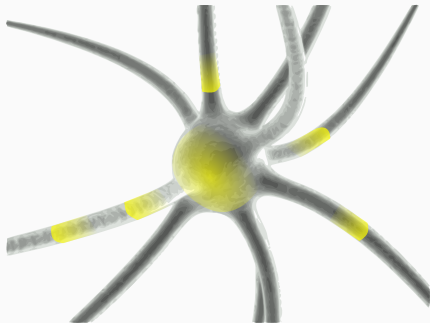# Biological neurons meet zeroth order optimisation

Wouter M. Koolen

CWI and University of Twente

Dutch Day on Optimization, November 7, 2024

## Warm Thanks


Johannes Schmidt-Hieber

**Menu**

# Neuroscience Motivation

**Motivation**

We seek to understand biological neural networks.

E.g. the brain.

And how they learn

Today in particular: connections with optimisation
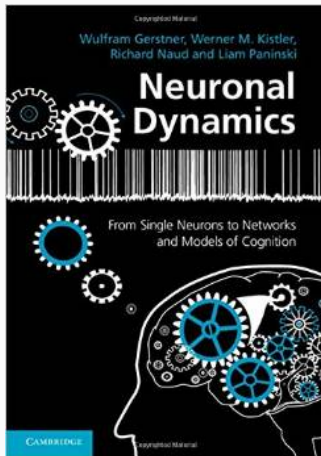
# Why is this interesting



Artificial NN

VS
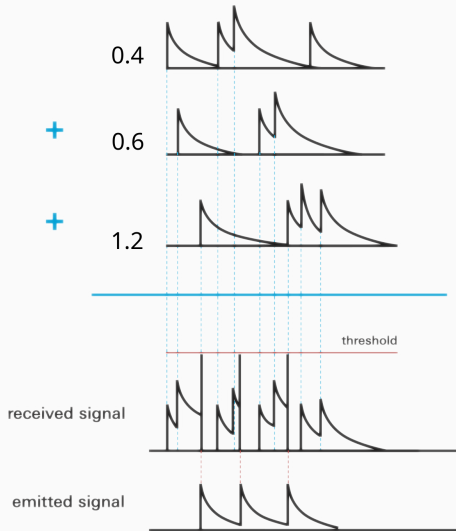
Biological NN

# Literature
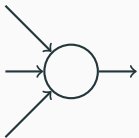


Gerstner, Kistler, Naud, and Paninski, 2014
Chapter 19: Synaptic Plasticity and Learning



Schmidt-Hieber, 2023

## Spiking Neurons



0.4

+ 0.6

+ 1.2

threshold

received signal

emitted signal

It takes multiple (20–50) concurrent incoming spikes to pass the threshold.

## Updating the Weights: Hebbian / Spike Timing-Dependent Plasticity



Let $w_{ij}$ denote the weight from a sending (presynaptic) neuron $j$ to a receiving (postsynaptic) neuron $i$. Then

$$\Delta w_{ij} = Cw_{ij}e^{-c|\Delta t|} \text{ at } t_{post} \text{ for } t_{pre} < t_{post}$$
$$\Delta w_{ij} = -Cw_{ij}e^{-c|\Delta t|} \text{ at } t_{pre} \text{ for } t_{pre} > t_{post}$$

where $\Delta t := |t_{post} - t_{pre}|$.

## Simplifying the weight update

Upon a presynaptic spike at time $\tau$, with previous and next postsynaptic spikes at $T_-$ and $T_+$,

$$w_{ij} \leftarrow w_{ij} + w_{ij} C(-e^{-c(\tau - T_-)} + e^{-c(T_+ - \tau)}).$$

## Simplifying the weight update

Upon a presynaptic spike at time $\tau$, with previous and next postsynaptic spikes at $T_-$ and $T_+$,

$$w_{ij} \leftarrow w_{ij} + w_{ij}C(-e^{-c(\tau-T_-)} + e^{-c(T_+-\tau)}).$$

For **supervised learning**, updates modulated by actual loss $L = L(\mathbf{w})$ minus anticipated loss $\bar{L}$

$$w_{ij} \leftarrow w_{ij} + w_{ij}\alpha(L - \bar{L})(e^{-c(\tau-T_-)} - e^{-c(T_+-\tau)}).$$

