

Training neural networks

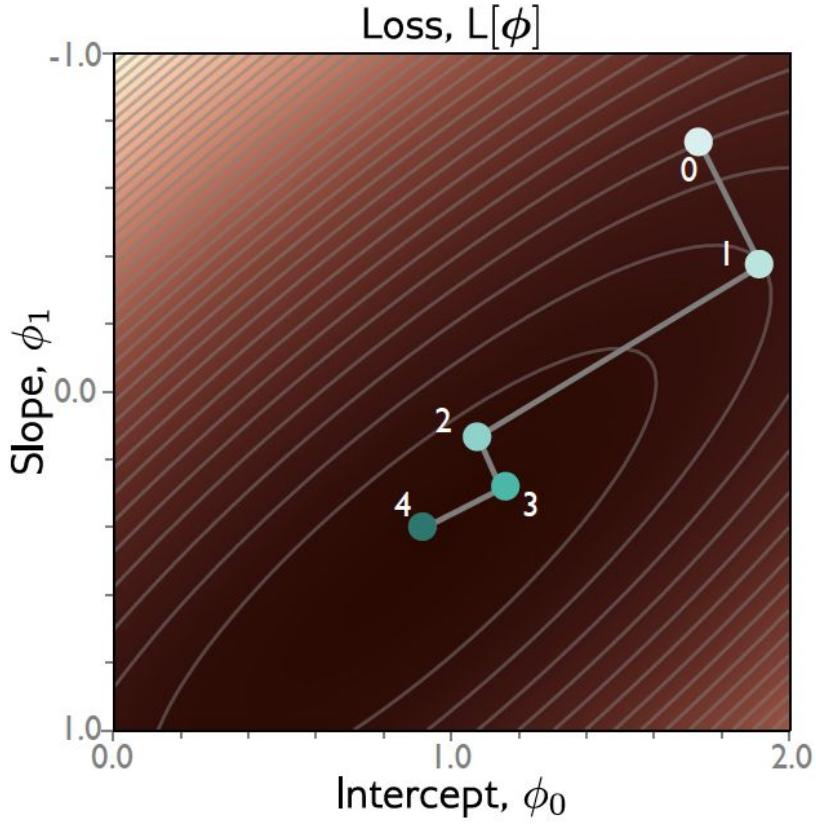
Florin Gogianu
Bitdefender ML team

* using slides, illustrations and ideas from [prof. Simon Prince UDLB, Stanford CS231n, LeCun, Canziani, NYU Deep Learning](#)

Reading

- Lecture 1: chapters 1, 2, 8.1-3, 9.1, Prince, *Understanding Deep Learning*
- Lecture 2: chapters 3 and 4, UDL. LeCun, Bengio, Hinton, *Deep Learning* (paper)
- Lecture 3:
 - chapters 5, 6 and 7, UDL.
 - Goh, [*Why Momentum Really Works?*](#)

Training a simple model



1. Define a loss function
2. Compute the change in parameters required to make the loss smaller
3. Apply the change and get new parameters
4. Repeat from (2)

Loss functions

Loss function

- Training dataset of I pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I$$

- Loss function or cost function measures how bad model is:

$$L \left[\underbrace{\phi, f[\mathbf{x}, \phi]}_{\text{model}}, \underbrace{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I}_{\text{train data}} \right]$$

Loss function

- Training dataset of I pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I$$

- Loss function or cost function measures how bad model is:

$$L[\phi]$$

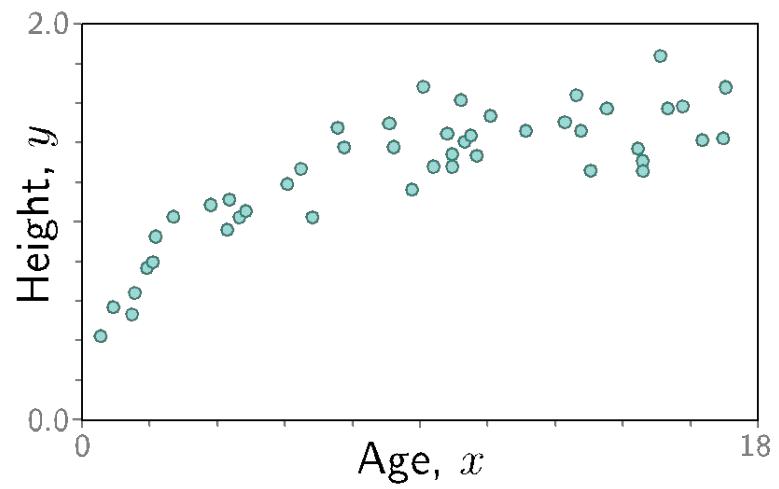
Returns a scalar that is smaller
when model maps inputs to
outputs better

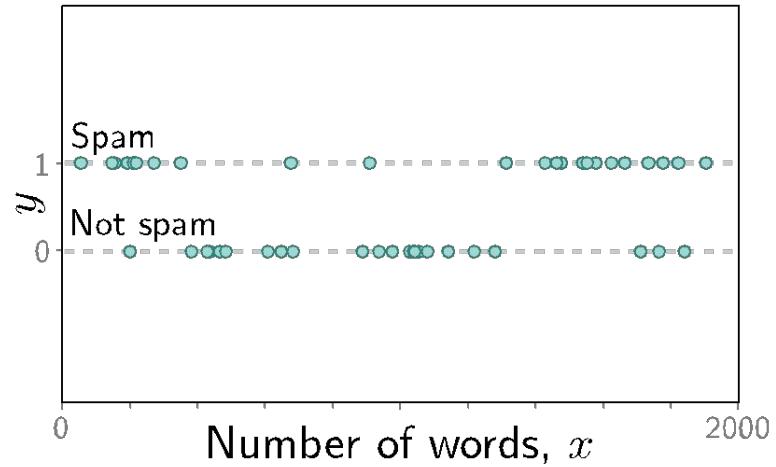
Loss function as an optimization objective

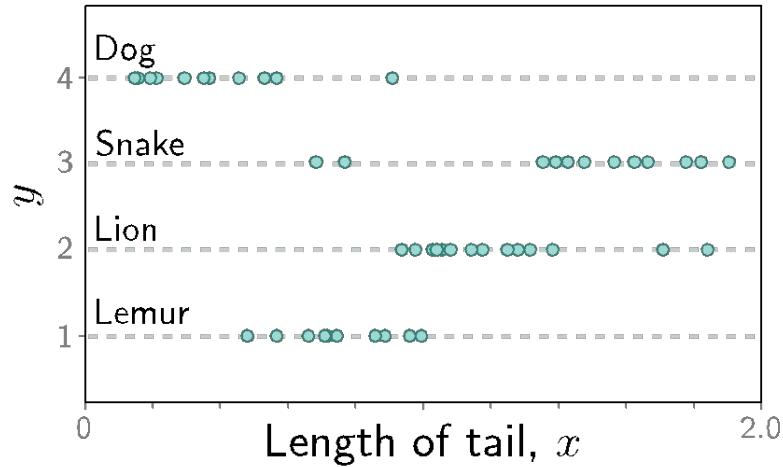
$$L \left[\underbrace{\phi, f[\mathbf{x}, \phi]}_{\text{model}}, \underbrace{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^I}_{\text{train data}} \right]$$

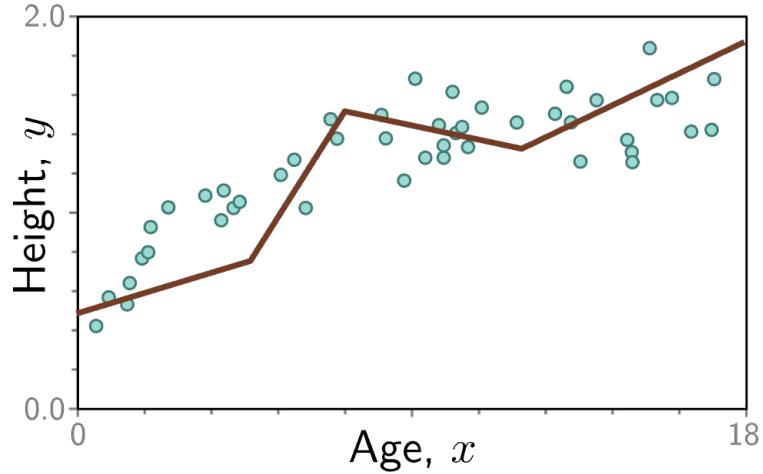
Find the parameters that minimize the loss:

$$\phi = \operatorname{argmin}_{\phi} [L(\phi)]$$



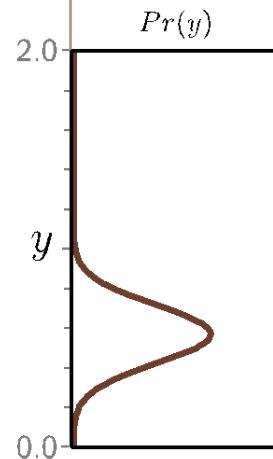
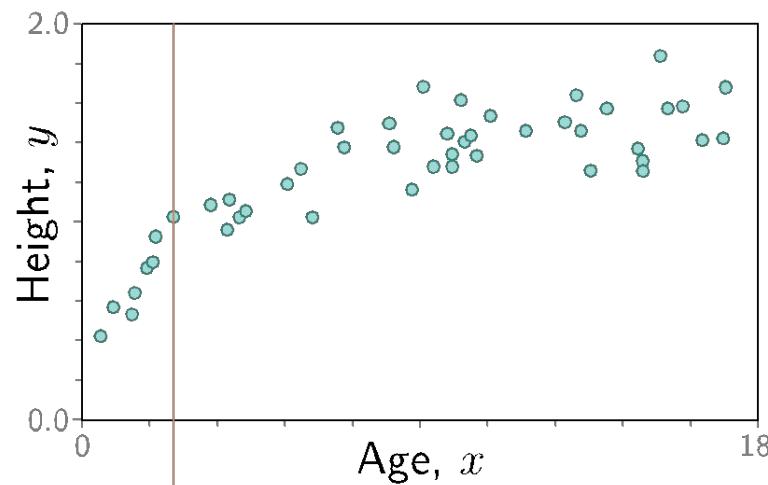


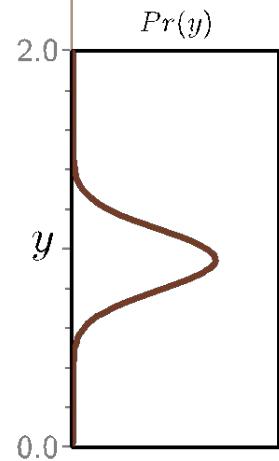
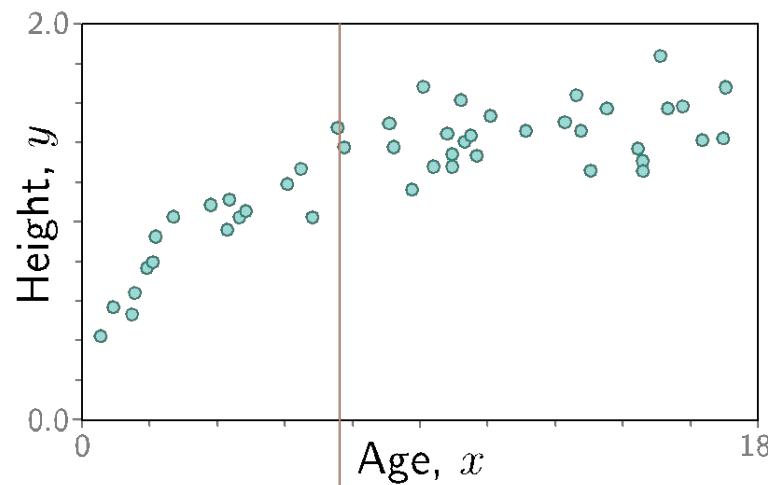


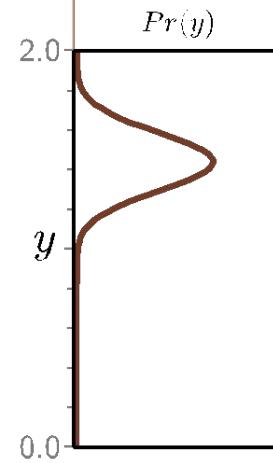
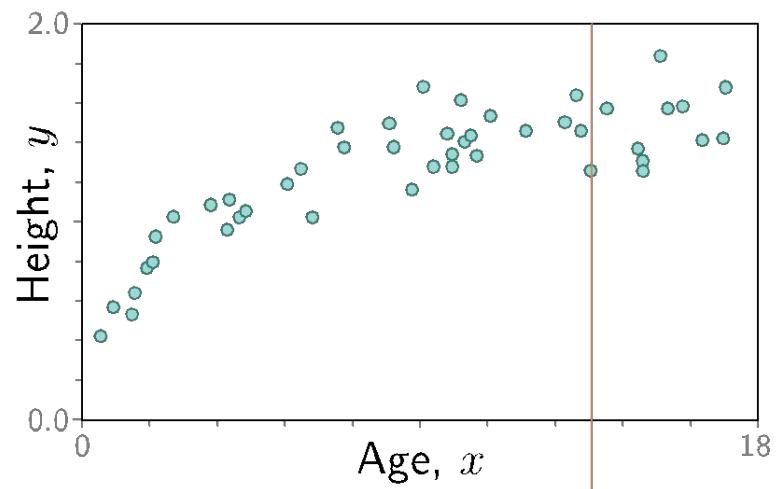


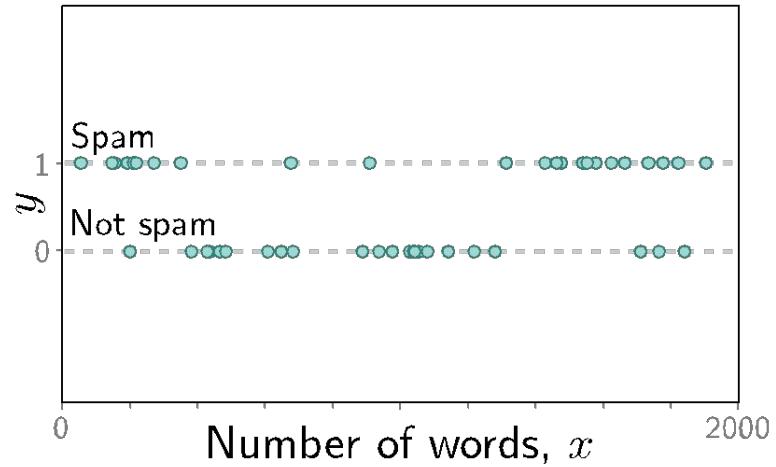
?

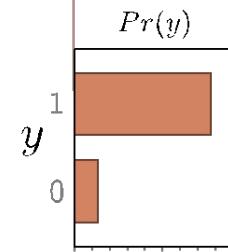
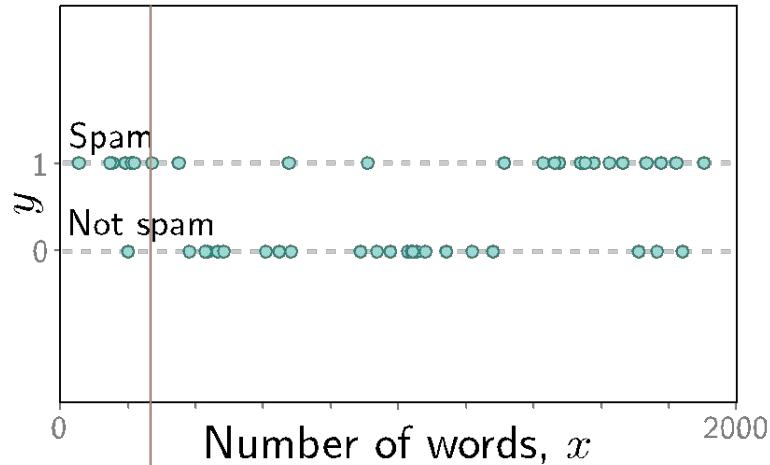
What do we mean,
probabilistically,
about fitting $y = f(x)$?

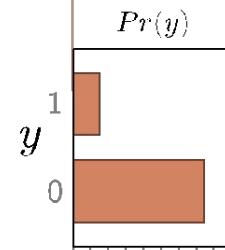
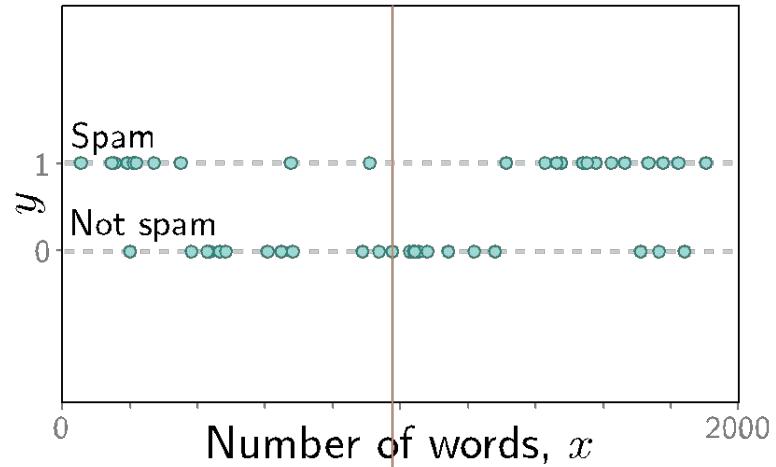


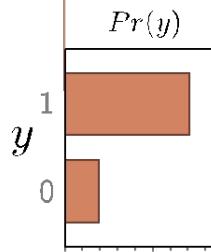
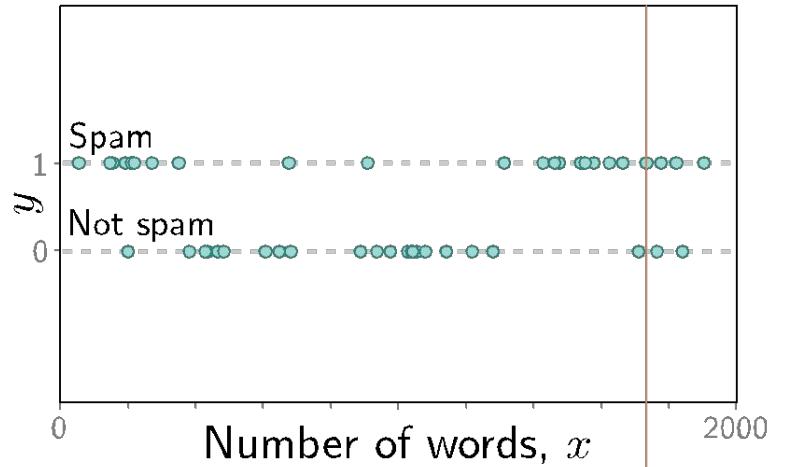


