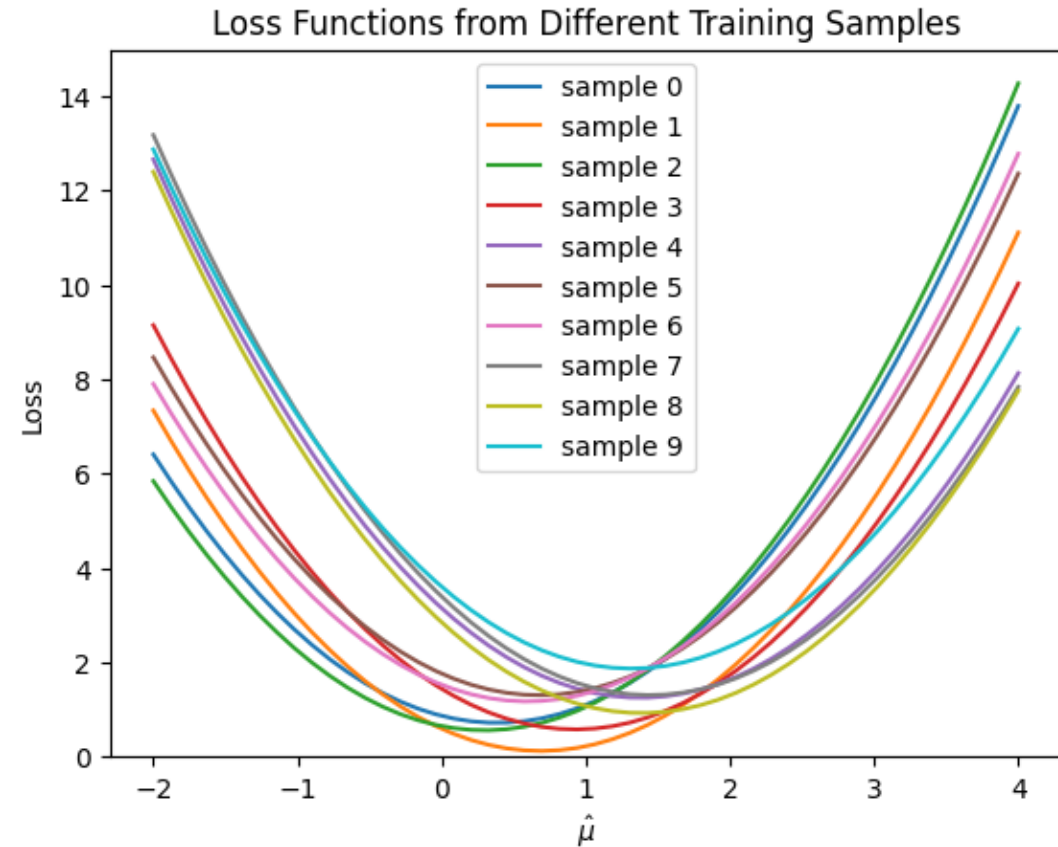# Loss Functions for Different Training Samples

- If we collect different training data sets, will we get different models?