## Simplifying the weight update

Upon a presynaptic spike at time $\tau$, with previous and next postsynaptic spikes at $T_-$ and $T_+$,

$$w_{ij} \leftarrow w_{ij} + w_{ij} C(-e^{-c(\tau - T_-)} + e^{-c(T_+ - \tau)}).$$

For **supervised learning**, updates modulated by actual loss $L = L(\mathbf{w})$ minus anticipated loss $\bar{L}$

$$w_{ij} \leftarrow w_{ij} + w_{ij} \alpha(L - \bar{L})(e^{-c(\tau - T_-)} - e^{-c(T_+ - \tau)}).$$

Spike times $\tau$ are roughly **uniform** in window $[T_+ - T_-]$, which we assume is of fixed length $2A$. So with $U$ independent uniform from $[\pm A]$:

$$w_{ij} \leftarrow w_{ij} + w_{ij} \alpha(L - \bar{L})(e^{-c(A+U)} - e^{-c(A-U)}).$$

## Simplifying the weight update

Upon a presynaptic spike at time $\tau$, with previous and next postsynaptic spikes at $T_-$ and $T_+$,

$$w_{ij} \leftarrow w_{ij} + w_{ij}C(-e^{-c(\tau-T_-)} + e^{-c(T_+-\tau)}).$$

For **supervised learning**, updates modulated by actual loss $L = L(\mathbf{w})$ minus anticipated loss $\bar{L}$

$$w_{ij} \leftarrow w_{ij} + w_{ij}\alpha(L - \bar{L})(e^{-c(\tau-T_-)} - e^{-c(T_+-\tau)}).$$

Spike times $\tau$ are roughly **uniform** in window $[T_+ - T_-]$, which we assume is of fixed length $2A$. So with $U$ independent uniform from $[\pm A]$:

$$w_{ij} \leftarrow w_{ij} + w_{ij}\alpha(L - \bar{L})(e^{-c(A+U)} - e^{-c(A-U)}).$$

Go to **logarithmic scale**, $\theta_{ij} := \ln w_{ij}$, and absorb constants

$$\theta_{ij} \leftarrow \theta_{ij} + \ln\left(1 + \alpha(L - \bar{L})(e^{-U} - e^{-U})\right)$$
$$\approx \theta_{ij} + \alpha(L - \bar{L})(e^{-U} - e^{-U})$$

## Fixing the details

We have arrived at the update rule from (Schmidt-Hieber, 2023). Here in vector form.
After round $k$ with data $\mathbf{X}_k, Y_k$, update to

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k \big( L(\boldsymbol{\theta}_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - \overline{L}_k \big) \big( e^{-\mathbf{U}_k} - e^{\mathbf{U}_k} \big).$$

## Fixing the details

We have arrived at the update rule from (Schmidt-Hieber, 2023). Here in vector form.

After round $k$ with data $\mathbf{X}_k, Y_k$, update to

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k \big( L(\boldsymbol{\theta}_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - \overline{L}_k \big) \big( e^{-\mathbf{U}_k} - e^{\mathbf{U}_k} \big).$$

Remains to choose the anticipated loss $\overline{L}_k$

- Zero?
- Loss of last round?
- Average loss so far?
- Treat as separate prediction task?
- Loss at $\boldsymbol{\theta}_{k-1} - \mathbf{U}_k$ (in theory a very good idea)
- ...

## Simplified/concrete/tractable form

Use loss with independent realisation $\mathbf{U}'$ of noise as the anticipated loss $\bar{L}_k$.

$$\boldsymbol{\theta}_k \;=\; \boldsymbol{\theta}_{k-1} + \alpha_k \Big( L(\boldsymbol{\theta}_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - L(\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k) \Big) \big( e^{-\mathbf{U}_k} - e^{\mathbf{U}_k} \big).$$

**What is this scheme?**

Let's stare at it:

$$\theta_k = \theta_{k-1} + \alpha_k \Big( L(\theta_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - L(\theta_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k) \Big) \big( e^{-\mathbf{U}_k} - e^{\mathbf{U}_k} \big).$$

- Zeroth order: evaluates loss $L(\theta_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k)$, **no derivatives**
- Two-point scheme: for each data item $\mathbf{X}_k, Y_k$ evaluate loss of *two* parameters $\theta_{k-1} + \mathbf{U}_k$ and $\theta_{k-1} + \mathbf{U}'_k$

# Benchmark Task: Zeroth order Linear Regression

## Linear Regression

Model ($\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2$): well-specified linear regression with random design.