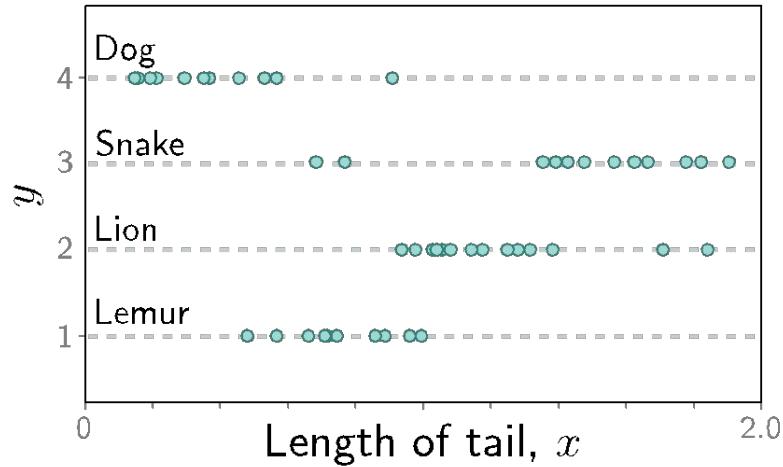


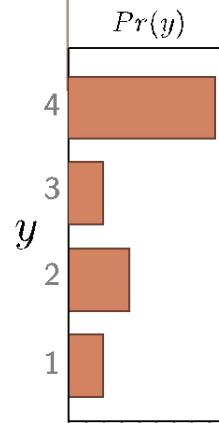
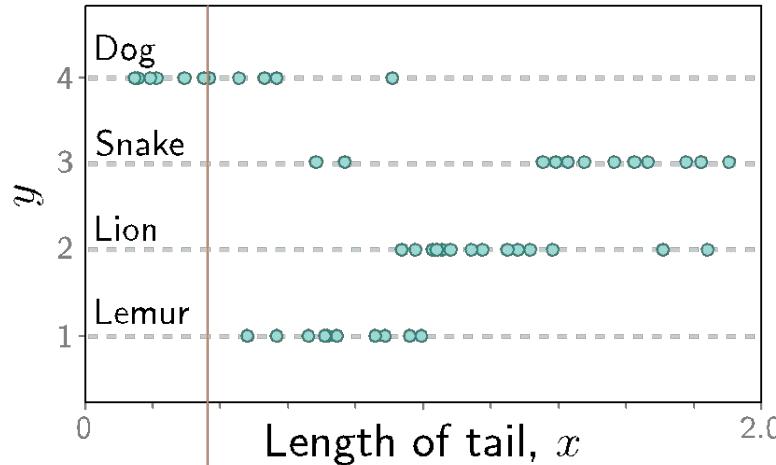


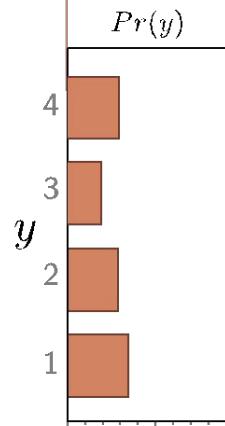
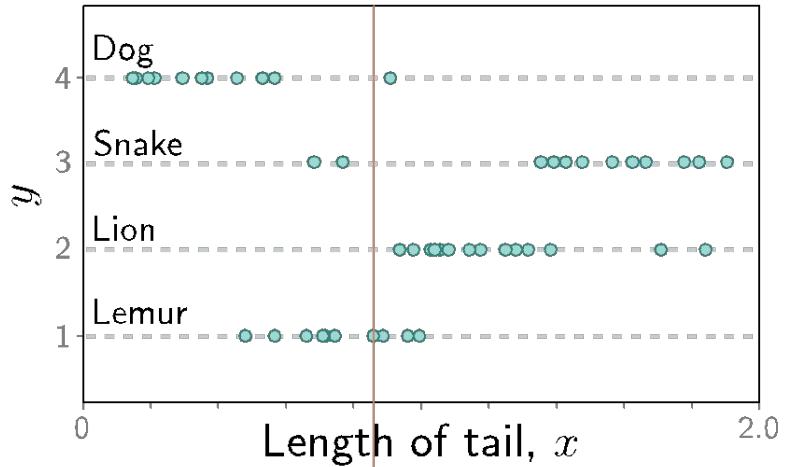


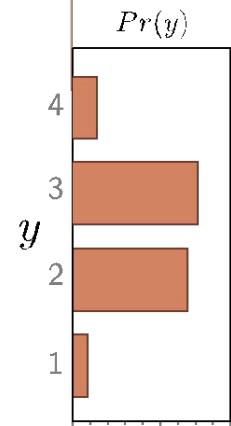
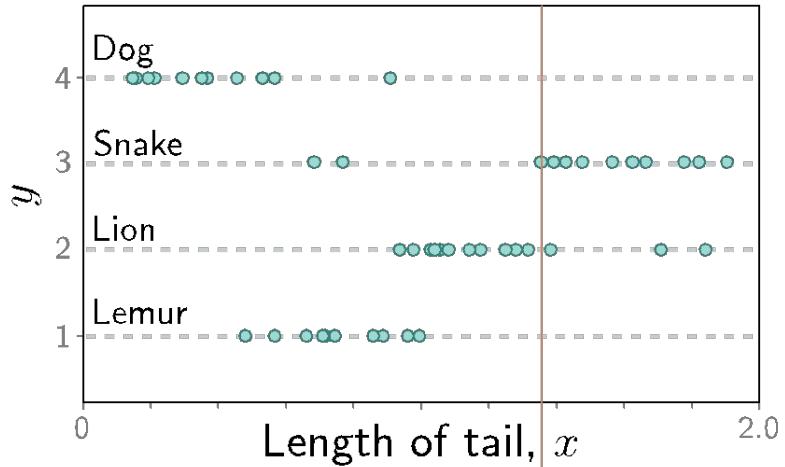


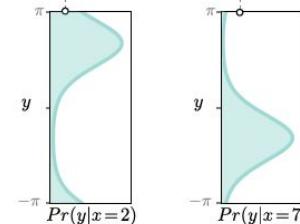
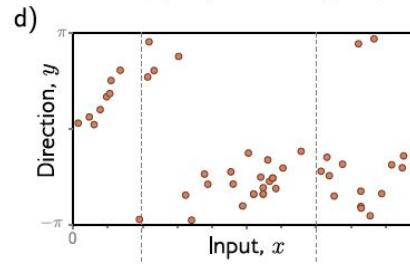
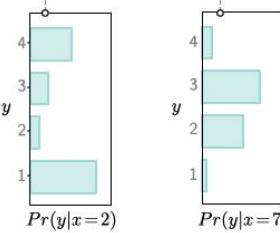
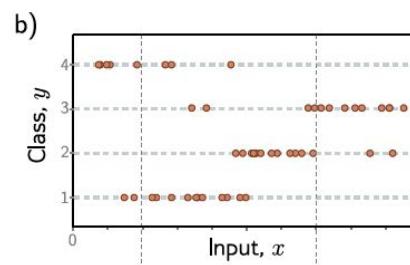
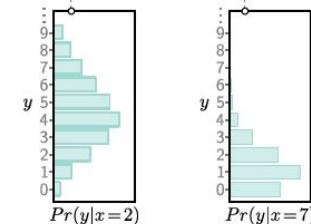
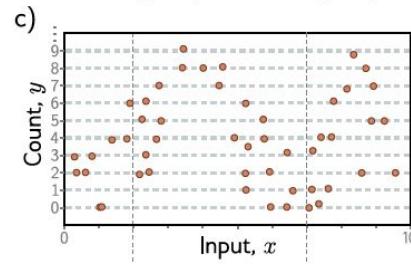
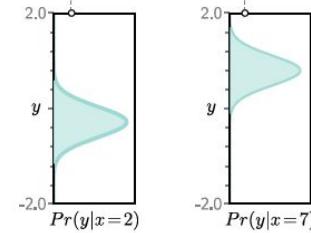
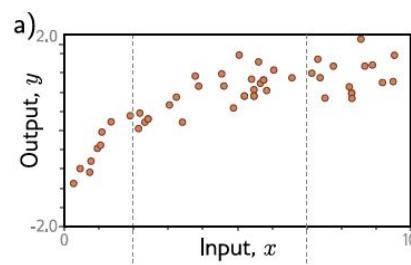












Encourage each prediction (model output) y_i
to have high *probability conditioned* on (model input) x_i

A general recipe for constructing loss functions

- ~~Model predicts output y given x~~
- Model predicts a *conditional probability* distribution:

$$Pr(y | x)$$

- Loss function aims to make the outputs have high probability under the conditional distribution

From model output to probability distribution

- Pick a known distribution (eg.: Normal) to model output y with parameters $\theta = \{\mu, \sigma^2\}$

$$Pr(y | x)$$

- Fit the model so that it predicts the parameters μ and σ of the distribution.

Maximum likelihood criterion

$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{x}_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \boldsymbol{\theta}_i) \right] \\ &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right].\end{aligned}$$

When we consider this probability as a function of the parameters ϕ , we call it a likelihood.

Problem:

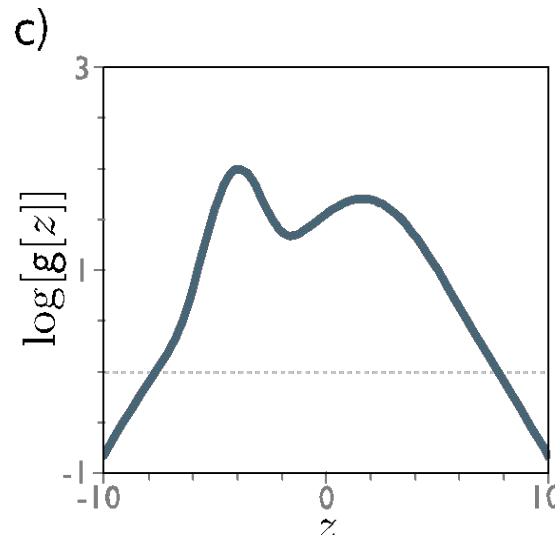
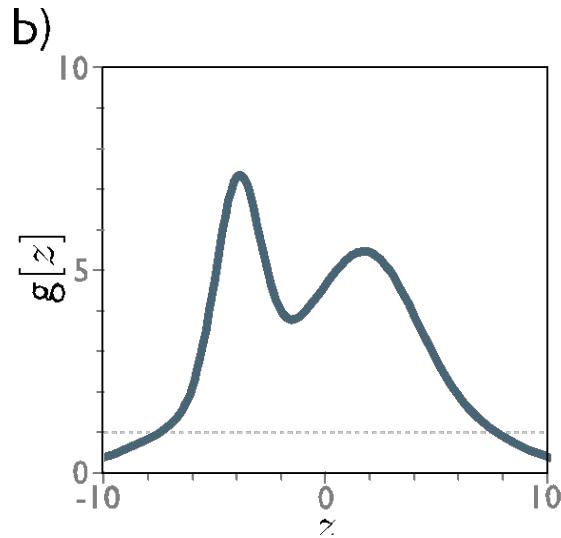
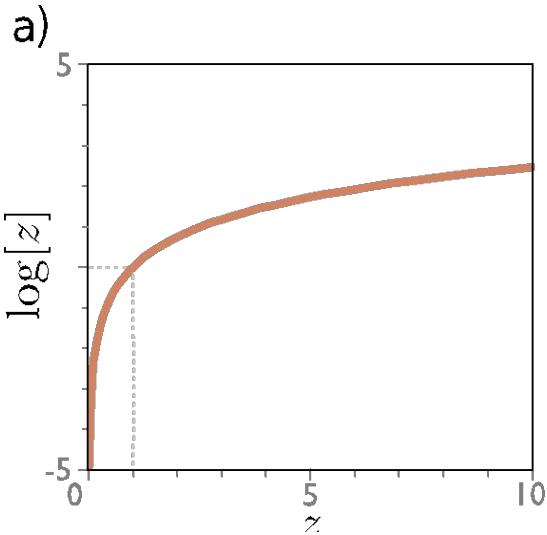
$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right].$$

- The terms in this product might all be small
- The product might get so small that we can't easily represent it

Maximum log likelihood

$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \\ &= \operatorname{argmax}_{\phi} \left[\log \left[\prod_{i=1}^I Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\ &= \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log \left[Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]\end{aligned}$$

The log function is monotonic



Minimizing negative log likelihood

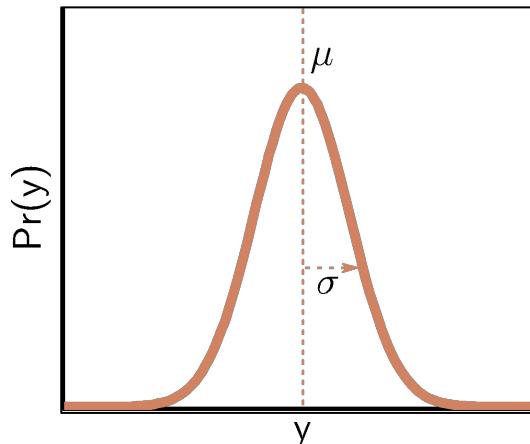
By convention, we minimize things (i.e., a loss)

$$\begin{aligned}\hat{\phi} &= \operatorname{argmax}_{\phi} \left[\sum_{i=1}^I \log \left[Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[- \sum_{i=1}^I \log \left[Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right] \\ &= \operatorname{argmin}_{\phi} \left[L[\phi] \right]\end{aligned}$$

Inference (prediction)

- But now we predict a probability distribution
- We need an actual prediction (point estimate)
- Find the peak of the probability distribution (i.e., mean for normal)

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} [Pr(\mathbf{y}|\mathbf{f}[\mathbf{x}, \boldsymbol{\phi}])]$$



From MLE to common loss functions

Go through 5.3 to 5.5 to convince yourselves that *minimizing the negative log likelihood* amounts to minimizing:

- Mean squared error (univariate regression / Gaussian)
- Binary cross-entropy (binary classification / Bernoulli)
- Cross-entropy (multiclass classification / Categorical)

Example: 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 $\boldsymbol{\theta}$.

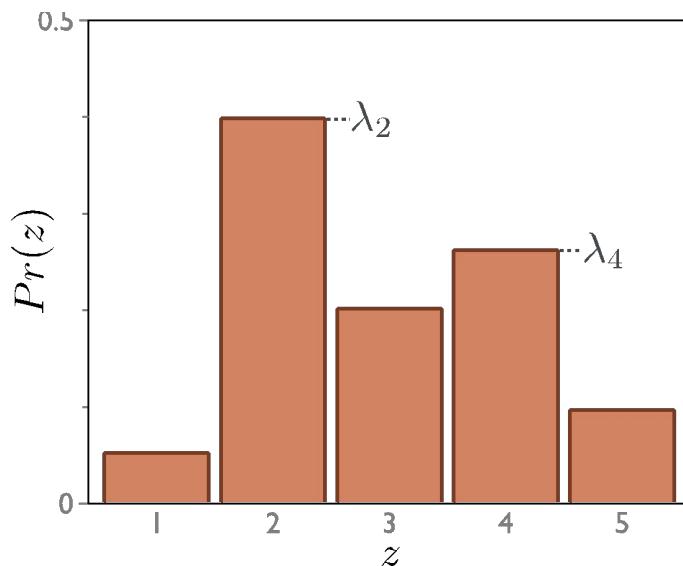
- **Domain:** $y \in \{1, 2, \dots, K\}$

- **Categorical distribution**

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

- **Sum of all parameters = 1**

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



Example: multiclass classification

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

Problem:

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

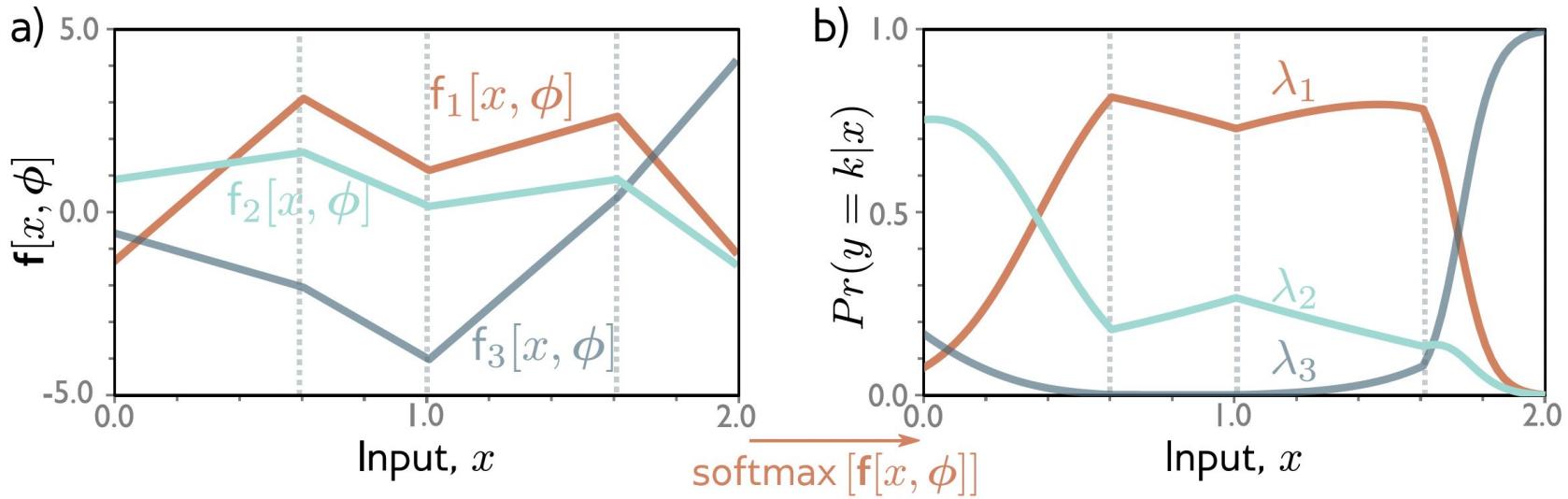
Solution:

- Pass through function that maps “anything” to $[0,1]$, sum to one

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

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

Example: multiclass classification



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

Example: multiclass classification

3. To train the model, find the network parameters $\hat{\phi}$ that minimize the negative log-likelihood loss function over the training dataset pairs $\{\mathbf{x}_i, \mathbf{y}_i\}$:

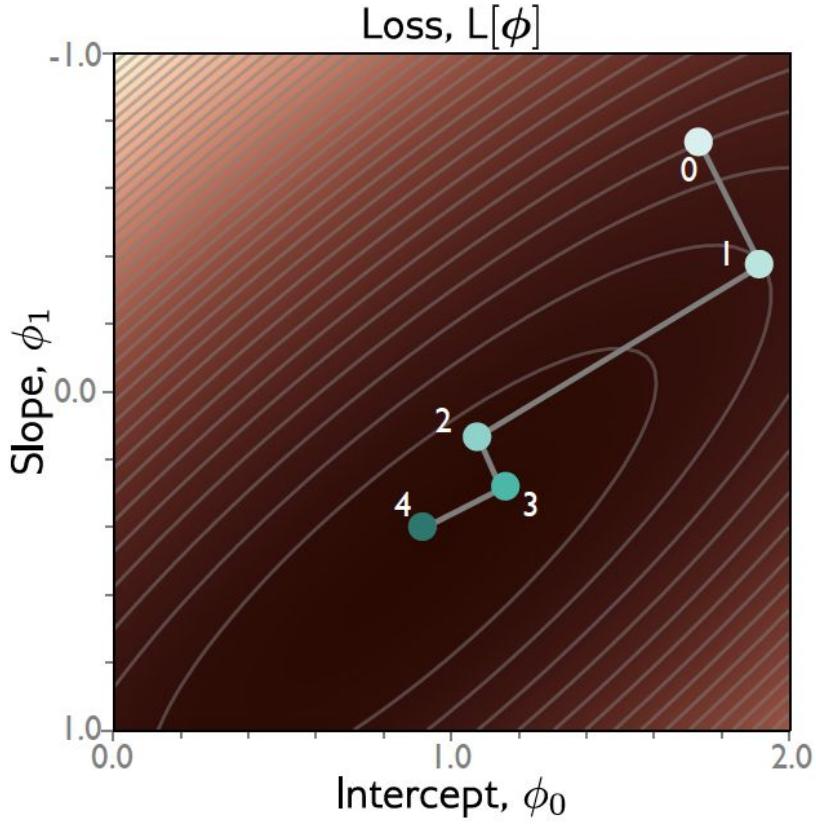
$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L[\phi]] = \underset{\phi}{\operatorname{argmin}} \left[- \sum_{i=1}^I \log \left[Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]. \quad (5.7)$$