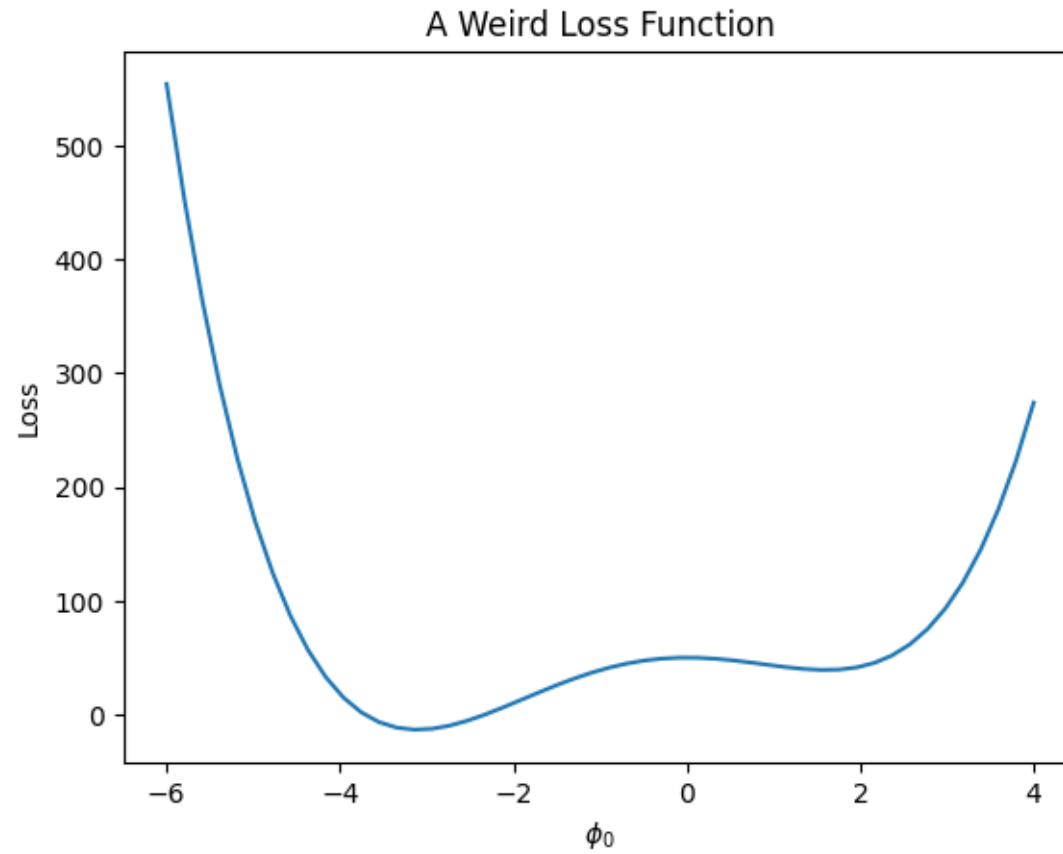# Loss Functions for Samples of the Training Set

- If we sample the training data, will we get different models?

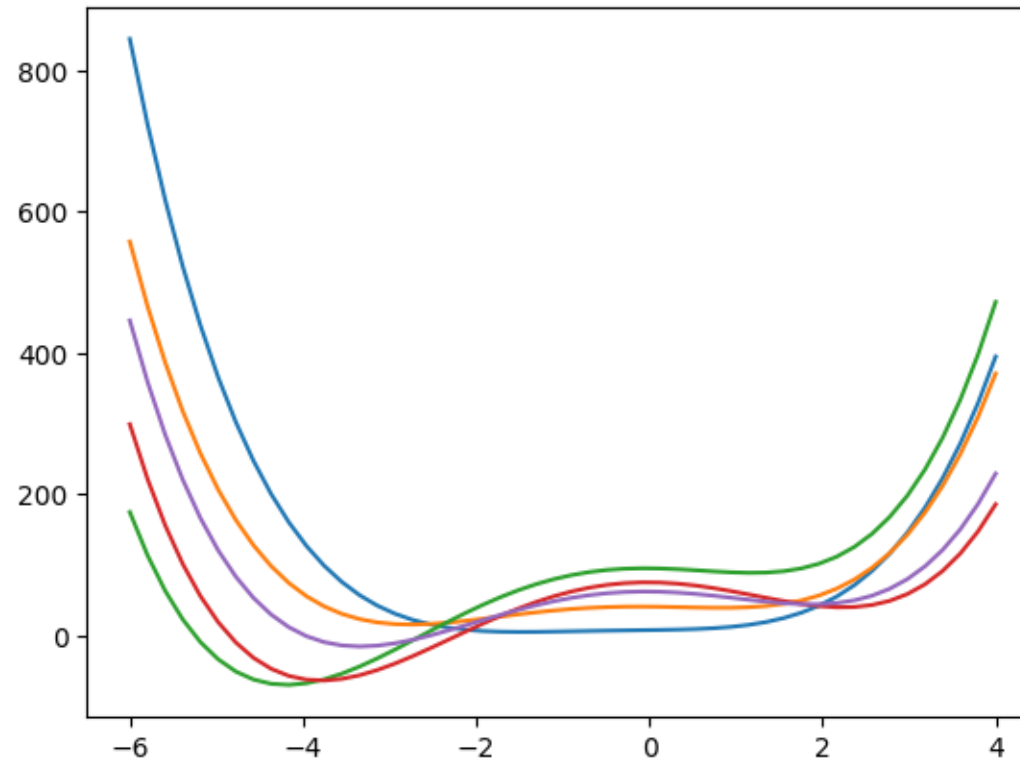# Comparing Models with Different Training Samples