## Linear Regression

Model $(\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\boldsymbol{\theta}^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

## Linear Regression

Model $(\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\boldsymbol{\theta}^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

- Covariates $\mathbf{X}_1, \mathbf{X}_2, \ldots$ are drawn i.i.d. from $\mathbb{P}$.

## Linear Regression

Model $(\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\boldsymbol{\theta}^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

- Covariates $\mathbf{X}_1, \mathbf{X}_2, \ldots$ are drawn i.i.d. from $\mathbb{P}$.
- Response variables are $Y_k := \mathbf{X}_k^\mathsf{T} \boldsymbol{\theta}^\star + \epsilon_k$ with independent Gaussian noise $\epsilon_k \sim \mathcal{N}(0, \sigma^2)$.

## Linear Regression

Model $(\theta^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\theta^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

- Covariates $\mathbf{X}_1, \mathbf{X}_2, \ldots$ are drawn i.i.d. from $\mathbb{P}$.
- Response variables are $Y_k := \mathbf{X}_k^\mathsf{T}\theta^\star + \epsilon_k$ with independent Gaussian noise $\epsilon_k \sim \mathcal{N}(0, \sigma^2)$.
- Loss of parameter $\theta$ on data item $\mathbf{X}, Y$ is the **square loss**

$$L(\theta, \mathbf{X}, Y) := (\mathbf{X}^\mathsf{T}\theta - Y)^2$$

## Linear Regression

Model $(\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\boldsymbol{\theta}^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

- Covariates $\mathbf{X}_1, \mathbf{X}_2, \dots$ are drawn i.i.d. from $\mathbb{P}$.
- Response variables are $Y_k := \mathbf{X}_k^\mathsf{T}\boldsymbol{\theta}^\star + \epsilon_k$ with independent Gaussian noise $\epsilon_k \sim \mathcal{N}(0, \sigma^2)$.
- Loss of parameter $\boldsymbol{\theta}$ on data item $\mathbf{X}, Y$ is the **square loss**
$$L(\boldsymbol{\theta}, \mathbf{X}, Y) := (\mathbf{X}^\mathsf{T}\boldsymbol{\theta} - Y)^2$$
- Risk (expected loss) of parameter $\boldsymbol{\theta}$ is
$$\mathbb{E}\big[L(\boldsymbol{\theta}, \mathbf{X}, Y)\big] = \mathbb{E}\big[(\mathbf{X}^\mathsf{T}(\boldsymbol{\theta} - \boldsymbol{\theta}^\star) - \epsilon)^2\big] = \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_Q^2 + \sigma^2$$
where we write $Q := \mathbb{E}[\mathbf{X}\mathbf{X}^\mathsf{T}] \succ 0$ for the (uncentred) covariance matrix of the covariates.

## Linear Regression

Model $(\boldsymbol{\theta}^\star, \mathbb{P}, \sigma^2)$: well-specified linear regression with random design.

unknown **true regression coefficient** $\boldsymbol{\theta}^\star \in \mathbb{R}^d$
unknown **covariate distribution** $\mathbb{P}$ on $\mathbb{R}^d$ and
known **noise level** $\sigma > 0$.

- Covariates $\mathbf{X}_1, \mathbf{X}_2, \ldots$ are drawn i.i.d. from $\mathbb{P}$.
- Response variables are $Y_k := \mathbf{X}_k^\mathsf{T}\boldsymbol{\theta}^\star + \epsilon_k$ with independent Gaussian noise $\epsilon_k \sim \mathcal{N}(0, \sigma^2)$.
- Loss of parameter $\boldsymbol{\theta}$ on data item $\mathbf{X}, Y$ is the **square loss**

$$L(\boldsymbol{\theta}, \mathbf{X}, Y) := (\mathbf{X}^\mathsf{T}\boldsymbol{\theta} - Y)^2$$

- Risk (expected loss) of parameter $\boldsymbol{\theta}$ is

$$\mathbb{E}\big[L(\boldsymbol{\theta}, \mathbf{X}, Y)\big] = \mathbb{E}\big[(\mathbf{X}^\mathsf{T}(\boldsymbol{\theta} - \boldsymbol{\theta}^\star) - \epsilon)^2\big] = \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_Q^2 + \sigma^2$$

where we write $Q := \mathbb{E}[\mathbf{X}\mathbf{X}^\mathsf{T}] \succ 0$ for the (uncentred) covariance matrix of the covariates.