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

*Multiclass cross-entropy loss

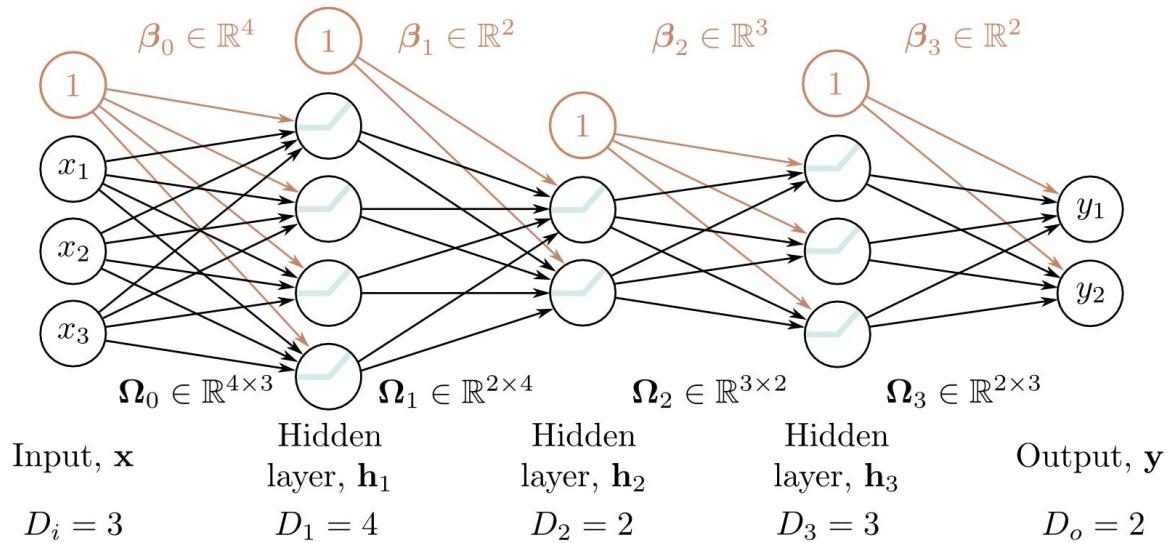
Computing gradients

Training a simple model



1. Define a loss function
2. Compute the change in parameters required to make the loss smaller
3. Apply the change and get new parameters
4. Repeat from (2)

Neural network



$$\mathbf{h}_1 = \mathbf{a}[\beta_0 + \Omega_0 \mathbf{x}]$$

$$\mathbf{h}_2 = \mathbf{a}[\beta_1 + \Omega_1 \mathbf{h}_1]$$

$$\mathbf{h}_3 = \mathbf{a}[\beta_2 + \Omega_2 \mathbf{h}_2]$$

$$\mathbf{f}[\mathbf{x}, \phi] = \beta_3 + \Omega_3 \mathbf{h}_3$$

Setup

Loss, sum of individual terms:

$$L[\phi] = \sum_{i=1}^I \ell_i = \sum_{i=1}^I l[f[\mathbf{x}_i, \phi], y_i]$$

SGD algorithm

$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi}$$

Parameters

$$\phi = \{\beta_0, \Omega_0, \beta_1, \Omega_1, \beta_2, \Omega_2, \beta_3, \Omega_3\}$$

How to compute gradients?

$$\frac{\partial \ell_i}{\partial \beta_k} \quad \text{and} \quad \frac{\partial \ell_i}{\partial \Omega_k}$$

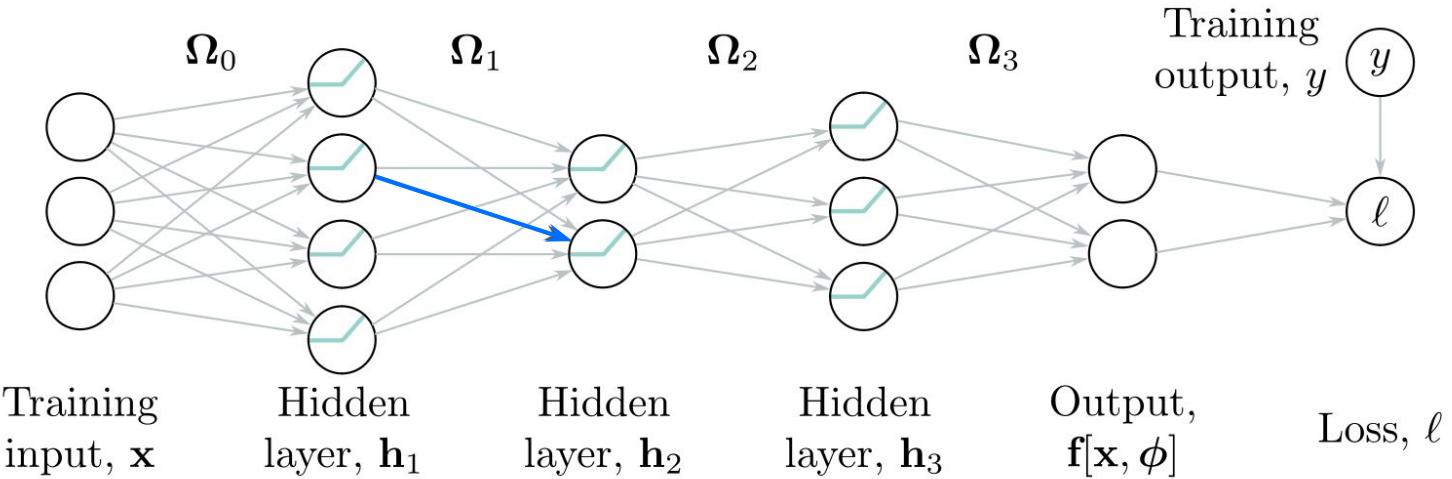
Big deal?

$$\begin{aligned}y' = & \phi'_0 + \phi'_1 a [\psi_{10} + \psi_{11} a[\theta_{10} + \theta_{11}x] + \psi_{12} a[\theta_{20} + \theta_{21}x] + \psi_{13} a[\theta_{30} + \theta_{31}x]] \\& + \phi'_2 a [\psi_{20} + \psi_{21} a[\theta_{10} + \theta_{11}x] + \psi_{22} a[\theta_{20} + \theta_{21}x] + \psi_{23} a[\theta_{30} + \theta_{31}x]] \\& + \phi'_3 a [\psi_{30} + \psi_{31} a[\theta_{10} + \theta_{11}x] + \psi_{32} a[\theta_{20} + \theta_{21}x] + \psi_{33} a[\theta_{30} + \theta_{31}x]]\end{aligned}$$

Huge equation, and we need to compute derivatives:

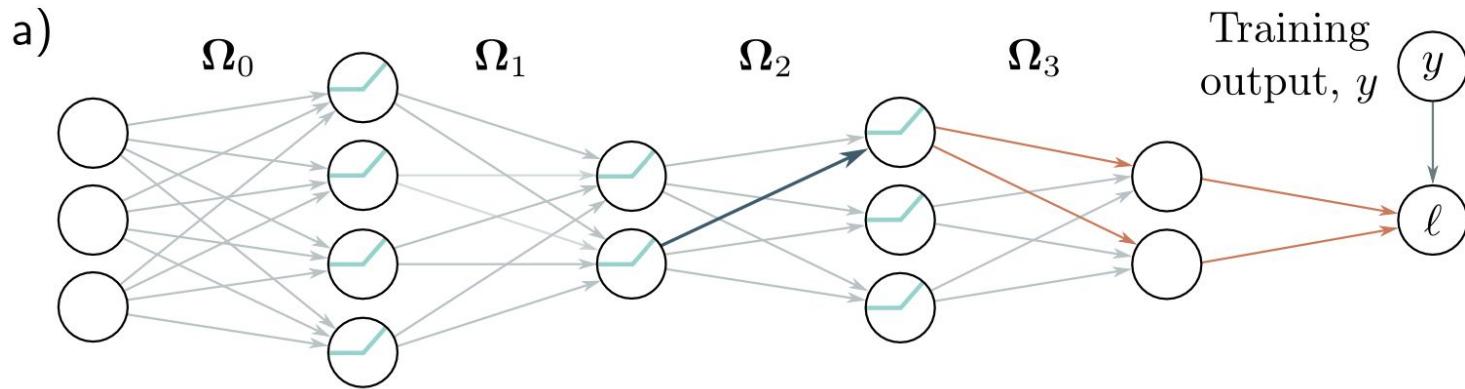
- for every parameter
- for every point in the batch
- for every iteration of SGD

Backpropagation algorithm. Forward pass



- Blue weight multiplies activation (ReLU output) in previous layer
- We want to know how change in **blue weight** affects loss
- If we double activation in previous layer, weight will have twice the effect
- Conclusion: **we need to know the activations at each layer.**

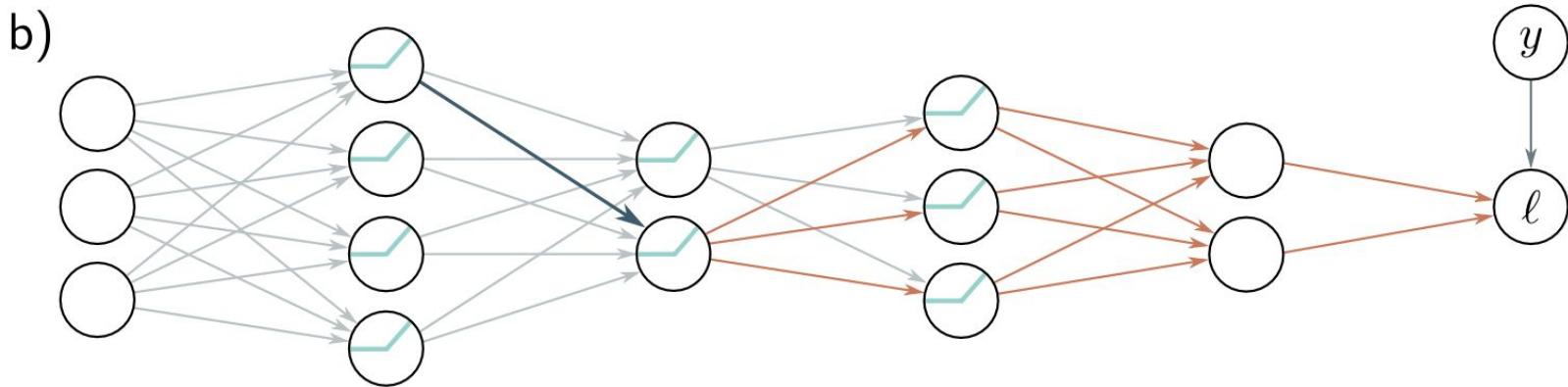
Backpropagation algorithm. Backward pass



To compute how a small change in a weight feeding into \mathbf{h}_3 modifies the loss, we need:

- How \mathbf{h}_3 changes the model output
- How the output changes the loss

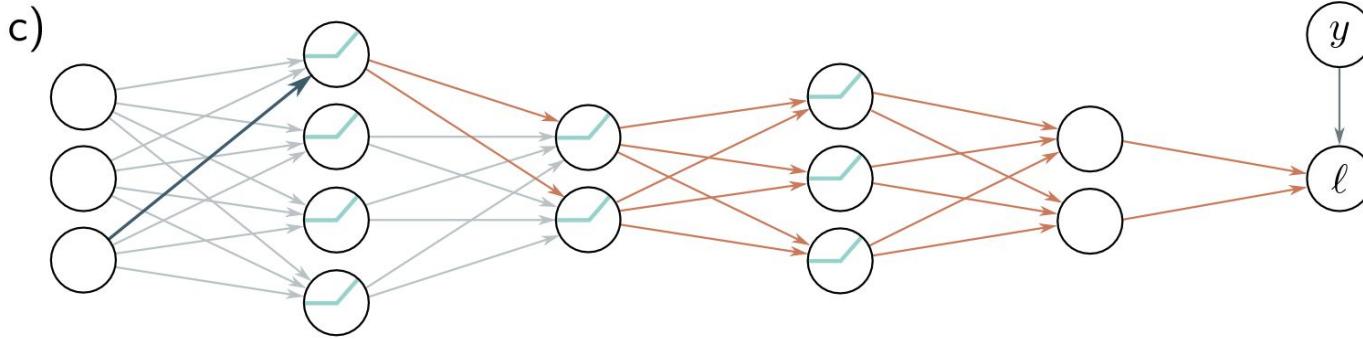
Backpropagation algorithm. Backward pass



To compute how a small change in a weight feeding into \mathbf{h}_2 modifies the loss, we need:

- How a change in layer \mathbf{h}_2 affects \mathbf{h}_3
- How \mathbf{h}_3 changes the model output
- How the output changes the loss

Backpropagation algorithm. Backward pass



To compute how a small change in a weight feeding into \mathbf{h}_2 modifies the loss, we need:

- How a change in layer \mathbf{h}_1 affects \mathbf{h}_2
- How a change in layer \mathbf{h}_2 affects \mathbf{h}_3
- How \mathbf{h}_3 changes the model output
- How the output changes the loss

Toy example

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x] \right] \right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

- A series of functions composed with each other
- Unlike in neural networks it consists of scalars and not vectors and matrices
- The “activation functions” are just \sin , \exp , \cos

Toy example

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos\left[\beta_2 + \omega_2 \cdot \exp\left[\beta_1 + \omega_1 \cdot \sin[\beta_0 + \omega_0 \cdot x]\right]\right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

Derivatives of the activation functions

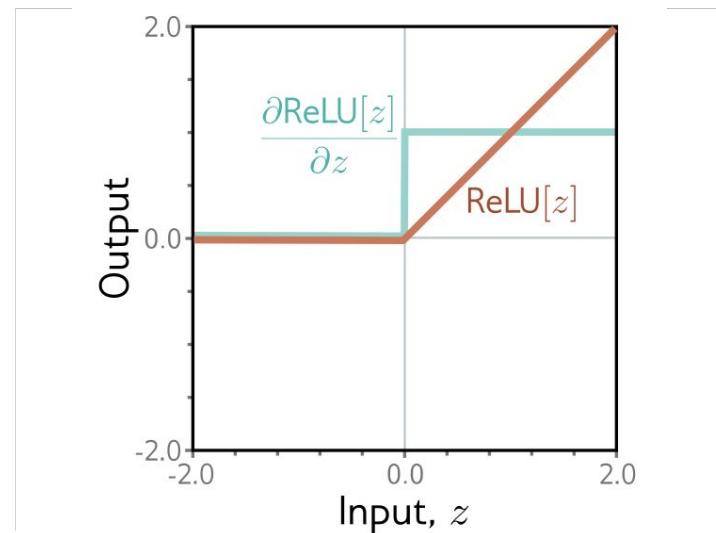
$$\frac{\partial \cos[z]}{\partial z} = -\sin[z] \quad \frac{\partial \exp[z]}{\partial z} = \exp[z] \quad \frac{\partial \sin[z]}{\partial z} = -\cos[z]$$



Warmup! Derivative of ReLU

$$a[z] = \text{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}.$$

Rectified Linear Unit



$$\mathbb{I}[z > 0]$$



Toy example

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x] \right] \right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

We want to compute:

$$\frac{\partial \ell_i}{\partial \beta_0}, \quad \frac{\partial \ell_i}{\partial \omega_0}, \quad \frac{\partial \ell_i}{\partial \beta_1}, \quad \frac{\partial \ell_i}{\partial \omega_1}, \quad \frac{\partial \ell_i}{\partial \beta_2}, \quad \frac{\partial \ell_i}{\partial \omega_2}, \quad \frac{\partial \ell_i}{\partial \beta_3}, \quad \text{and} \quad \frac{\partial \ell_i}{\partial \omega_3}$$

How does a small change in
 β_2 change the loss ℓ_i for the
i'th example?

Gradients of composed functions

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x] \right] \right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

Calculating expressions by hand:

- Some expressions are very complicated
- There are some redundancies

$$\begin{aligned} \frac{\partial \ell_i}{\partial \omega_0} = & -2 \left(\beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] \right] - y_i \right) \\ & \cdot \omega_1 \omega_2 \omega_3 \cdot x_i \cdot \cos [\beta_0 + \omega_0 \cdot x_i] \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] \\ & \cdot \sin \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin [\beta_0 + \omega_0 \cdot x_i] \right] \right] \end{aligned}$$

Forward pass

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin[\beta_0 + \omega_0 \cdot x] \right] \right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

1. Write it as a series of intermediate calculations

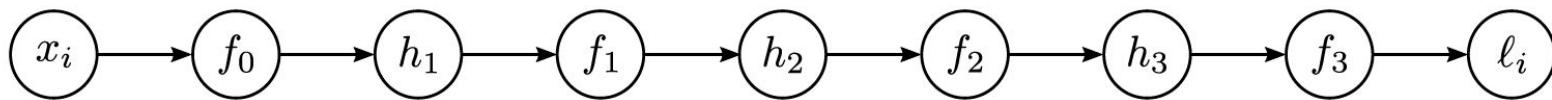
$$f_0 = \beta_0 + \omega_0 \cdot x_i \qquad \qquad f_2 = \beta_2 + \omega_2 \cdot h_2$$

$$h_1 = \sin[f_0] \qquad \qquad h_3 = \cos[f_2]$$

2. Compute these intermediate quantities

$$f_1 = \beta_1 + \omega_1 \cdot h_1 \qquad \qquad f_3 = \beta_3 + \omega_3 \cdot h_3$$

$$h_2 = \exp[f_1] \qquad \qquad \ell_i = (f_3 - y_i)^2.$$



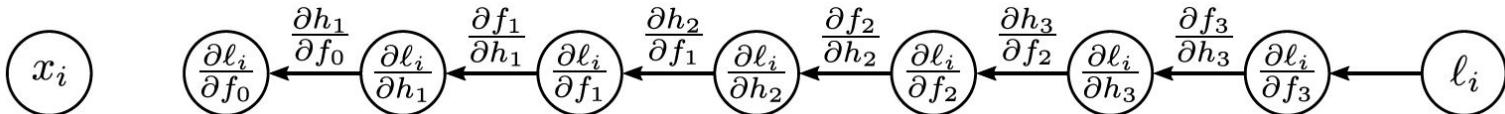
Backward pass

$$f[x, \phi] = \beta_3 + \omega_3 \cdot \cos \left[\beta_2 + \omega_2 \cdot \exp \left[\beta_1 + \omega_1 \cdot \sin[\beta_0 + \omega_0 \cdot x] \right] \right]$$

$$\ell_i = (f[x_i, \phi] - y_i)^2$$

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\frac{\partial \ell_i}{\partial f_3}, \quad \frac{\partial \ell_i}{\partial h_3}, \quad \frac{\partial \ell_i}{\partial f_2}, \quad \frac{\partial \ell_i}{\partial h_2}, \quad \frac{\partial \ell_i}{\partial f_1}, \quad \frac{\partial \ell_i}{\partial h_1}, \quad \text{and} \quad \frac{\partial \ell_i}{\partial f_0}$$



B

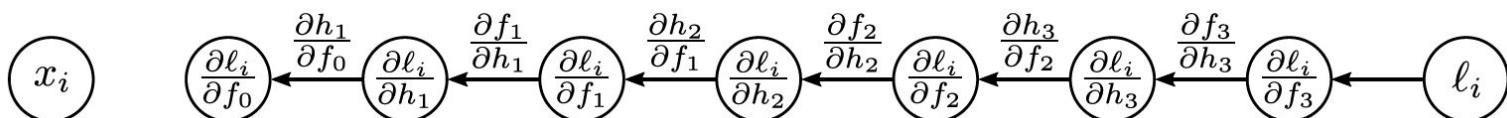
Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

The first are easy:

$$\frac{\partial \ell_i}{\partial f_3} = 2(f_3 - y_i)$$



Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

The rest are computed using the *chain rule*.

$$\frac{\partial \ell_i}{\partial h_3} = \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3}$$

How does a small change in \mathbf{h}_3 change ℓ_i ?

