

Effective Theory of Deep Learning

Beyond the Infinite-Width Limit

Dan Roberts^a and Sho Yaida^b

^aMIT, IAIFI, & Salesforce, ^bFacebook AI Research

Deep Learning Theory Summer School at Princeton
July 27, 2021 – August 8, 2021

Course Plan

~~Lecture 1~~ **Initialization, Linear Models**

▶ ~~§0 + §7.1 + §10.4~~

Lecture 2 Quadratic Models & Nearly-Kernel Methods

▶ §11.4 (+ §7.2) + § ∞ .2.2

Lecture 3 The Principle of Sparsity (Recurring)

▶ §4, §8, §11.2, § ∞ .3

Lecture 4 The Principle of Criticality

▶ §5, §9, §11.3, § ∞ .1, +§10.3

Lecture 5 The End of Training & More

▶ § ∞ .2.3 + Maybe { §A, §B, ... }

Linear Models and Kernel Methods

Two forms of a solution for a **linear model**:

- ▶ *parameter space – linear regression*

$$z_i(x_{\hat{\beta}}; \theta^*) = \sum_{j=0}^{n_f} W_{ij}^* \phi_j(x_{\hat{\beta}})$$

- ▶ *sample space – kernel methods*

$$z_i(x_{\hat{\beta}}; \theta^*) = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} k_{\hat{\beta}\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2} .$$

Linear Models and Kernel Methods

Two forms of a solution for a **linear model**:

- ▶ *parameter space – linear regression*

$$z_i(x_{\dot{\beta}}; \theta^*) = \sum_{j=0}^{n_f} W_{ij}^* \phi_j(x_{\dot{\beta}})$$

- ▶ *sample space – kernel methods*

$$z_i(x_{\dot{\beta}}; \theta^*) = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} k_{\dot{\beta}\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2}.$$

Features of this model, expressed as $\phi_j(x)$ or $k_{\delta_1\delta_2}$, are *fixed*.
The goal of this lecture find a model where they evolve: a **minimal model of representation learning**.

Nonlinear Models

To go beyond the linear paradigm, let's slightly *deform* it to get a **nonlinear model**, specifically a **quadratic model**:

$$z_{i;\delta}(\theta) = \sum_{j=0}^{n_f} W_{ij} \phi_j(x_\delta) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_\delta)$$

Nonlinear Models

To go beyond the linear paradigm, let's slightly *deform* it to get a **nonlinear model**, specifically a **quadratic model**:

$$z_{i;\delta}(\theta) = \sum_{j=0}^{n_f} W_{ij} \phi_j(x_\delta) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_\delta)$$

- ▶ It's nonlinear because it's quadratic in the weights: $W_{ij_1} W_{ij_2}$.

Nonlinear Models

To go beyond the linear paradigm, let's slightly *deform* it to get a **nonlinear model**, specifically a **quadratic model**:

$$z_{i;\delta}(\theta) = \sum_{j=0}^{n_f} W_{ij} \phi_j(x_\delta) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_\delta)$$

- ▶ It's nonlinear because it's quadratic in the weights: $W_{ij_1} W_{ij_2}$.
- ▶ $\epsilon \ll 1$ is small parameter that controls the size of the deformation.

Nonlinear Models

To go beyond the linear paradigm, let's slightly *deform* it to get a **nonlinear model**, specifically a **quadratic model**:

$$z_{i;\delta}(\theta) = \sum_{j=0}^{n_f} W_{ij} \phi_j(x_\delta) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_\delta)$$

- ▶ It's nonlinear because it's quadratic in the weights: $W_{ij_1} W_{ij_2}$.
- ▶ $\epsilon \ll 1$ is small parameter that controls the size of the deformation.
- ▶ We've introduced $(n_f + 1)(n_f + 2)/2$ **meta feature functions**, $\psi_{j_1 j_2}(x)$, with *two* feature indices.

Quadratic Models

To familiarize ourselves with this model, let's make a small change in the model parameters $W_{ij} \rightarrow W_{ij} + dW_{ij}$:

$$z_i(x_\delta; \theta + d\theta) = z_i(x_\delta; \theta) + \sum_{j=0}^{n_f} dW_{ij} \left[\phi_j(x_\delta) + \epsilon \sum_{j_1=0}^{n_f} W_{ij_1} \psi_{j_1 j}(x_\delta) \right] \\ + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} dW_{ij_1} dW_{ij_2} \psi_{j_1 j_2}(x_\delta).$$

Quadratic Models

To familiarize ourselves with this model, let's make a small change in the model parameters $W_{ij} \rightarrow W_{ij} + dW_{ij}$:

$$z_i(x_\delta; \theta + d\theta) = z_i(x_\delta; \theta) + \sum_{j=0}^{n_f} dW_{ij} \left[\phi_j(x_\delta) + \epsilon \sum_{j_1=0}^{n_f} W_{ij_1} \psi_{j_1 j}(x_\delta) \right] \\ + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} dW_{ij_1} dW_{ij_2} \psi_{j_1 j_2}(x_\delta).$$

Let us make a shorthand for the quantity in the square bracket,

$$\phi_{ij}^E(x_\delta; \theta) \equiv \frac{dz_i(x_\delta; \theta)}{dW_{ij}} = \phi_j(x_\delta) + \epsilon \sum_{k=0}^{n_f} W_{ik} \psi_{kj}(x_\delta),$$

which is an **effective feature function**.