- Excess risk of $\boldsymbol{\theta}$ over risk minimiser $\boldsymbol{\theta}^\star$ is

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_Q^2$$

**Interaction Protocol: Two-Point Zeroth-order Stochastic Optimization**

For $k = 1, 2, \ldots$

1. Learner picks *two* query points $\theta_{k-1}^{(1)}$ and $\theta_{k-1}^{(2)}$
2. Data item $\mathbf{X}_k, Y_k$ is drawn from linear regression model behind the scenes
3. Learner observes losses $L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k)$ and $L(\theta_{k-1}^{(2)}, \mathbf{X}_k, Y_k)$ of the two query points
4. Learner recommends evaluation point $\theta_k$

## Interaction Protocol: Two-Point Zeroth-order Stochastic Optimization

For $k = 1, 2, \ldots$

1. Learner picks *two* query points $\theta_{k-1}^{(1)}$ and $\theta_{k-1}^{(2)}$
2. Data item $\mathbf{X}_k, Y_k$ is drawn from linear regression model behind the scenes
3. Learner observes losses $L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k)$ and $L(\theta_{k-1}^{(2)}, \mathbf{X}_k, Y_k)$ of the two query points
4. Learner recommends evaluation point $\theta_k$

NB: Learner has **no access** to data $\mathbf{X}_k, Y_k$ or gradient $\nabla_\theta L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k), \ldots$

## Interaction Protocol: Two-Point Zeroth-order Stochastic Optimization

For $k = 1, 2, \ldots$

1. Learner picks *two* query points $\theta_{k-1}^{(1)}$ and $\theta_{k-1}^{(2)}$
2. Data item $\mathbf{X}_k, Y_k$ is drawn from linear regression model behind the scenes
3. Learner observes losses $L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k)$ and $L(\theta_{k-1}^{(2)}, \mathbf{X}_k, Y_k)$ of the two query points
4. Learner recommends evaluation point $\theta_k$

NB: Learner has **no access** to data $\mathbf{X}_k, Y_k$ or gradient $\nabla_\theta L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k), \ldots$

We are interested in the excess risk of the evaluation point $\theta_k$ as a function of time $k$.

## Interaction Protocol: Two-Point Zeroth-order Stochastic Optimization

For $k = 1, 2, \ldots$

1. Learner picks *two* query points $\theta_{k-1}^{(1)}$ and $\theta_{k-1}^{(2)}$
2. Data item $\mathbf{X}_k, Y_k$ is drawn from linear regression model behind the scenes
3. Learner observes losses $L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k)$ and $L(\theta_{k-1}^{(2)}, \mathbf{X}_k, Y_k)$ of the two query points
4. Learner recommends evaluation point $\theta_k$

NB: Learner has **no access** to data $\mathbf{X}_k, Y_k$ or gradient $\nabla_\theta L(\theta_{k-1}^{(1)}, \mathbf{X}_k, Y_k), \ldots$

We are interested in the excess risk of the evaluation point $\theta_k$ as a function of time $k$.

The evaluation point $\theta_k$ is random due to random data $\mathbf{X}_1, Y_1, \ldots$ (and randomised queries)
So we evaluate a strategy for Learner by its expected excess risk after $k$ rounds

$$\mathbb{E}_{(\theta_0^{(1)}, \theta_0^{(2)}, \mathbf{X}_1, Y_1) \ldots (\theta_{k-1}^{(1)}, \theta_{k-1}^{(2)}, \mathbf{X}_k, Y_k)} \left[ \|\theta_k - \theta^\star\|_Q^2 \right]$$

## Impact of the Query model for Linear Regression

If we query at $\theta$, we see the scalar loss

$$L = (\mathbf{X}^\mathsf{T}\theta - Y)^2 = (\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon)^2$$

**Impact of the Query model for Linear Regression**

If we query at $\theta$, we see the scalar loss

$$L = (\mathbf{X}^\mathsf{T}\theta - Y)^2 = (\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon)^2$$

If further $\mathbf{X} \sim \mathbb{P} = \mathcal{N}(0, I)$ for simplicity, we have

$$\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon \quad \sim \quad \mathcal{N}\left(0, \|\theta - \theta^\star\|^2 + \sigma^2\right)$$

so that the loss is scaled chi-squared

$$L = (\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon)^2 \quad \sim \quad \left(\|\theta - \theta^\star\|^2 + \sigma^2\right)\chi_1^2$$