How does a small change in \mathbf{h}_3 change f_3 ?

How does a small change in f_3 change ℓ_i ?

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

The rest are computed using the *chain rule*.

$$\frac{\partial \ell_i}{\partial h_3} = \frac{\partial f_3}{\partial h_3} \frac{\partial \ell_i}{\partial f_3}$$

How does a small change in h_3 change ℓ_i ?

Already computed!

?

w_3

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

The rest are computed using the *chain rule*.

$$\frac{\partial \ell_i}{\partial f_2} = \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \right) \frac{\partial h_3}{\partial f_2}$$

?

Algebraically computed!

$-\sin(h_2)$

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$f_0 = \beta_0 + \omega_0 \cdot x_i$$

$$h_1 = \sin[f_0]$$

$$f_1 = \beta_1 + \omega_1 \cdot h_1$$

$$h_2 = \exp[f_1]$$

$$f_2 = \beta_2 + \omega_2 \cdot h_2$$

$$h_3 = \cos[f_2]$$

$$f_3 = \beta_3 + \omega_3 \cdot h_3$$

$$\ell_i = (f_3 - y_i)^2.$$

$$\frac{\partial \ell_i}{\partial f_2} = \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \right) \frac{\partial h_3}{\partial f_2}$$

The rest are computed using the *chain rule*.

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

The rest are computed using the *chain rule*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

$$\begin{aligned}\frac{\partial \ell_i}{\partial f_2} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \right) \frac{\partial h_3}{\partial f_2} \\ \frac{\partial \ell_i}{\partial h_2} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \frac{\partial h_3}{\partial f_2} \right) \frac{\partial f_2}{\partial h_2}\end{aligned}$$

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

The rest are computed using the *chain rule*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

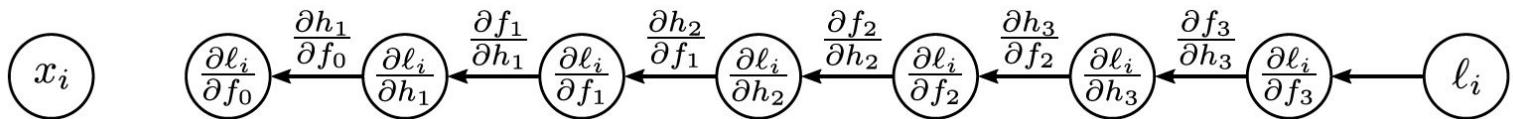
$$\begin{aligned} \frac{\partial \ell_i}{\partial f_2} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \right) \frac{\partial h_3}{\partial f_2} \\ \frac{\partial \ell_i}{\partial h_2} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \frac{\partial h_3}{\partial f_2} \right) \frac{\partial f_2}{\partial h_2} \\ \frac{\partial \ell_i}{\partial f_1} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \frac{\partial h_3}{\partial f_2} \frac{\partial f_2}{\partial h_2} \right) \frac{\partial h_2}{\partial f_1} \\ \frac{\partial \ell_i}{\partial h_1} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \frac{\partial h_3}{\partial f_2} \frac{\partial f_2}{\partial h_2} \frac{\partial h_2}{\partial f_1} \right) \frac{\partial f_1}{\partial h_1} \\ \frac{\partial \ell_i}{\partial f_0} &= \left(\frac{\partial \ell_i}{\partial f_3} \frac{\partial f_3}{\partial h_3} \frac{\partial h_3}{\partial f_2} \frac{\partial f_2}{\partial h_2} \frac{\partial h_2}{\partial f_1} \frac{\partial f_1}{\partial h_1} \right) \frac{\partial h_1}{\partial f_0}. \end{aligned}$$

Backward pass

1. Compute the derivatives of the loss with respect to these intermediate quantities, in *reverse order*.

$$\begin{array}{ll} f_0 = \beta_0 + \omega_0 \cdot x_i & f_2 = \beta_2 + \omega_2 \cdot h_2 \\ h_1 = \sin[f_0] & h_3 = \cos[f_2] \\ f_1 = \beta_1 + \omega_1 \cdot h_1 & f_3 = \beta_3 + \omega_3 \cdot h_3 \\ h_2 = \exp[f_1] & \ell_i = (f_3 - y_i)^2. \end{array}$$

Chain rule all the way down!

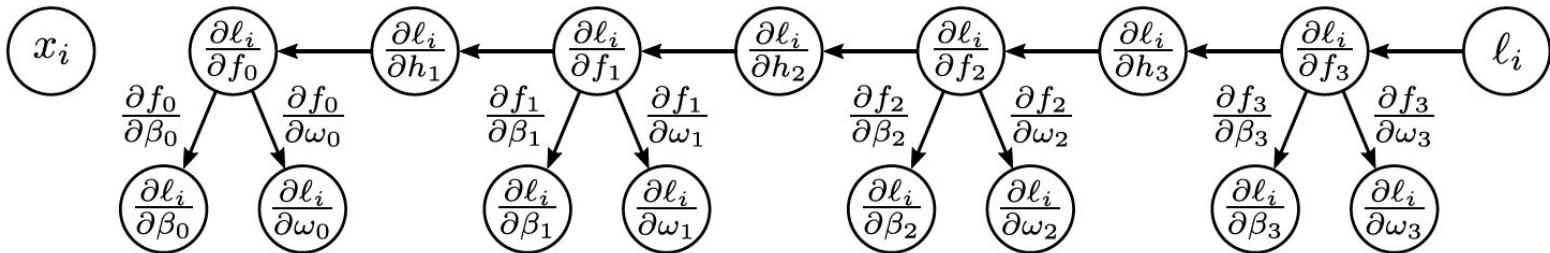


Backward pass

2. Find how the loss changes as a function of the parameters β and ω .

$$\begin{aligned}f_0 &= \beta_0 + \omega_0 \cdot x_i & f_2 &= \beta_2 + \omega_2 \cdot h_2 \\h_1 &= \sin[f_0] & h_3 &= \cos[f_2] \\f_1 &= \beta_1 + \omega_1 \cdot h_1 & f_3 &= \beta_3 + \omega_3 \cdot h_3 \\h_2 &= \exp[f_1] & \ell_i &= (f_3 - y_i)^2.\end{aligned}$$

Chain rule all the way down! Same recipe for weight and bias terms too!



Backward pass

2. Find how the loss changes as a function of the parameters β and ω .

$$f_0 = \beta_0 + \omega_0 \cdot x_i$$

$$h_1 = \sin[f_0]$$

$$f_1 = \beta_1 + \omega_1 \cdot h_1$$

$$h_2 = \exp[f_1]$$

$$f_2 = \beta_2 + \omega_2 \cdot h_2$$

$$h_3 = \cos[f_2]$$

$$f_3 = \beta_3 + \omega_3 \cdot h_3$$

$$\ell_i = (f_3 - y_i)^2.$$

Chain rule all the way down! Same recipe for weight and bias terms too!

$$\frac{\partial \ell_i}{\partial \omega_k} = \frac{\partial f_k}{\partial \omega_k} \frac{\partial \ell_i}{\partial f_k}$$

Reading suggestion

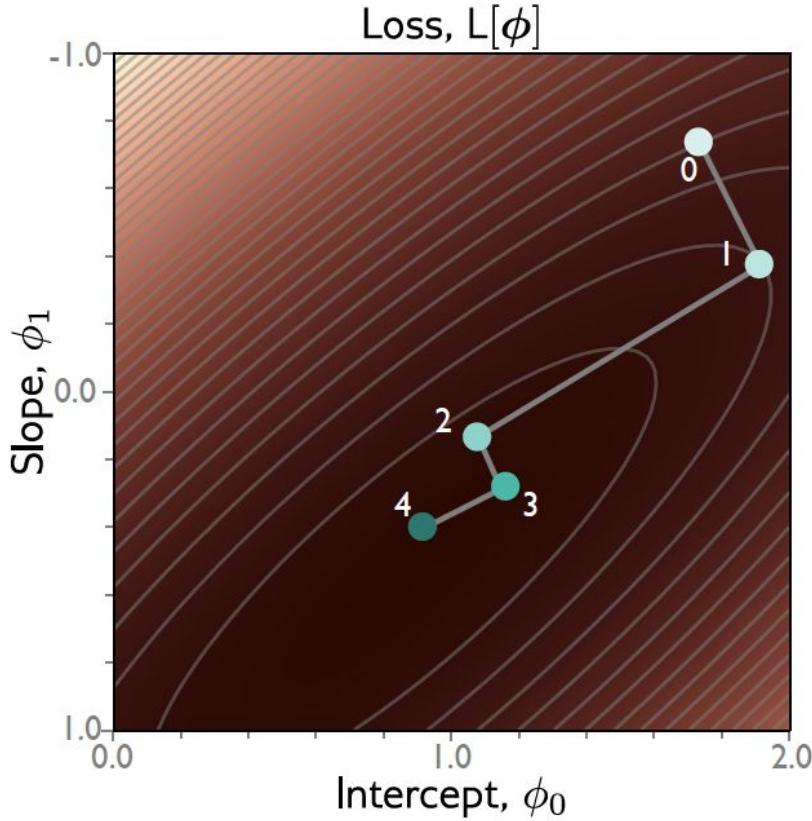
Check 7.4 in your textbook for a discussion on the extension to matrix calculus.

Automatic differentiation

- Modern deep learning frameworks compute derivatives **automatically**
- You just have to specify the model and the loss
- How?
 - Each component knows how to compute its own derivative
 - ReLU knows how to compute deriv of output w.r.t. input
 - Linear function knows how to compute deriv of output w.r.t. input
 - Linear function knows how to compute deriv of output w.r.t. parameter
 - You specify the order of the components
 - It can compute the chain of derivatives
- Works with branches as long as it's still an acyclic graph

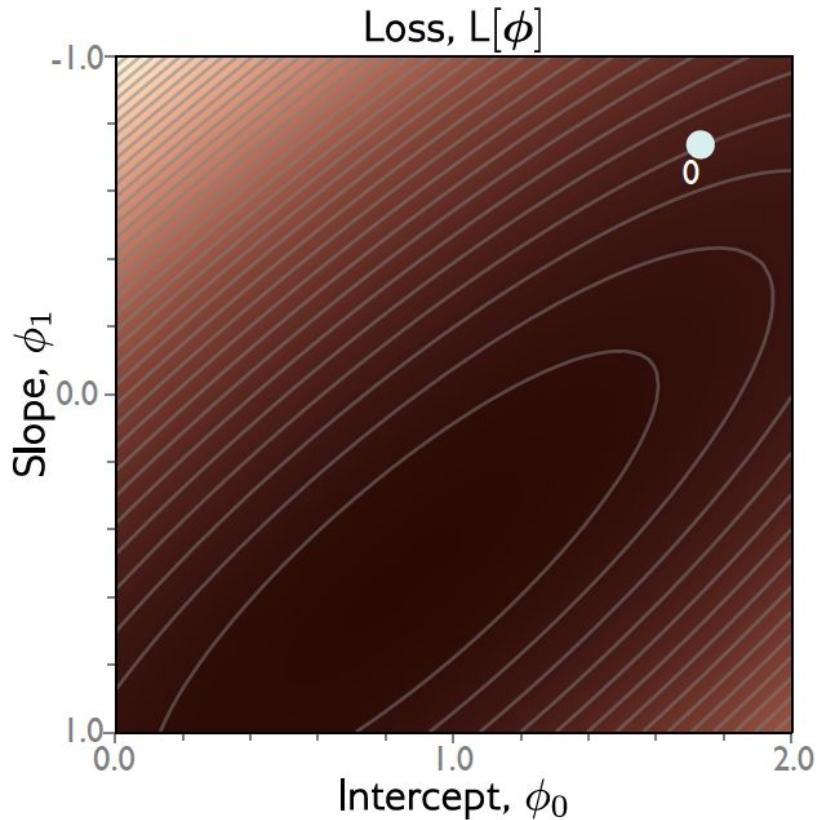
Optimisation

Training a simple model



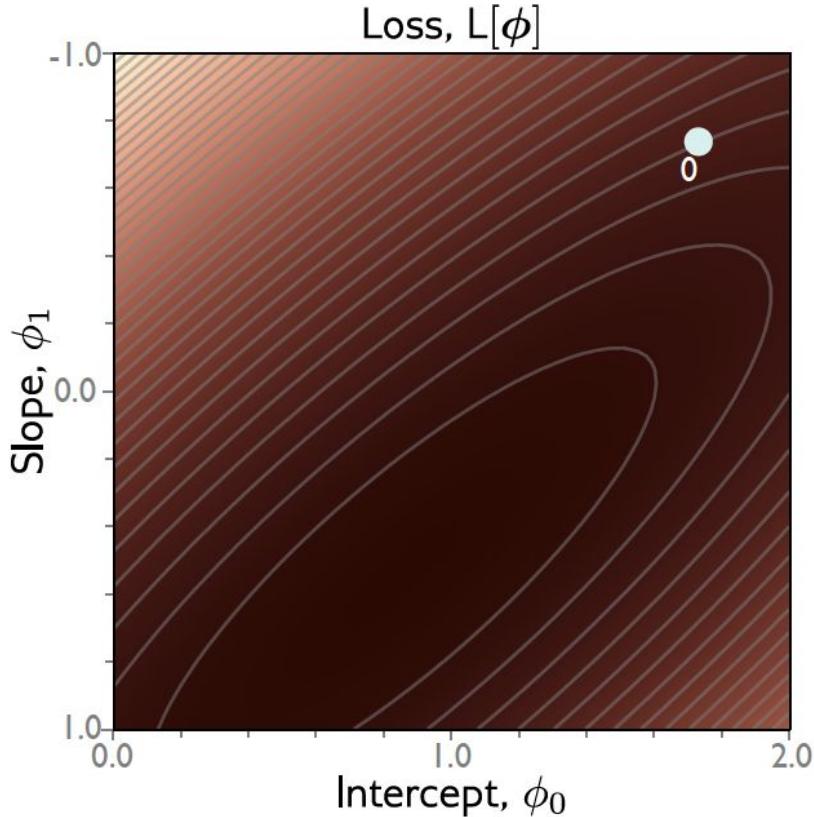
1. Define a loss function
2. Compute the change in parameters required to make the loss smaller
3. **Apply the change and get new parameters**
4. Repeat from (2)

Gradient descent



$$\begin{aligned} L[\phi] &= \sum_{i=1}^I \ell_i = \sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

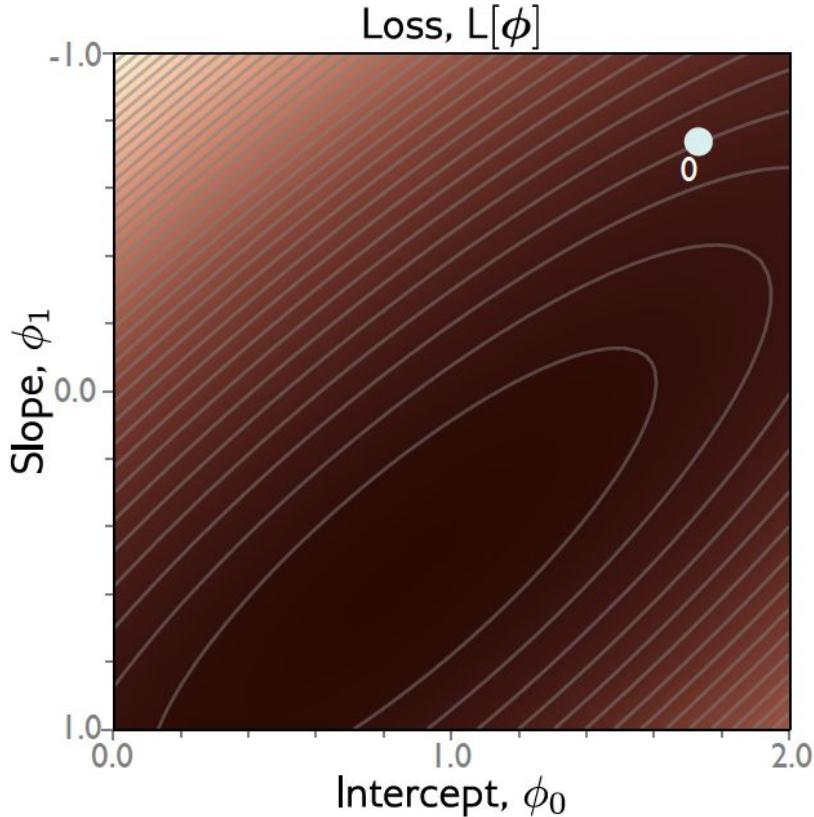
Gradient descent



$$\begin{aligned} L[\phi] &= \sum_{i=1}^I \ell_i = \sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

Gradient descent

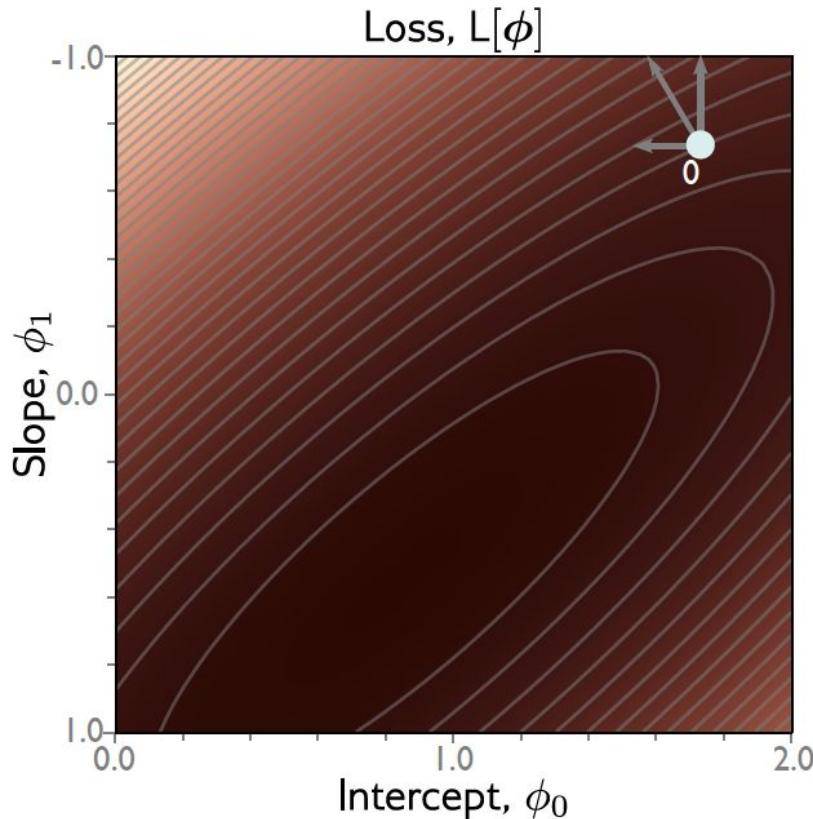


$$\begin{aligned} L[\phi] &= \sum_{i=1}^I \ell_i = \sum_{i=1}^I (\mathbf{f}[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