- How far apart are the models of these samples?

- Where do they agree and disagree?



Loss Functions from Different Training Samples

# A Weird Loss Function

# Local Minima vs Samples of the Training Set

# Stochastic Gradient Descent

Idea: Run gradient descent with "mini batches" instead of the full training set.

- E.g. pick a random partition of data into 10 equal-sized batches.
- One epoch = running through all the data once.
  - Vanilla gradient → one parameter update.
  - Stochastic gradient descent → one parameter update per mini batch.

# Variation in Sampled Gradients

- Expected mini batch gradient = whole training set gradient.
  - On average, they agree.
  - But with noise from sampling.


- But remember, just taking one step with each mini batch.
  - Not optimizing to mini batch minimum loss.

# Local Minima vs Stochastic Gradient Descent

- When far from a local minima, mini batches tend to agree on gradient direction.


- When close to a local minima, mini batches disagree more.
  - Sampling noise.
  - Explore the flat area around the minima.

# Speed of Stochastic Gradient Descent

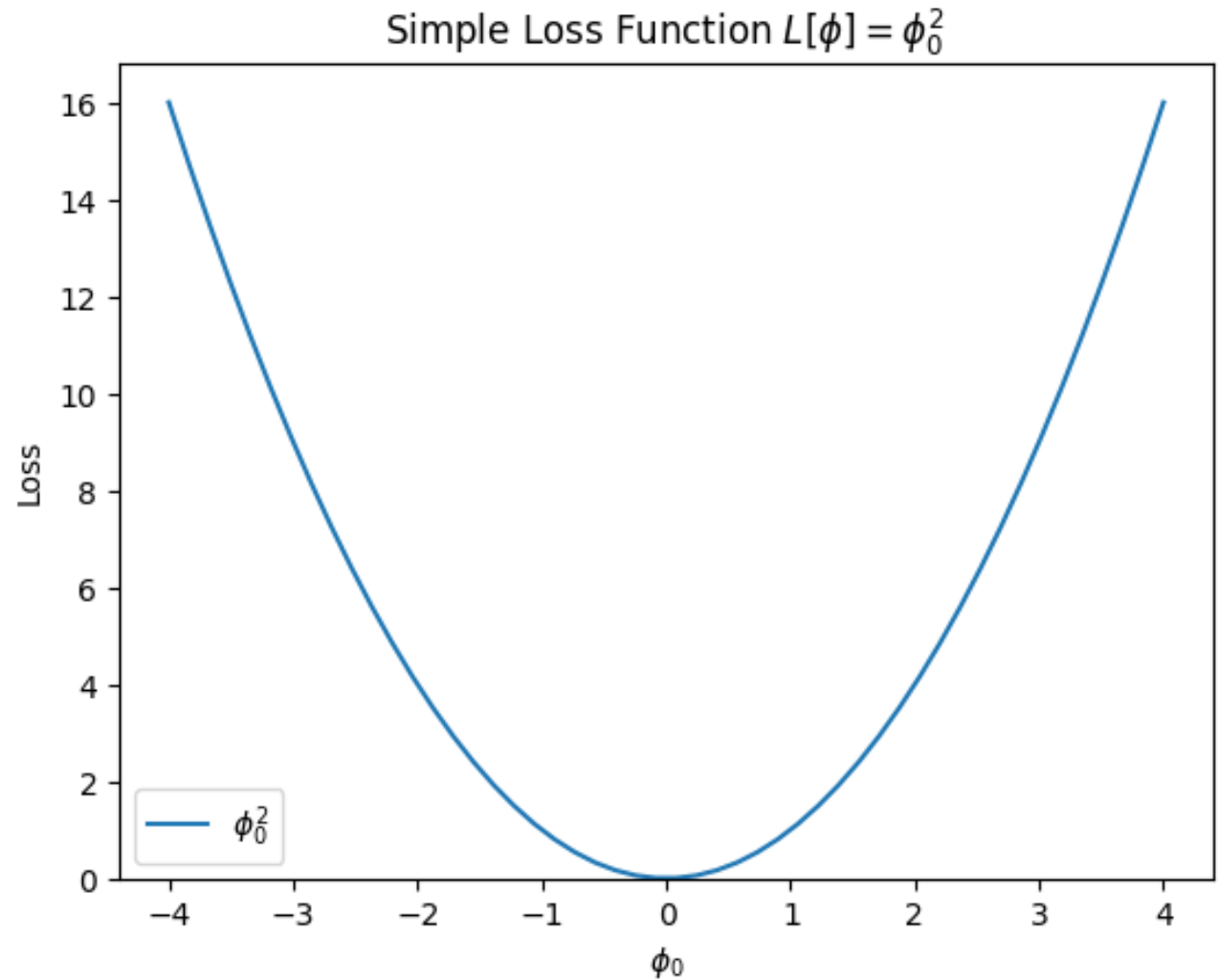- How fast is this compared to vanilla gradient descent?