## Impact of the Query model for Linear Regression

If we query at $\theta$, we see the scalar loss

$$L = (\mathbf{X}^\mathsf{T}\theta - Y)^2 = (\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon)^2$$

If further $\mathbf{X} \sim \mathbb{P} = \mathcal{N}(0, I)$ for simplicity, we have

$$\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon \;\sim\; \mathcal{N}\left(0, \|\theta - \theta^\star\|^2 + \sigma^2\right)$$

so that the loss is scaled chi-squared

$$L = (\mathbf{X}^\mathsf{T}(\theta - \theta^\star) - \epsilon)^2 \;\sim\; \left(\|\theta - \theta^\star\|^2 + \sigma^2\right)\chi_1^2$$

Multiplicative noise. Very different from additive noise $L \sim \left(\|\theta - \theta^\star\|^2 + \sigma^2\right) + \mathcal{N}(0, \text{const})$.

## How hard is this task?

Minimax lower bound for any two-point scheme $\mathcal{V}_k, \widehat{\boldsymbol{\theta}}$.

**Theorem**

If $d \geq 3$ and $k \geq d^2$, then,

$$\inf_{\mathcal{V}_k, \widehat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^\star \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}^\star, \mathcal{V}_k} \left[ \| \widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star \|^2 \right] \geq \frac{1}{162} \left( 1 - \frac{1}{\sqrt{2}} \right) \left( R^2 \wedge \frac{d^2}{k} \sigma^2 \right).$$

## How hard is this task?

Minimax lower bound for any two-point scheme $\mathcal{V}_k, \widehat{\boldsymbol{\theta}}$.

**Theorem**

If $d \geq 3$ and $k \geq d^2$, then,

$$\inf_{\mathcal{V}_k, \widehat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^\star \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}^\star, \mathcal{V}_k} \left[ \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star\|^2 \right] \geq \frac{1}{162} \left( 1 - \frac{1}{\sqrt{2}} \right) \left( R^2 \wedge \frac{d^2}{k} \sigma^2 \right).$$

Minimax excess risk lower bound for non-adaptive two-point schemes

**Theorem**

If $d \geq 6$, then for any $k = 1, 2, \ldots$

$$\inf_{\mathcal{V}_k \in \mathcal{M}_k, \widehat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^\star \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}^\star, \mathcal{V}_k} \left[ \|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}^\star\|^2 \right] \geq 2^{-18} \left( R^2 \wedge \frac{d^2}{k} (R^2 \vee \sigma^2) \right).$$

# BNN meets Linear Regression

**Slogan**

For our combination of loss and update, (almost) everything is fully explicit **linear/quadratic**.

## Connecting BNN to 2P-0O-StochOpt

We query at

$$\theta^{(1)}_{k-1} \;=\; \theta_{k-1} + \mathbf{U}_k \qquad\qquad\qquad \theta^{(2)}_{k-1} \;=\; \theta_{k-1} + \mathbf{U}'_k$$

## Connecting BNN to 2P-0O-StochOpt

We query at

$$\theta_{k-1}^{(1)} \;=\; \theta_{k-1} + \mathbf{U}_k \qquad\qquad \theta_{k-1}^{(2)} \;=\; \theta_{k-1} + \mathbf{U}_k'$$

and update using

$$
\begin{aligned}
\theta_k \;&=\; \theta_{k-1} + \alpha_k\Big((\mathbf{X}_k^\mathsf{T}(\theta_{k-1} + \mathbf{U}_k) - Y_k)^2 - (\mathbf{X}_k^\mathsf{T}(\theta_{k-1} + \mathbf{U}_k') - Y_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big) \\
&=\; \theta_{k-1} + \alpha_k\Big((\mathbf{X}_k^\mathsf{T}(\theta_{k-1} - \theta^\star + \mathbf{U}_k) - \epsilon_k)^2 - (\mathbf{X}_k^\mathsf{T}(\theta_{k-1} - \theta^\star + \mathbf{U}_k') - \epsilon_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).
\end{aligned}
$$