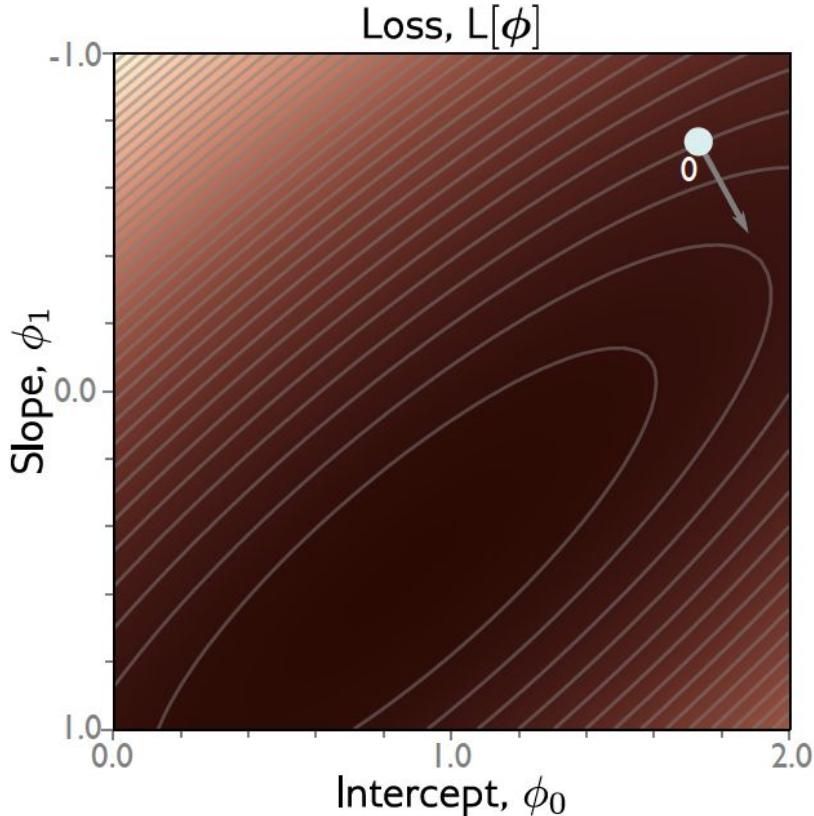
Gradient descent



$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

Gradient descent



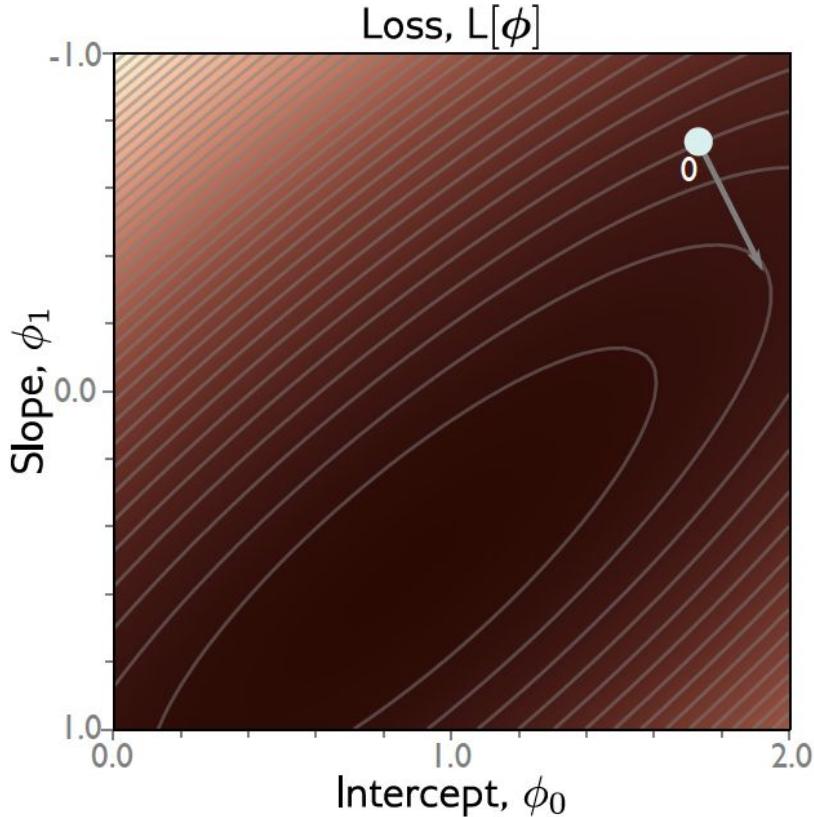
$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

$$\phi \leftarrow \phi - \alpha \frac{\partial L}{\partial \phi}$$

α = step size or learning rate if fixed

Gradient descent



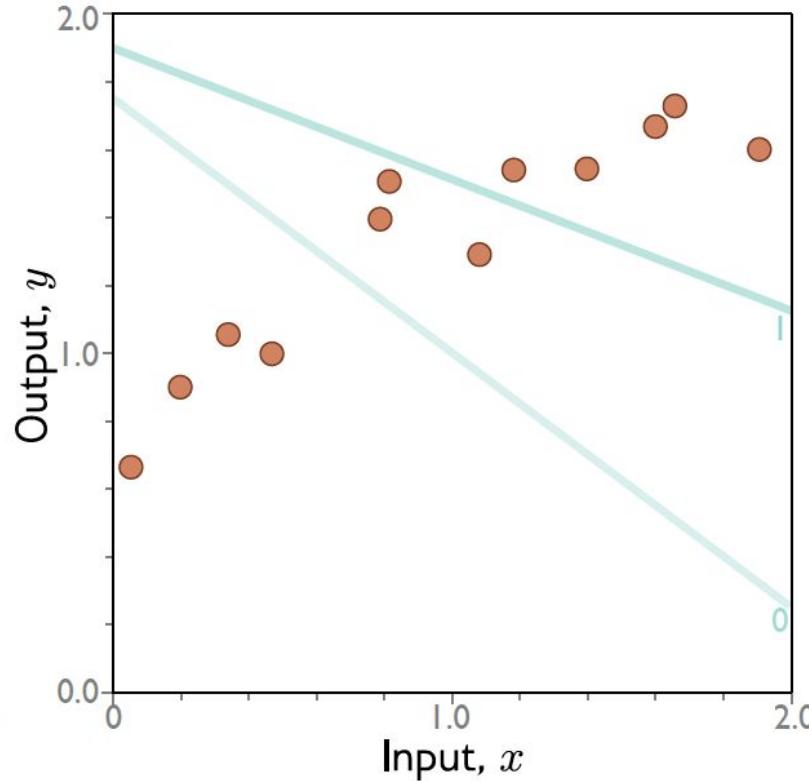
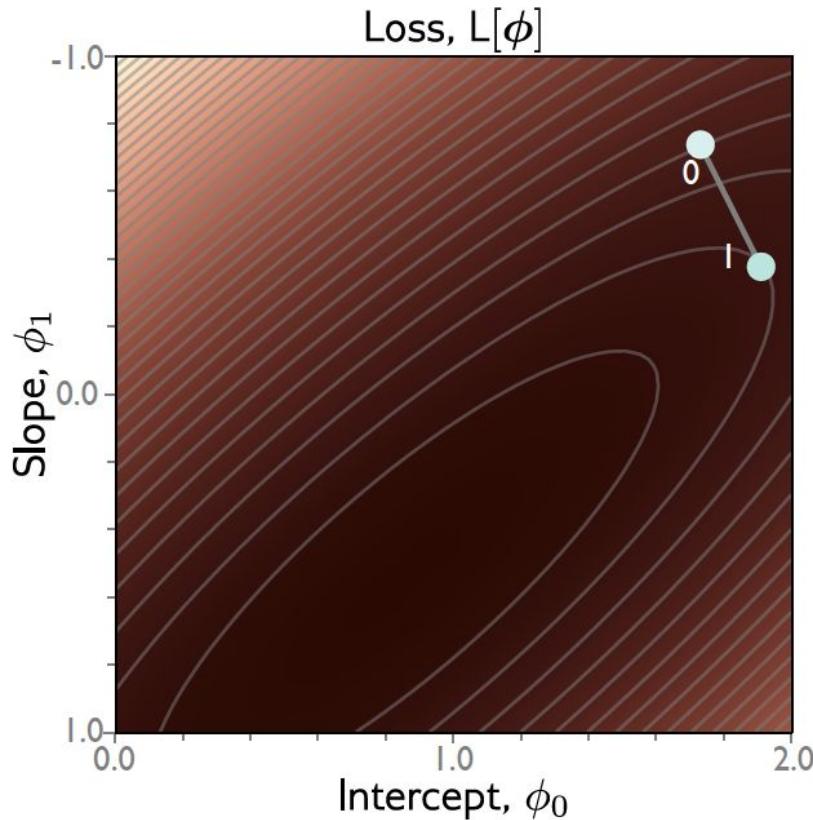
$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

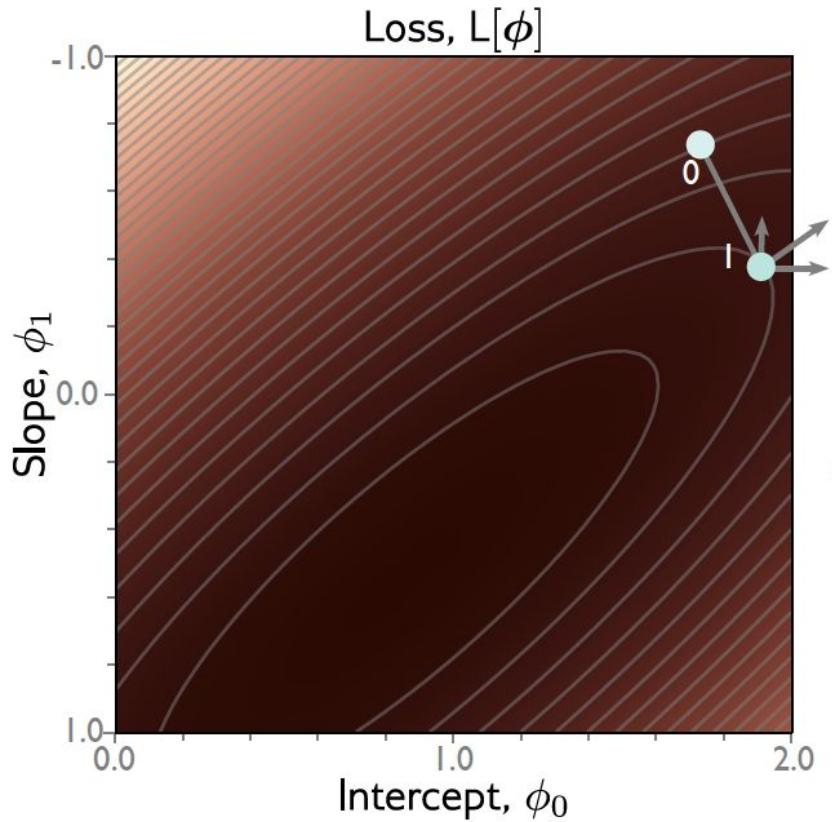
$$\phi \leftarrow \phi - \alpha \frac{\partial L}{\partial \phi}$$

α = step size

Gradient descent



Gradient descent



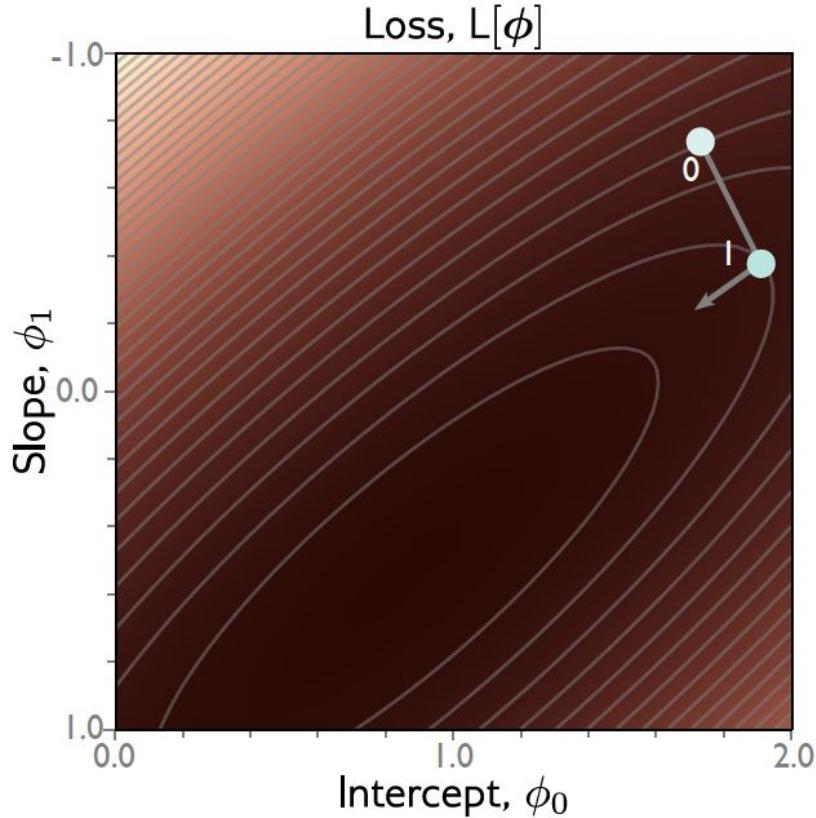
$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

$$\phi \leftarrow \phi - \alpha \frac{\partial L}{\partial \phi}$$

α = step size

Gradient descent



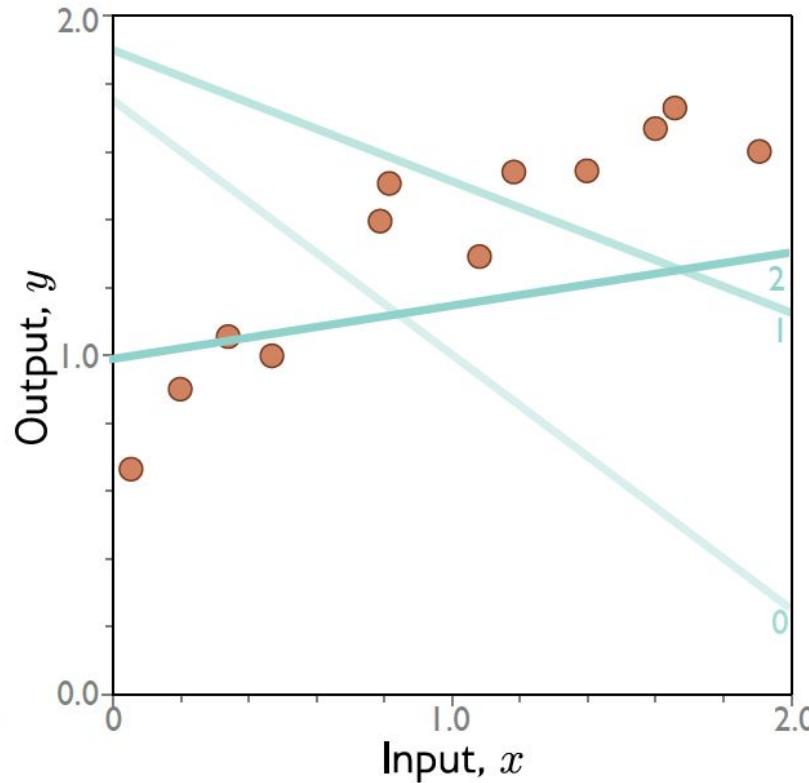
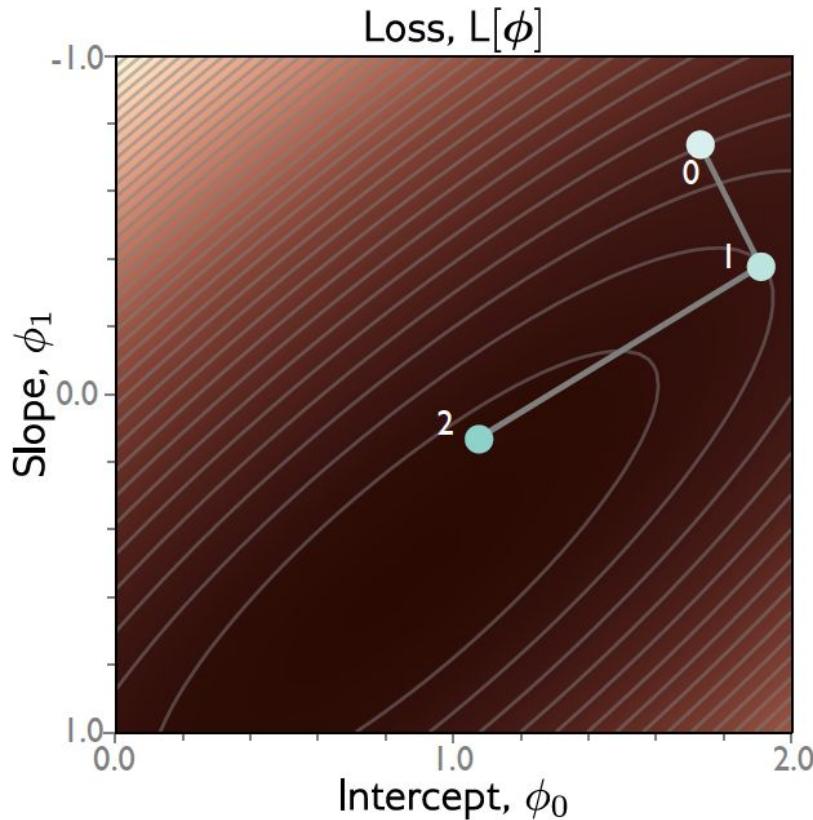
$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