# Any questions?

# Gradient Descent as a Universal Algorithm

- What's the catch?

# How do we pick Learning Rate?

- Remember, $\alpha = 1$ gives an infinite loop.

- Also, be impatient.



Simple Loss Function $L[\phi] = \phi_0^2$

# Really Bad Linear Regression

- $f(x) = f_1\left(f_2\left(f_3(f_4(x))\right)\right)$

- $f_1(x) = a_1 x + b_1$
- $f_2(x) = a_2 x + b_2$
- $f_3(x) = a_3 x + b_3$
- $f_4(x) = a_4 x + b_4$

- $f(x)$ is just a linear function?

# Really Bad Linear Regression (part 2)

- $f(x) = f_1\left(f_2\left(f_3(f_4(x))\right)\right)$

- $f_1(x) = a_1 x + b_1$
- $f_2(x) = a_2 x + b_2$
- $f_3(x) = a_3 x + b_3$
- $f_4(x) = a_4 x + b_4$

- $f(x)$ is just a linear function?

- Initialize all parameters to zero.

- What are the gradients?

# Really Bad Linear Regression (part 3)

- $f(x) = f_1\left(f_2\left(f_3(f_4(x))\right)\right)$

- $f_1(x) = a_1 x + b_1$
- $f_2(x) = a_2 x + b_2$
- $f_3(x) = a_3 x + b_3$
- $f_4(x) = a_4 x + b_4$

- $f(x)$ is just a linear function?

- Initialize all parameters to 100.

- What are the gradients?

# Really Bad Linear Regression (part 4)

- $f(x) = f_1\left(f_2\left(f_3(f_4(x))\right)\right)$

- $f_1(x) = a_1 x + b_1$
- $f_2(x) = a_2 x + b_2$
- $f_3(x) = a_3 x + b_3$
- $f_4(x) = a_4 x + b_4$

- $f(x)$ is just a linear function?

- We will see both these problems with neural networks if we use the wrong initialization.

# Any questions?