## Connecting BNN to 2P-0O-StochOpt

We query at

$$\theta_{k-1}^{(1)} = \theta_{k-1} + \mathbf{U}_k \qquad\qquad \theta_{k-1}^{(2)} = \theta_{k-1} + \mathbf{U}_k'$$

and update using

$$\theta_k = \theta_{k-1} + \alpha_k\Big((\mathbf{X}_k^\mathsf{T}(\theta_{k-1} + \mathbf{U}_k) - Y_k)^2 - (\mathbf{X}_k^\mathsf{T}(\theta_{k-1} + \mathbf{U}_k') - Y_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big)$$

$$= \theta_{k-1} + \alpha_k\Big((\mathbf{X}_k^\mathsf{T}(\theta_{k-1} - \theta^\star + \mathbf{U}_k) - \epsilon_k)^2 - (\mathbf{X}_k^\mathsf{T}(\theta_{k-1} - \theta^\star + \mathbf{U}_k') - \epsilon_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).$$

So with $\delta_k := \theta_k - \theta^\star$, we get the recurrence

$$\delta_k = \delta_{k-1} + \alpha_k\Big((\mathbf{X}_k^\mathsf{T}(\delta_{k-1} + \mathbf{U}_k) - \epsilon_k)^2 - (\mathbf{X}_k^\mathsf{T}(\delta_{k-1} + \mathbf{U}_k') - \epsilon_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).$$

$$= \delta_{k-1} + \alpha_k\Big(2(\mathbf{X}_k^\mathsf{T}\delta_{k-1} - \epsilon_k)\mathbf{X}_k^\mathsf{T}(\mathbf{U}_k' - \mathbf{U}_k) + (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k)^2 - (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k')^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).$$

$$= \Big(I + 2\alpha_k\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big)(\mathbf{U}_k' - \mathbf{U}_k)^\mathsf{T}\mathbf{X}_k\mathbf{X}_k^\mathsf{T}\Big)\delta_{k-1}$$

$$+ \alpha_k\Big(-2\epsilon_k\mathbf{X}_k^\mathsf{T}(\mathbf{U}_k' - \mathbf{U}_k) + (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k)^2 - (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k')^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big)$$

### Does it even make sense on average?

We expressed our update rule in Stochastic Approximation form

$$\delta_k \; = \; (I - \alpha_k \mathbf{A}_k)\delta_{k-1} + \alpha_k \mathbf{b}_k$$

for i.i.d. random matrix $\mathbf{A}_k$ and vector $\mathbf{b}_k$ given by

$$\mathbf{A}_k \; := \; -2\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big)(\mathbf{U}'_k - \mathbf{U}_k)^\mathsf{T}\mathbf{X}_k\mathbf{X}_k^\mathsf{T},$$

$$\mathbf{b}_k \; := \; \Big(-2\epsilon_k \mathbf{X}_k^\mathsf{T}(\mathbf{U}'_k - \mathbf{U}_k) + (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k)^2 - (\mathbf{X}_k^\mathsf{T}\mathbf{U}'_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).$$

## Does it even make sense on average?

We expressed our update rule in Stochastic Approximation form

$$\delta_k = (I - \alpha_k \mathbf{A}_k)\delta_{k-1} + \alpha_k \mathbf{b}_k$$

for i.i.d. random matrix $\mathbf{A}_k$ and vector $\mathbf{b}_k$ given by

$$\mathbf{A}_k := -2\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big)(\mathbf{U}'_k - \mathbf{U}_k)^\mathsf{T}\mathbf{X}_k\mathbf{X}_k^\mathsf{T},$$
$$\mathbf{b}_k := \Big(-2\epsilon_k\mathbf{X}_k^\mathsf{T}(\mathbf{U}'_k - \mathbf{U}_k) + (\mathbf{X}_k^\mathsf{T}\mathbf{U}_k)^2 - (\mathbf{X}_k^\mathsf{T}\mathbf{U}'_k)^2\Big)\big(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}\big).$$

We have $\mathbb{E}[\mathbf{b}] = 0$ and $\mathbb{E}[\mathbf{A}] = \eta Q$ with constant $\eta := 2\,\mathbb{E}[(e^{-U} - e^{U})U]$ depending on the scale $A$ of noise $\mathbf{U}$.

In expectation, our update gives

$$\mathbb{E}_k[\delta_k] = \Big(I - \alpha_k\eta Q\Big)\delta_{k-1}$$

That is **exactly** gradient descent on the **risk** $\|\delta\|_Q^2 + \sigma^2$, with learning rate $\frac{1}{2}\alpha_k\eta$.