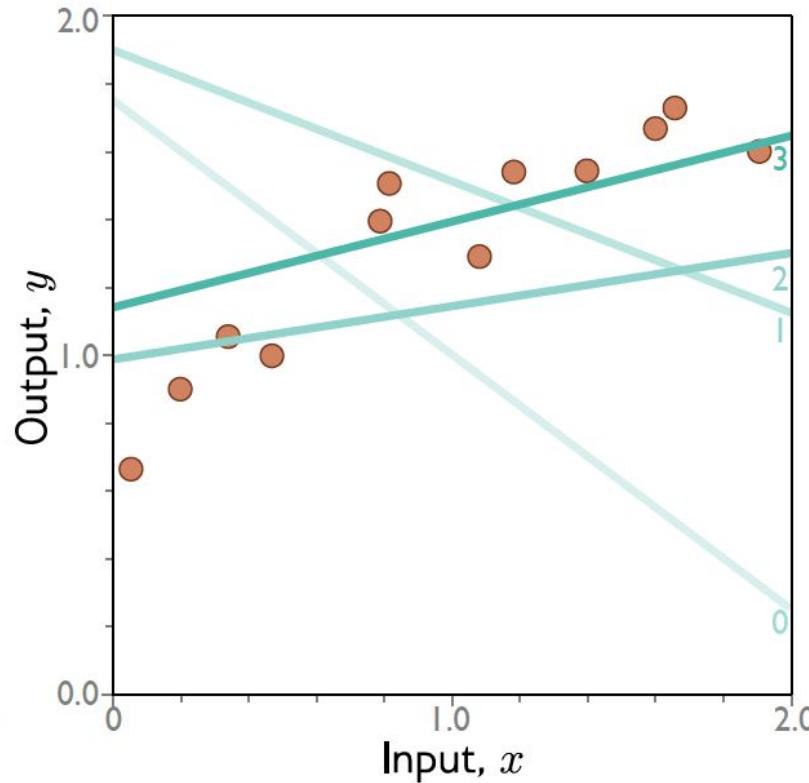
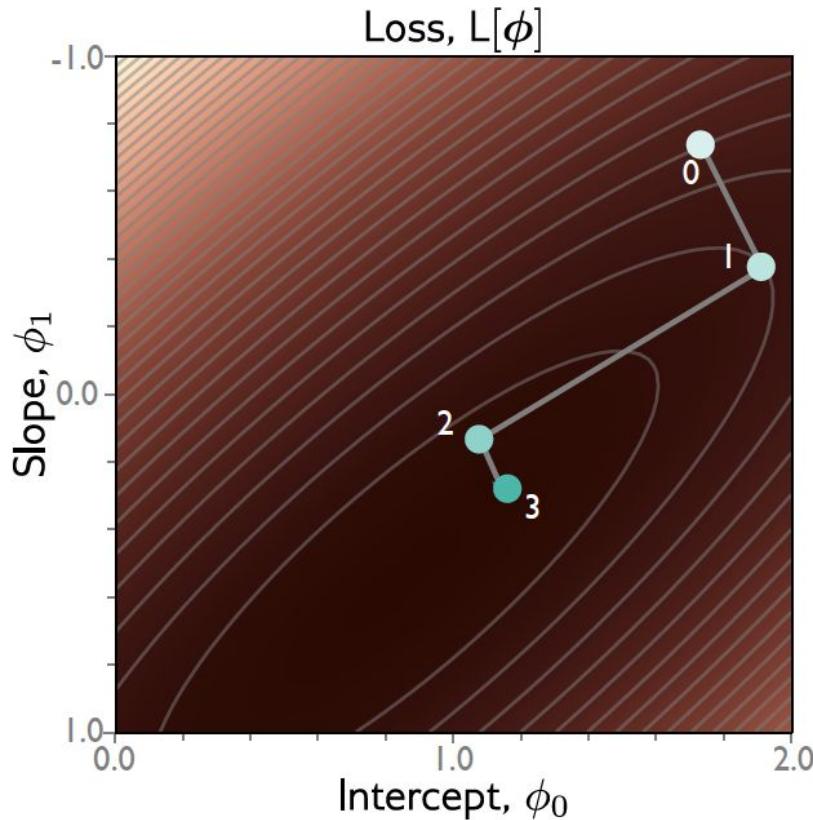
$$\phi \leftarrow \phi - \alpha \frac{\partial L}{\partial \phi}$$

α = step size

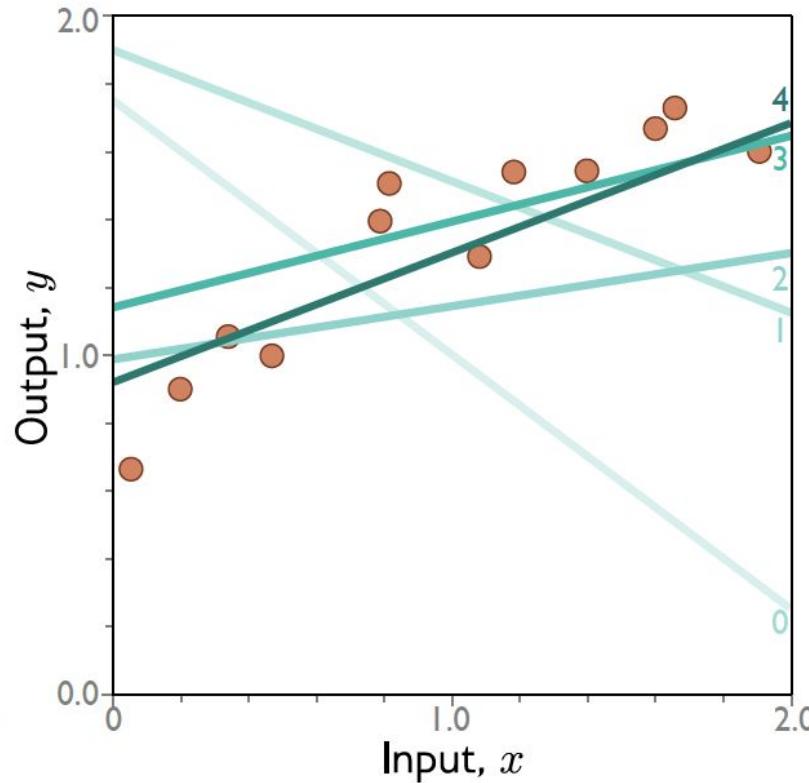
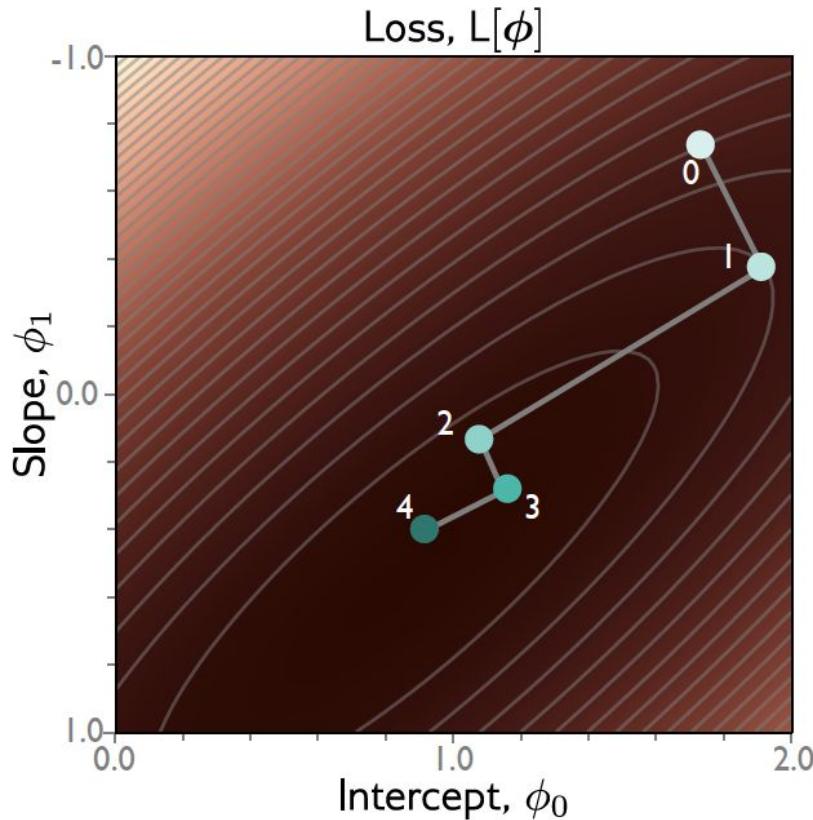
Gradient descent



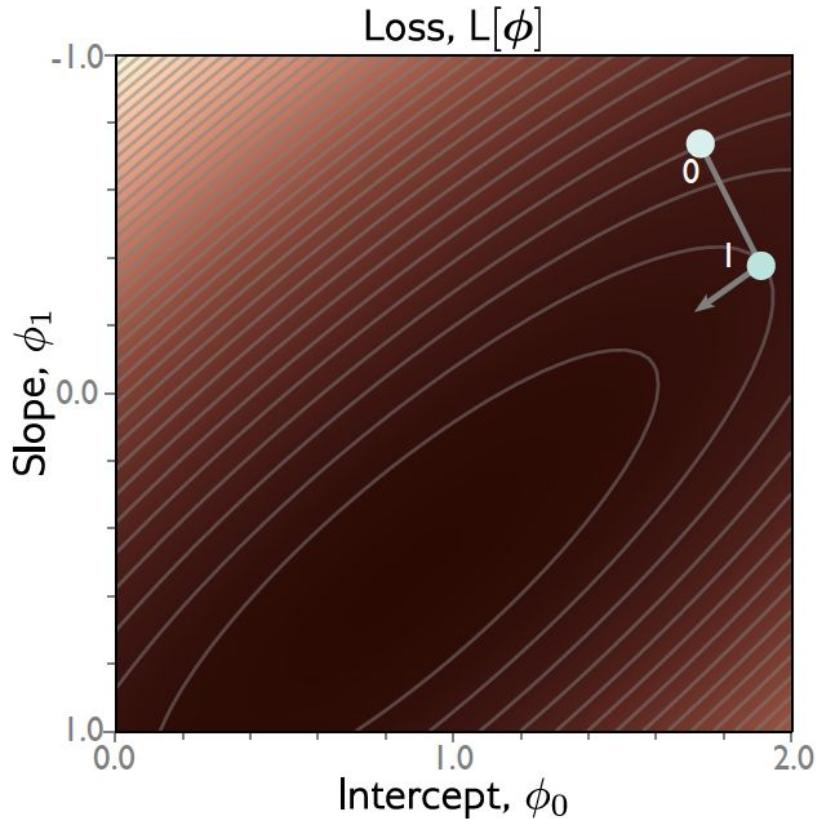
Gradient descent



Gradient descent



Line search



$$\frac{\partial L}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^I \ell_i = \sum_{i=1}^I \frac{\partial \ell_i}{\partial \phi}$$

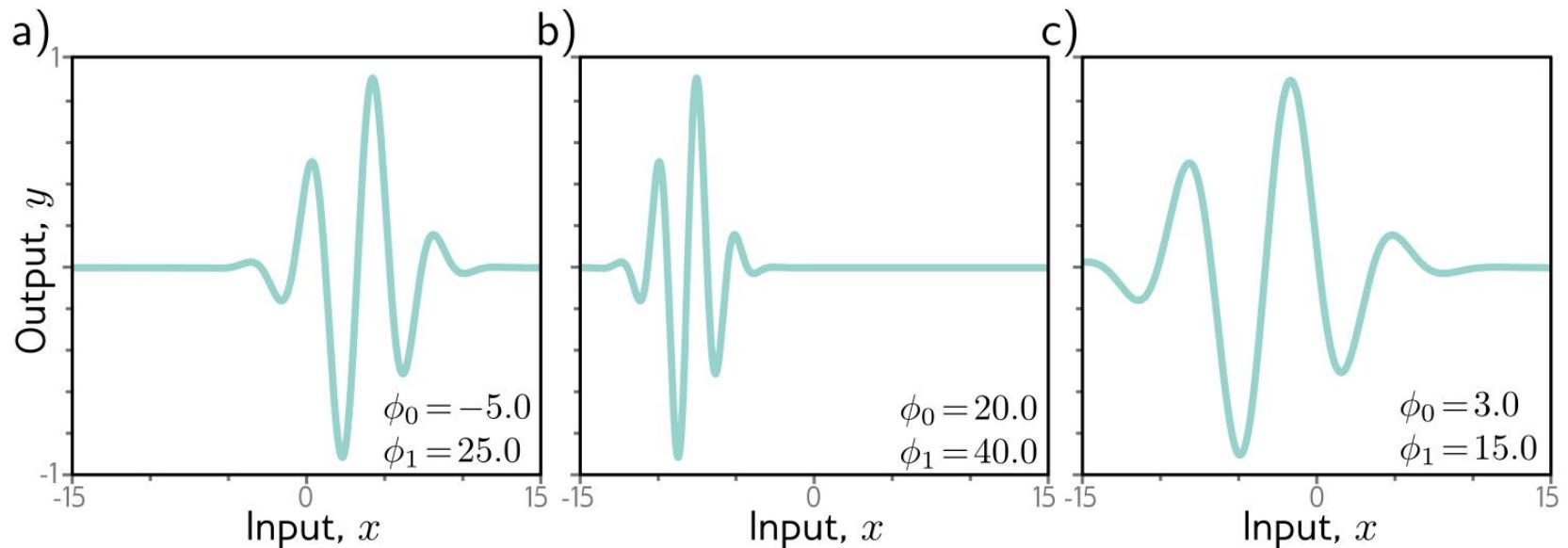
$$\frac{\partial \ell_i}{\partial \phi} = \begin{bmatrix} \frac{\partial \ell_i}{\partial \phi_0} \\ \frac{\partial \ell_i}{\partial \phi_1} \end{bmatrix} = \begin{bmatrix} 2(\phi_0 + \phi_1 x_i - y_i) \\ 2x_i(\phi_0 + \phi_1 x_i - y_i) \end{bmatrix}$$

$$\phi \leftarrow \phi - \alpha \frac{\partial L}{\partial \phi}$$

α = step size

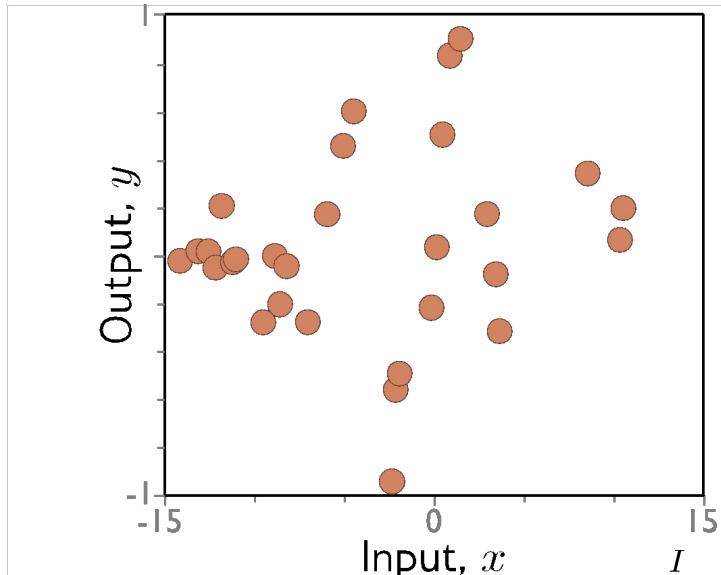
Non-convex case. Gabor model

$$f[x, \phi] = \sin[\phi_0 + 0.06 \cdot \phi_1 x] \cdot \exp\left(-\frac{(\phi_0 + 0.06 \cdot \phi_1 x)^2}{8.0}\right)$$

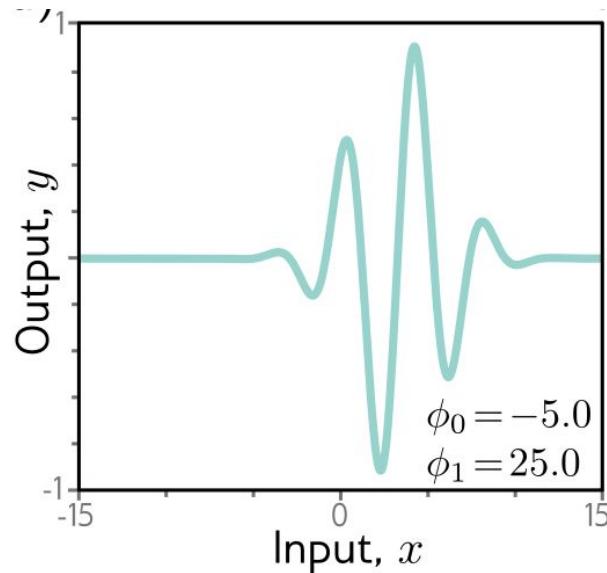


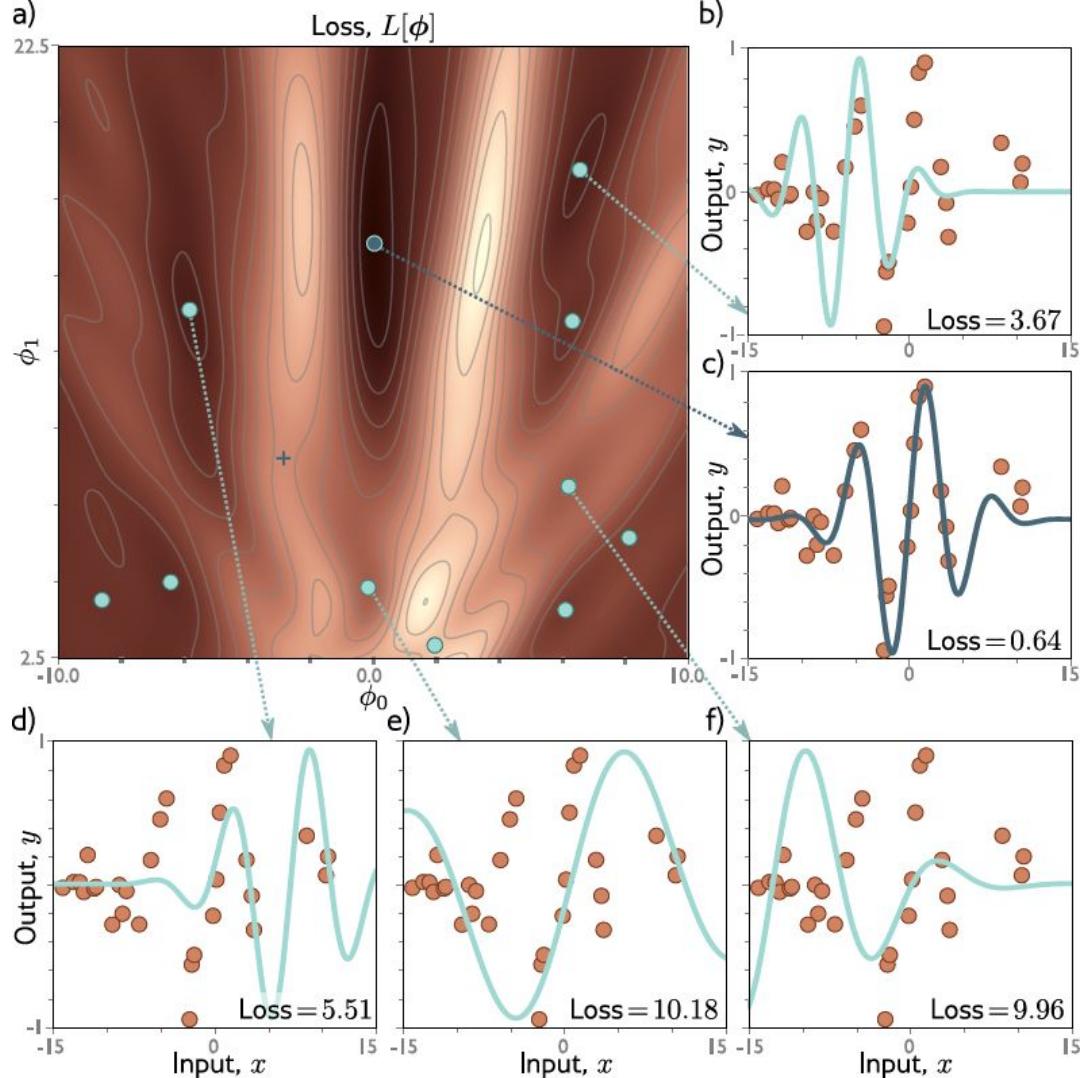
Gabor model

$$f[x, \phi] = \sin[\phi_0 + 0.06 \cdot \phi_1 x] \cdot \exp\left(-\frac{(\phi_0 + 0.06 \cdot \phi_1 x)^2}{8.0}\right)$$