**Case closed?**

So the average iterate $\mathbb{E}[\boldsymbol{\theta}_k] \to \boldsymbol{\theta}^\star$ converges to the risk minimiser. Exponentially fast.

**Case closed?**

So the average iterate $\mathbb{E}[\boldsymbol{\theta}_k] \to \boldsymbol{\theta}^\star$ converges to the risk minimiser. Exponentially fast.

The metric of interest is excess risk $\|\boldsymbol{\delta}\|_Q^2$. Variance matters!

**Case closed?**

So the average iterate $\mathbb{E}[\theta_k] \to \theta^\star$ converges to the risk minimiser. Exponentially fast.

The metric of interest is excess risk $\|\delta\|_Q^2$. Variance matters!

So let's work on the expected excess risk after $k$ rounds (whp bounds also interesting):

$$\Xi_k \coloneqq \mathbb{E}\left[\|\delta_k\|_Q^2\right] \qquad \text{where} \qquad Q = \mathbb{E}[\mathbf{X}\mathbf{X}^\mathsf{T}]$$

Can we get a recurrence for $\Xi_k$? Yes!

## Recurrence for excess risk

Recall our update rule is of the form

$$\delta_k = (I - \alpha_k \mathbf{A}_k)\delta_{k-1} + \alpha_k \mathbf{b}_k$$

for i.i.d. random matrix $\mathbf{A}_k$ and vector $\mathbf{b}_k$, with $\mathbb{E}[\mathbf{A}] = \eta Q$, and $\mathbb{E}[\mathbf{b}] = \mathbb{E}[\mathbf{A}^\mathsf{T} Q \mathbf{b}] = 0$.

**Recurrence for excess risk**

Recall our update rule is of the form

$$\boldsymbol{\delta}_k \;=\; (I - \alpha_k \mathbf{A}_k)\boldsymbol{\delta}_{k-1} + \alpha_k \mathbf{b}_k$$

for i.i.d. random matrix $\mathbf{A}_k$ and vector $\mathbf{b}_k$, with $\mathbb{E}[\mathbf{A}] = \eta Q$, and $\mathbb{E}[\mathbf{b}] = \mathbb{E}[\mathbf{A}^\mathsf{T} Q \mathbf{b}] = 0$.

So the excess risk satisfies

$$
\begin{aligned}
\Xi_k \;&=\; \mathbb{E}_k\left[\boldsymbol{\delta}_k^\mathsf{T} Q \boldsymbol{\delta}_k\right] \\
&=\; \mathbb{E}_k\left[((I - \alpha_k \mathbf{A}_k)\boldsymbol{\delta}_{k-1} + \alpha_k \mathbf{b}_k)^\mathsf{T} Q((I - \alpha_k \mathbf{A}_k)\boldsymbol{\delta}_{k-1} + \alpha_k \mathbf{b}_k)\right] \\
&=\; \boldsymbol{\delta}_{k-1}^\mathsf{T} \mathbb{E}_k\left[(I - \alpha_k \mathbf{A}_k)^\mathsf{T} Q(I - \alpha_k \mathbf{A}_k)\right]\boldsymbol{\delta}_{k-1} + \alpha_k^2\,\mathbb{E}_k\left[\mathbf{b}_k^\mathsf{T} Q \mathbf{b}_k\right] \\
&=\; \boldsymbol{\delta}_{k-1}^\mathsf{T}\left\{(I - \alpha_k \eta Q)^\mathsf{T} Q(I - \alpha_k \eta Q) + \alpha_k^2\,\mathbb{E}_k\left[(\mathbf{A}_k - \eta Q)^\mathsf{T} Q(\mathbf{A}_k - \eta Q)\right]\right\}\boldsymbol{\delta}_{k-1} + \alpha_k^2\,\mathbb{E}_k\left[\mathbf{b}_k^\mathsf{T} Q \mathbf{b}_k\right] \\
&\leq\; \left((1 - \alpha_k \eta \lambda_{\min}(Q))^2 + \alpha_k^2 \beta\right)\Xi_{k-1} + \alpha_k^2 \gamma
\end{aligned}
$$

abbreviating $\beta := \lambda_{\max}\left(\mathbb{E}_k\left[Q^{-1/2}(\mathbf{A}_k - \eta Q)^\mathsf{T} Q(\mathbf{A}_k - \eta Q)Q^{-1/2}\right]\right)$ and $\gamma := \mathbb{E}_k\left[\mathbf{b}_k^\mathsf{T} Q \mathbf{b}_k\right]$.