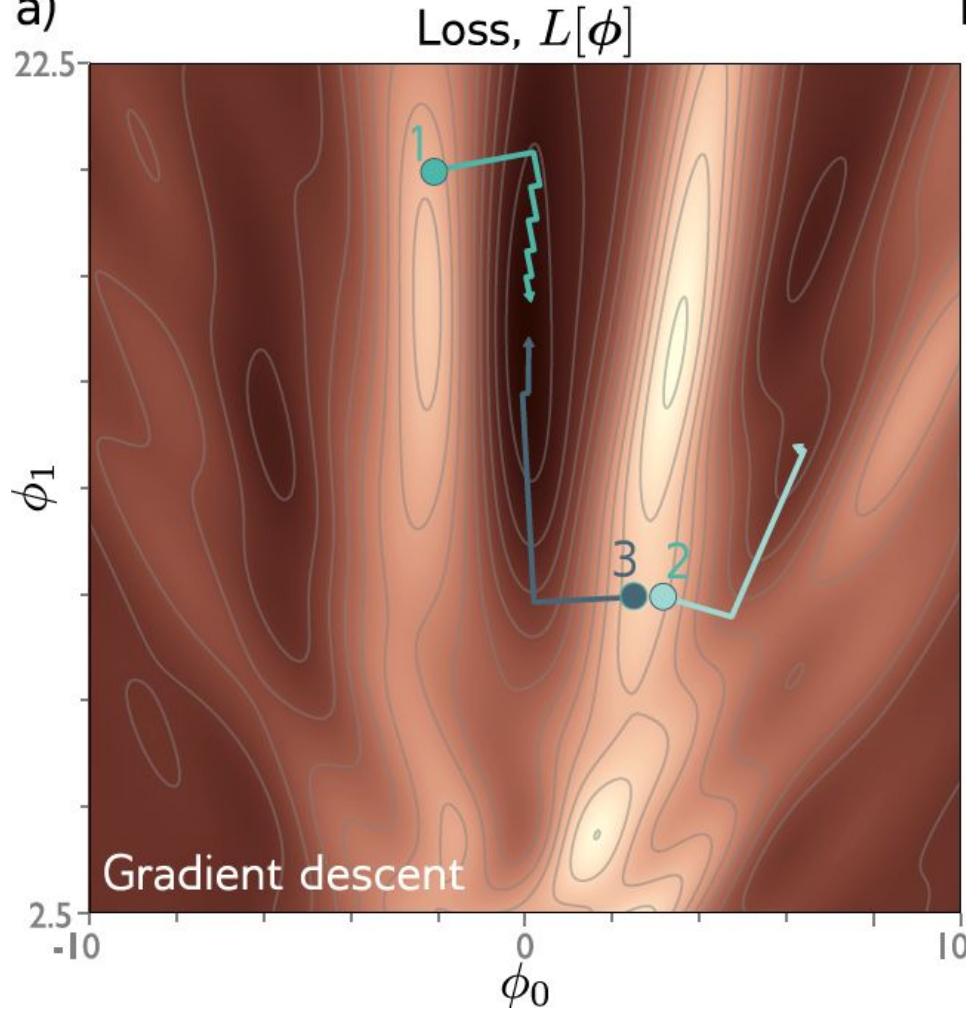


$$L[\phi] = \sum_{i=1}^I (f[x_i, \phi] - y_i)^2$$



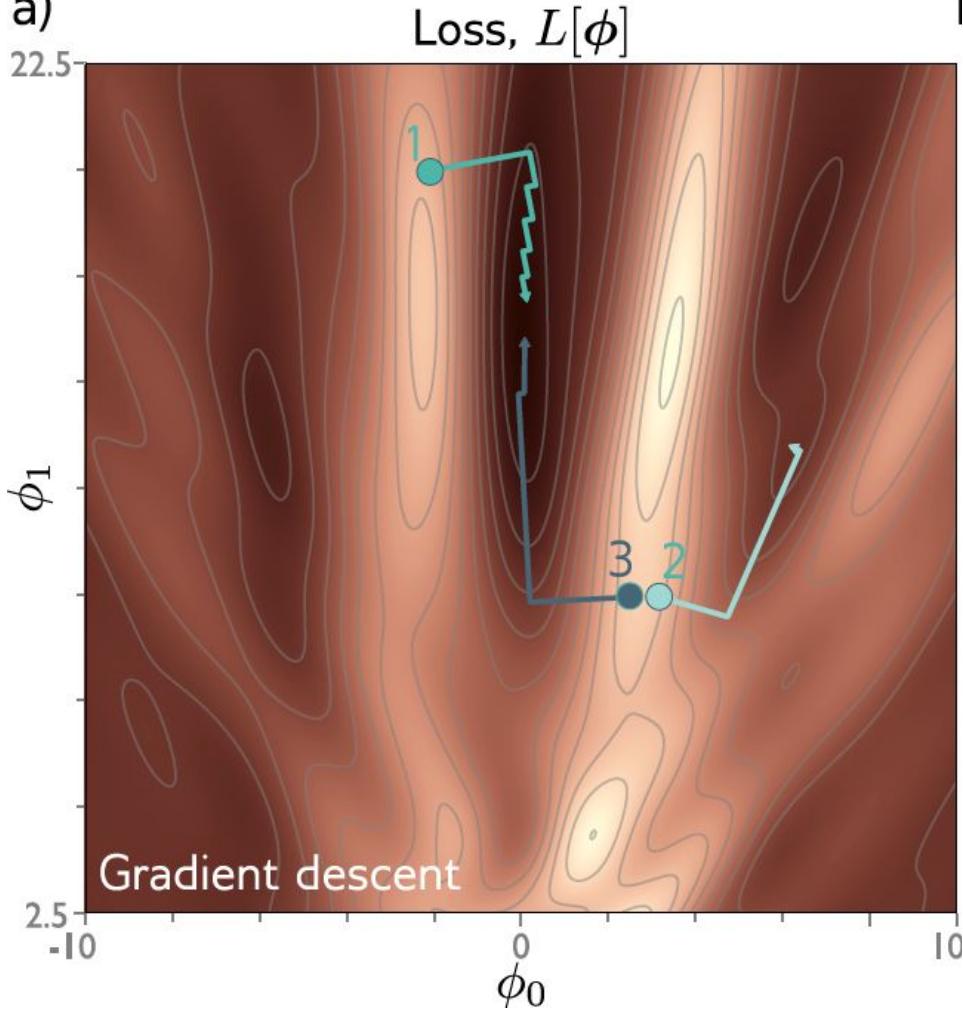


a)



- Gradient descent gets to the global minimum if we start in the right “valley”
- Otherwise, descent to a local minimum
- Or get stuck near a saddle point

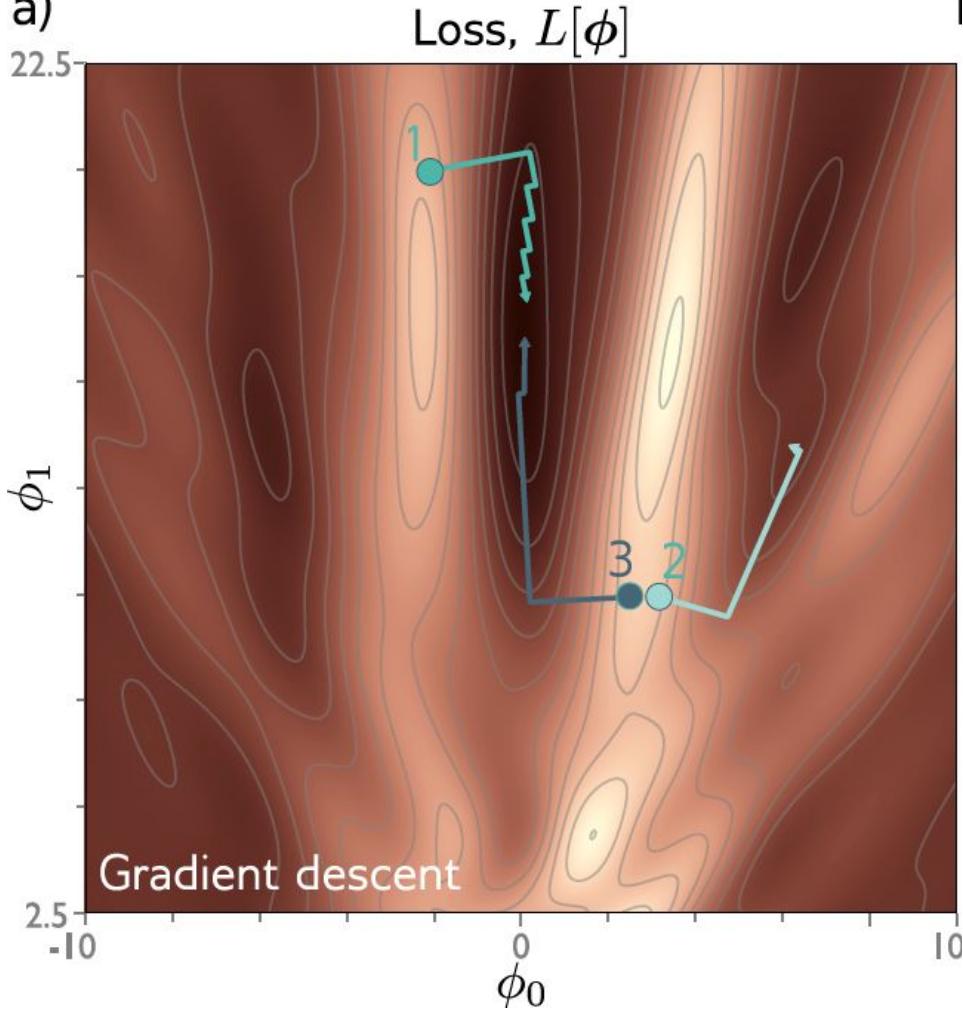
a)



Solution: add noise!

- Stochastic gradient descent
- Compute gradient based on only a subset of points – a **mini-batch**
- Work through dataset sampling without replacement
- One pass though the data is called an **epoch**

a)



Stochastic gradient descent

Before (full batch descent)

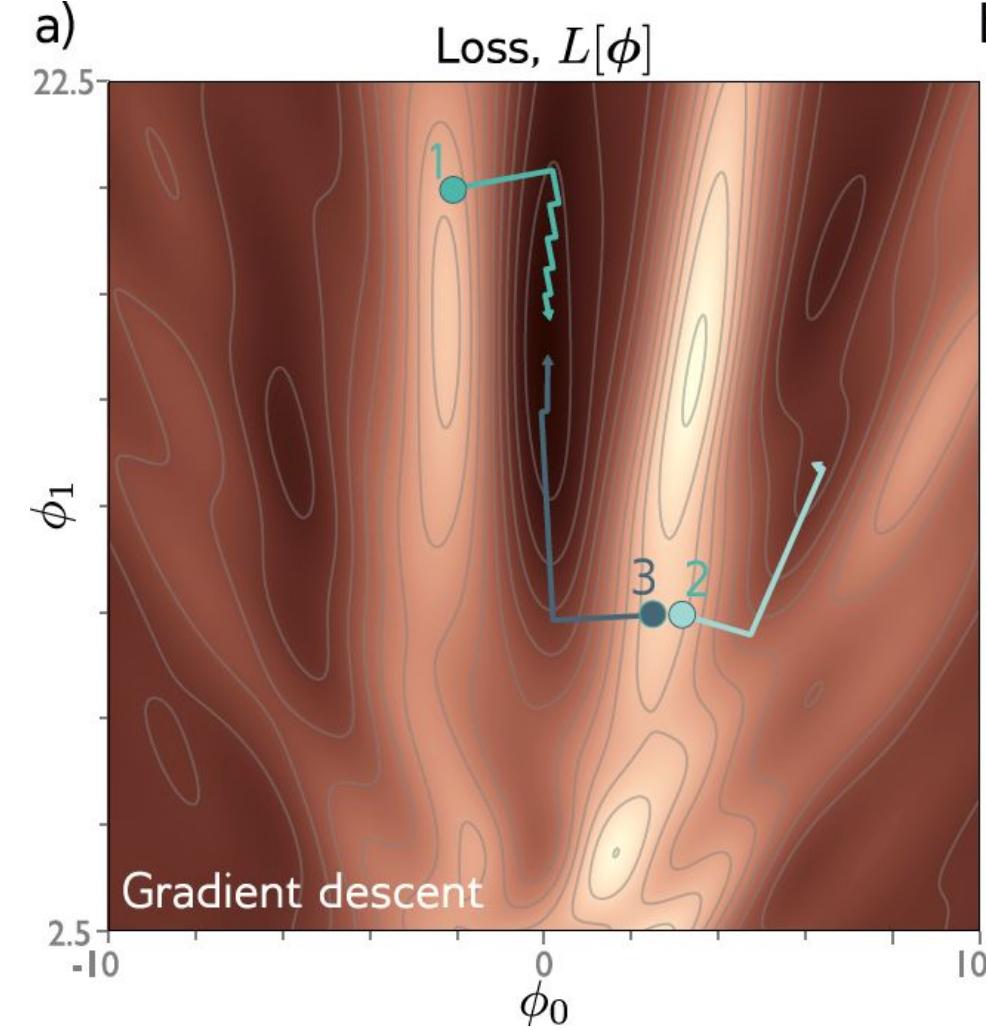
$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i=1}^I \frac{\partial \ell_i[\phi_t]}{\partial \phi},$$

After (SGD)

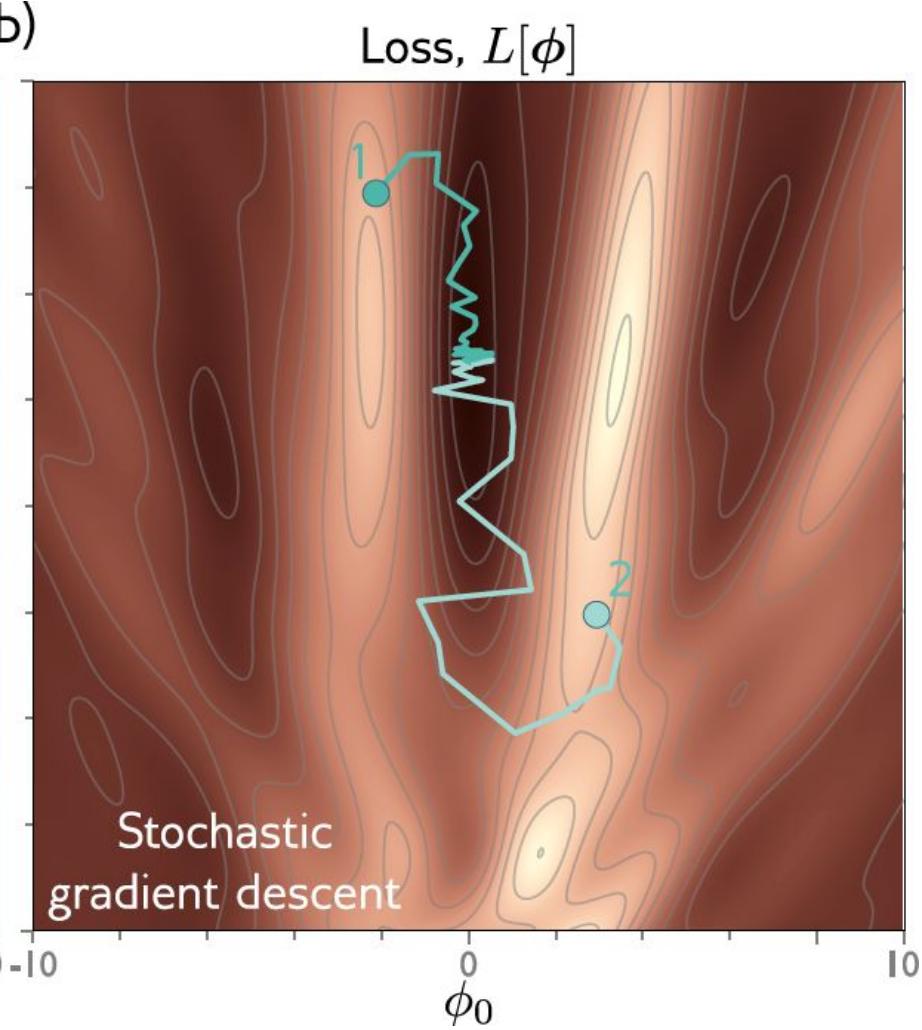
$$\phi_{t+1} \leftarrow \phi_t - \alpha \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi},$$

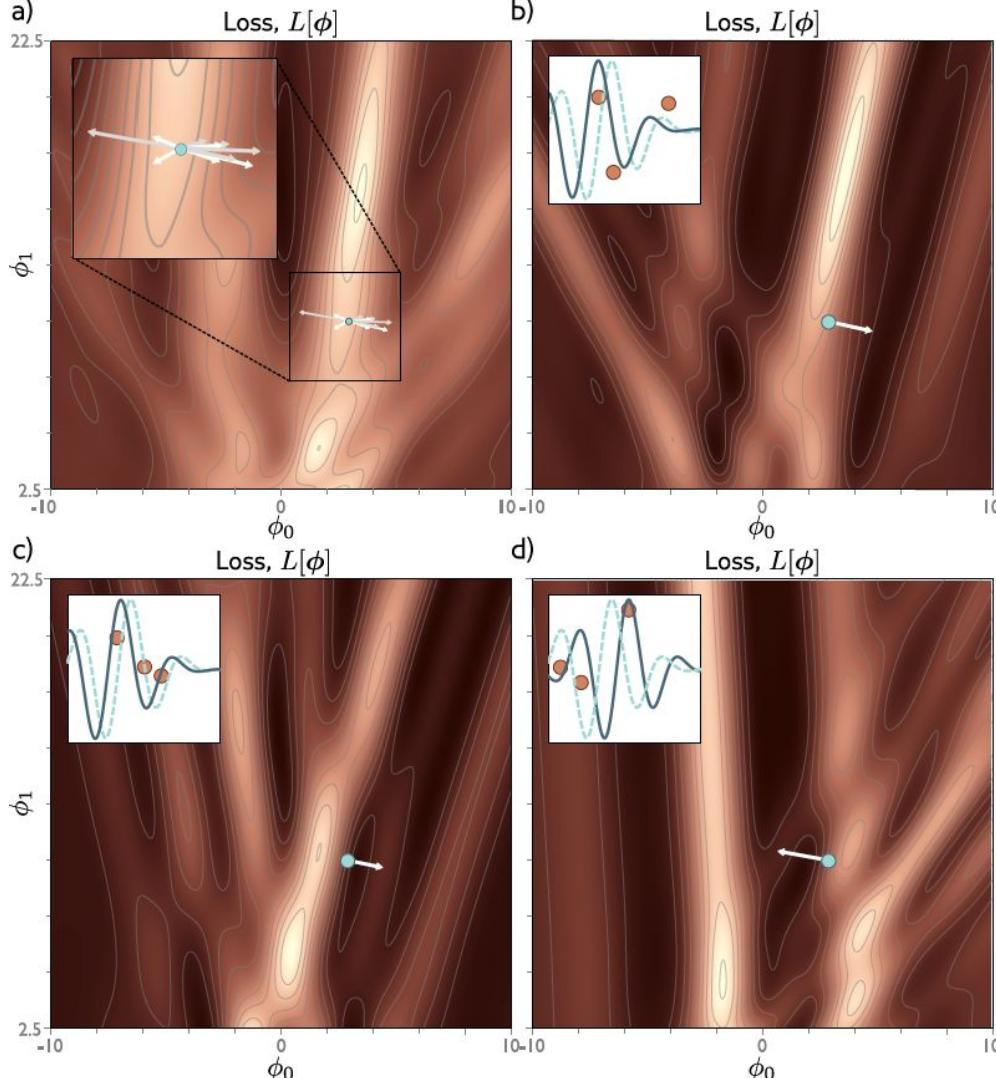
Fixed learning rate α

a)



b)

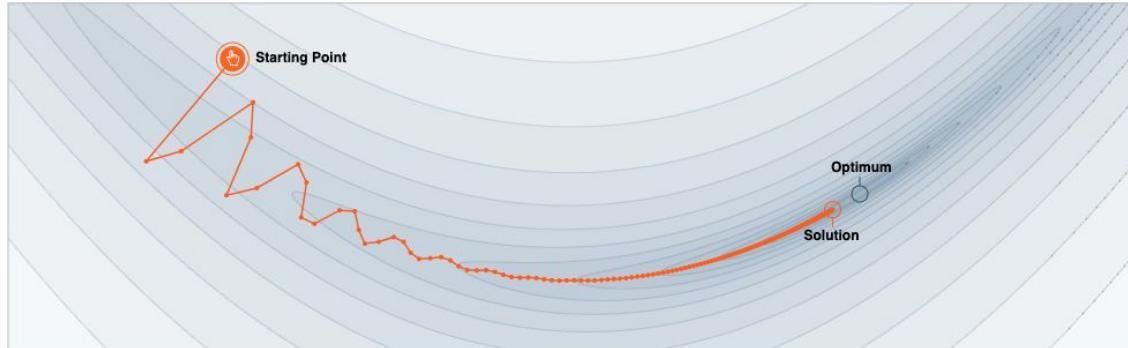




Properties of SGD

- Can escape from local minima
 - Adds noise, but still sensible updates as based on part of data
 - Uses all data equally
 - Less computationally expensive
 - Seems to find better solutions
-
- Doesn't converge in traditional sense
 - Learning rate schedule – decrease learning rate over time

Momentum



Step-size $\alpha = 0.02$

0 0.003 0.006

Momentum $\beta = 0.99$

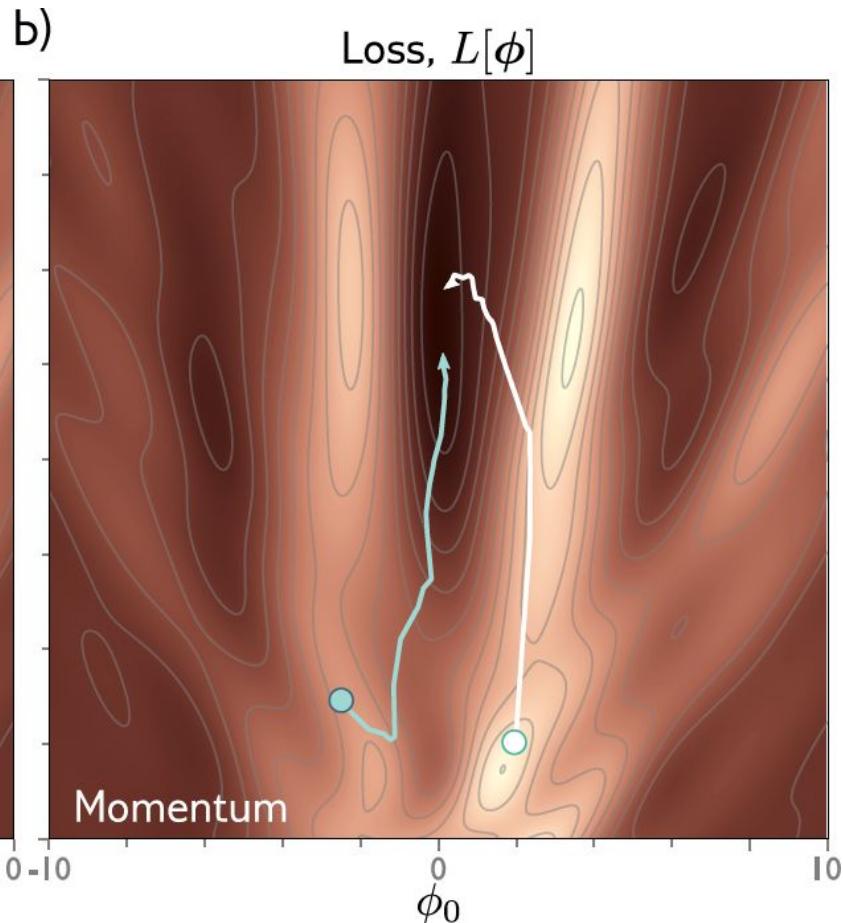
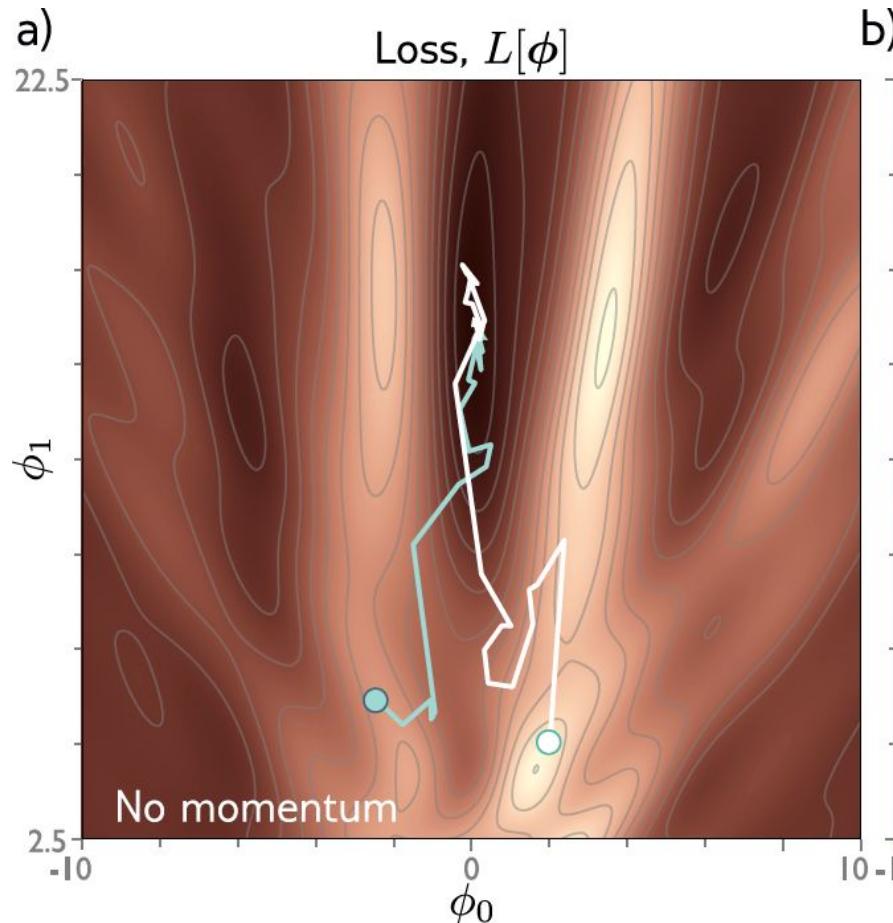
0.00 0.500 0.990

We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

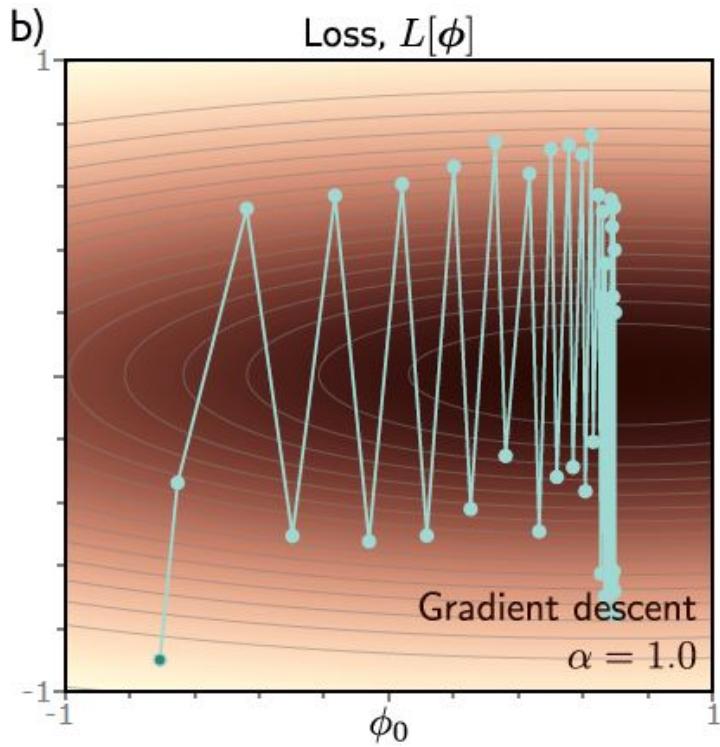
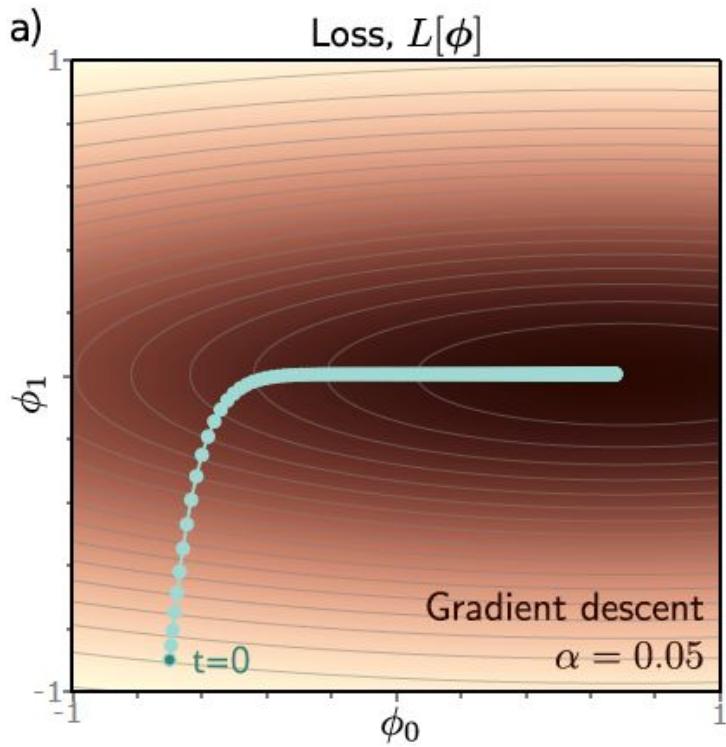
Momentum

Weighted sum of this gradient and previous gradient

$$\begin{aligned}\mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi} \\ \phi_{t+1} &\leftarrow \phi_t - \alpha \cdot \mathbf{m}_{t+1}\end{aligned}$$



Adaptive moment estimation. Adam



Normalized gradients

Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$

Normalized gradients

Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$

$$\mathbf{m}_{t+1} = \begin{bmatrix} 3.0 \\ -2.0 \\ 5.0 \end{bmatrix}$$

$$\mathbf{v}_{t+1} = \begin{bmatrix} 9.0 \\ 4.0 \\ 25.0 \end{bmatrix}$$

$$\frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon} = \begin{bmatrix} 1.0 \\ -1.0 \\ 1.0 \end{bmatrix}$$

Normalized gradients

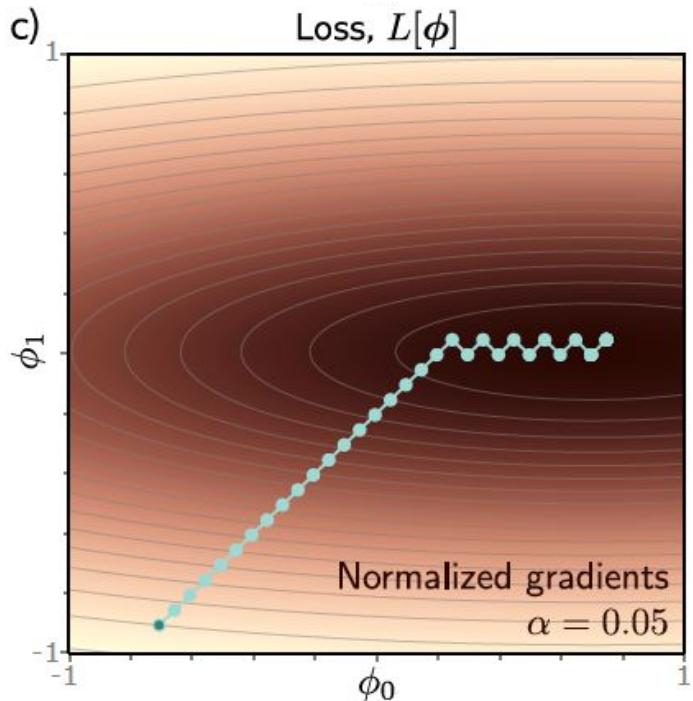
Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$



Normalized gradients

Compute mean and pointwise squared gradients with momentum

Moderate near start of the sequence
to prevent small estimates

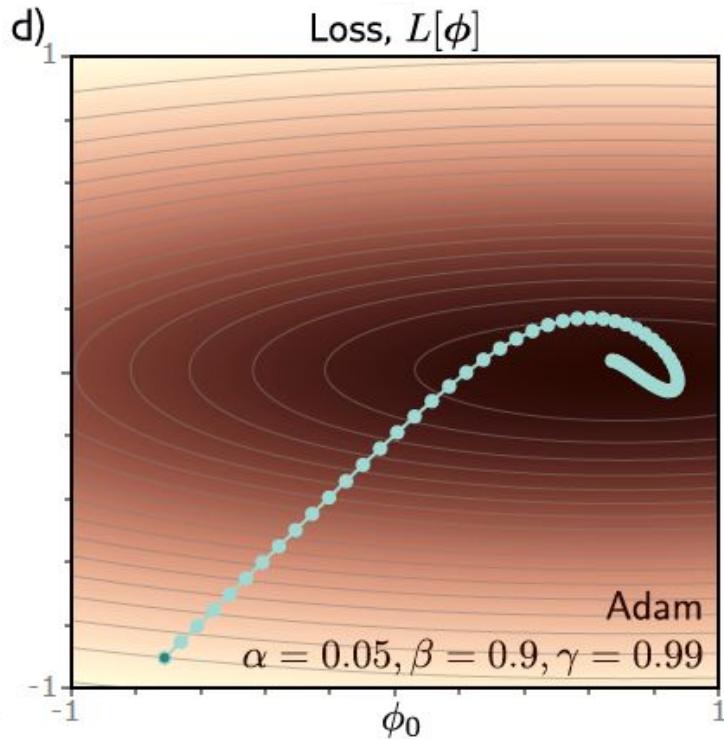
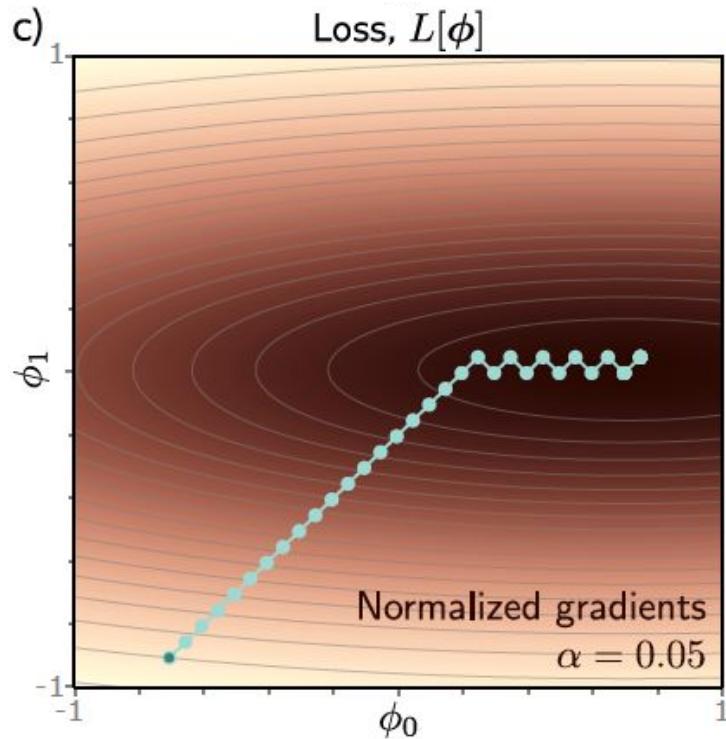
Update the parameters

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \frac{\partial L[\phi_t]}{\partial \phi}$$
$$\mathbf{v}_{t+1} \leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2$$

$$\tilde{\mathbf{m}}_{t+1} \leftarrow \frac{\mathbf{m}_{t+1}}{1 - \beta^{t+1}}$$
$$\tilde{\mathbf{v}}_{t+1} \leftarrow \frac{\mathbf{v}_{t+1}}{1 - \gamma^{t+1}}$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\tilde{\mathbf{m}}_{t+1}}{\sqrt{\tilde{\mathbf{v}}_{t+1}} + \epsilon}$$

Adaptive moment estimation. Adam



RMSprop (precursor to Adam)

Compute mean and pointwise squared gradients with momentum

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \frac{\partial L[\phi_t]}{\partial \phi}$$
$$\mathbf{v}_{t+1} \leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2$$

Update the parameters

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\tilde{\mathbf{m}}_{t+1}}{\sqrt{\tilde{\mathbf{v}}_{t+1}} + \epsilon} \frac{\partial L[\phi_t]}{\partial \phi}$$

How to reach me:

florin.gogianu@gmail.com

Please send unstructured feedback, since this is a new version of the lecture!

Post-lecture reading suggestions

- Maximum Likelihood Estimation:
 - Prince, UDL book, ch. 5 to 5.5
 - Deisenroth, [Mathematics for deep learning](#), ch. 8, 9
 - Bishop, [Pattern recognition and machine learning](#), ch. 1.2.5, 1.2.6, 3.1, 3.3 and 4.3
- Gradient checkpointing:
 - I actually believe the [first google result](#) is a good resource :), especially if you check the papers linked at the bottom of the post.
- Ensembling in the weight space, some good starting points:
 - Izmailov, [Averaging weights leads to wider optima and better generalization](#)
 - Wortsman, [Learning neural network subspaces](#)