**Inspecting where we are**

Our state of progress so far is

$$\Xi_k \leq \left((1 - \alpha_k \eta \lambda_{\min}(Q))^2 + \alpha_k^2 \beta\right)\Xi_{k-1} + \alpha_k^2 \gamma$$

for fixed $\eta$, $\lambda_{\min}(Q)$, $\beta$ and $\gamma$. The question is how to tune $\alpha_k$. This is now a scalar problem.

**Inspecting where we are**

Our state of progress so far is

$$\Xi_k \leq \left((1 - \alpha_k \eta \lambda_{\min}(Q))^2 + \alpha_k^2 \beta\right)\Xi_{k-1} + \alpha_k^2 \gamma$$

for fixed $\eta$, $\lambda_{\min}(Q)$, $\beta$ and $\gamma$. The question is how to tune $\alpha_k$. This is now a scalar problem.

Cancelling derivative reveals this bound is optimised in $\alpha_k$ at

$$\alpha_k^* = \frac{\eta \lambda_{\min}(Q)}{\eta^2 \lambda_{\min}(Q)^2 + \beta + \frac{\gamma}{\Xi_{k-1}}}$$

and at that point we obtain

$$\Xi_k \leq \left(\frac{\beta + \frac{\gamma}{\Xi_{k-1}}}{\eta^2 \lambda_{\min}(Q)^2 + \beta + \frac{\gamma}{\Xi_{k-1}}}\right)\Xi_{k-1}$$

## Cute ODE upper bound

We can write our recurrence so far as a difference equation

$$\frac{\Xi_k - \Xi_{k-1}}{\Xi_{k-1}} \leq -\frac{\eta^2 \lambda_{\min}(Q)^2}{\eta^2 \lambda_{\min}(Q)^2 + \beta + \frac{\gamma}{\Xi_{k-1}}}$$

and solve the corresponding **differential equation** with equality to find

$$\frac{\Xi_k}{\Xi_1} \leq \frac{y}{W(ye^{y+xk})} \quad \text{with} \quad x := \frac{\eta^2 \lambda_{\min}(Q)^2}{\eta^2 \lambda_{\min}(Q)^2 + \beta} \quad \text{and} \quad y := \frac{\gamma/\Xi_1}{\eta^2 \lambda_{\min}(Q)^2 + \beta}$$

so that all in all the excess risk decays as $\Xi_k \cong \Xi_1/k$ and the **learning rate** as $\alpha_k^* \cong 1/k$.

**More precisely in terms of relevant problem-dependent constants**

We arrive at excess risk bound

**Theorem**

$$\Xi_k \; \leq \; \frac{121\kappa d^2}{2\lambda_{\min}(Q)} \frac{48\sigma^2 M_2 + 107A^2 dM_4}{k + C}$$

*where $\kappa = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)}$ is the condition number of $Q$, and $M_p$ bounds the $i$th moment of each entry of the covariate vector $\mathbf{X} \sim \mathbb{P}$.*

If $A^2 d$ is at most of order $\sigma^2$, this is $d^2/k$. Matching lower bounds.

# Reflections

**To think about**

- Is the optimal tuning $\alpha_k \cong 1/k$ biologically realistic?
- Learning rate $\alpha_k$ needs to decay. What decides a *new task* in the brain?
- Optimal tuning for $\alpha_k$ depends on zoo of unknowns. How are these estimated?
- Brutal tuning $\alpha_k = \frac{c}{C+k}$ may result in risk rising to $e^{\text{const}}$ before $1/k$ decay kicks in.
- Is the noise rate $A$ biologically small compared to $\sigma/\sqrt{d}$?
- Realism in the model
    - More than one neuron
    - Depth, architecture
    - Other tasks and losses

**Conclusion**

We saw a simple model for spiking neurons inspired by biology.

We saw a concrete rendering of resulting update rule.

We interpreted it as a zeroth-order two-point iterative scheme.

We evaluated this scheme on a linear regression task.

We derived a rate for the excess risk, and proved that it matches lower bounds.

# Let's talk!

# References i

📄 Gerstner, W., W. M. Kistler, R. Naud, and L. Paninski (2014). **Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition.** Cambridge University Press.

📄 Schmidt-Hieber, J. (2023). **"Interpreting learning in biological neural networks as zero-order optimization method".** In: *arXiv preprint*, arXiv:2301.11777.