Effective Feature Functions

The utility of this is as follows:

Effective Feature Functions

The utility of this is as follows:

- ▶ The *linear response* of $z_i(x_\delta; \theta)$ behaves *effectively* as if it has a parameter-dependent feature function, $\phi_{ij}^E(x_\delta; \theta)$.

Effective Feature Functions

The utility of this is as follows:

- ▶ The *linear response* of $z_i(x_\delta; \theta)$ behaves *effectively* as if it has a parameter-dependent feature function, $\phi_{ij}^E(x_\delta; \theta)$.
- ▶ The change in the $\phi_{ij}^E(x_\delta; \theta)$ given $W_{ik} \rightarrow W_{ik} + dW_{ik}$ is

$$\phi_{ij}^E(x_\delta; \theta + d\theta) = \phi_{ij}^E(x_\delta; \theta) + \epsilon \sum_{k=0}^{n_f} dW_{ik} \psi_{kj}(x_\delta),$$

Effective Feature Functions

The utility of this is as follows:

- ▶ The *linear response* of $z_i(x_\delta; \theta)$ behaves *effectively* as if it has a parameter-dependent feature function, $\phi_{ij}^E(x_\delta; \theta)$.
- ▶ The change in the $\phi_{ij}^E(x_\delta; \theta)$ given $W_{ik} \rightarrow W_{ik} + dW_{ik}$ is

$$\phi_{ij}^E(x_\delta; \theta + d\theta) = \phi_{ij}^E(x_\delta; \theta) + \epsilon \sum_{k=0}^{n_f} dW_{ik} \psi_{kj}(x_\delta),$$

- ▶ For comparison, for the linear model we had:

$$z_i(x_\delta; \theta + d\theta) = z_i(x_\delta; \theta) + \sum_{j=0}^{n_f} dW_{ij} \phi_j(x_\delta)$$

Effective Feature Functions

The utility of this is as follows:

- ▶ The *linear response* of $z_i(x_\delta; \theta)$ behaves *effectively* as if it has a parameter-dependent feature function, $\phi_{ij}^E(x_\delta; \theta)$.
- ▶ The change in the $\phi_{ij}^E(x_\delta; \theta)$ given $W_{ik} \rightarrow W_{ik} + dW_{ik}$ is

$$\phi_{ij}^E(x_\delta; \theta + d\theta) = \phi_{ij}^E(x_\delta; \theta) + \epsilon \sum_{k=0}^{n_f} dW_{ik} \psi_{kj}(x_\delta),$$

- ▶ For comparison, for the linear model we had:

$$z_i(x_\delta; \theta + d\theta) = z_i(x_\delta; \theta) + \sum_{j=0}^{n_f} dW_{ij} \phi_j(x_\delta)$$

Thus quadratic model has a *hierarchical structure*, where the features evolve as if they are described by a linear model and the model's output evolves in a more complicated nonlinear way.

Quadratic Regression

Supervised learning a quadratic model doesn't have a particular name, but if it did, we'd all probably agree that its name should be **quadratic regression**:

$$\mathcal{L}_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) - \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_{\tilde{\alpha}}) \right]^2 .$$

Quadratic Regression

Supervised learning a quadratic model doesn't have a particular name, but if it did, we'd all probably agree that its name should be **quadratic regression**:

$$\mathcal{L}_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) - \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_{\tilde{\alpha}}) \right]^2 .$$

The loss is now *quartic* in the parameters, and in general

$$0 = \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W=W^*} ,$$

doesn't give analytical solutions or a tractable practical method.

Quadratic Regression

Supervised learning a quadratic model doesn't have a particular name, but if it did, we'd all probably agree that its name should be **quadratic regression**:

$$\mathcal{L}_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) - \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_{\tilde{\alpha}}) \right]^2 .$$

The loss is now *quartic* in the parameters, but we can optimize with *gradient descent*:

$$W_{ij}(t+1) = W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} .$$

This will find a minimum in practice.

Aside: Gradient Descent

Gradient descent (GD) can be used to minimize the training loss:

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}} \right|_{\theta_{\mu}=\theta_{\mu}(t)} .$$

Aside: Gradient Descent

Gradient descent (GD) can be used to minimize the training loss:

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}} \right|_{\theta_{\mu}=\theta_{\mu}(t)} .$$

- ▶ GD is an *iterative* learning algorithm, and here t keeps track of the number of steps in the iterative training process.

Aside: Gradient Descent

Gradient descent (GD) can be used to minimize the training loss:

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}} \right|_{\theta_{\mu}=\theta_{\mu}(t)} .$$

- ▶ GD is an *iterative* learning algorithm, and here t keeps track of the number of steps in the iterative training process.
- ▶ $\eta > 0$ is a **training hyperparameter** called the **learning rate**, which controls the size of the step taken in *parameter space*.

Aside: Gradient Descent

Gradient descent (GD) can be used to minimize the training loss:

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}} \right|_{\theta_{\mu}=\theta_{\mu}(t)} .$$

- ▶ GD is an *iterative* learning algorithm, and here t keeps track of the number of steps in the iterative training process.
- ▶ $\eta > 0$ is a **training hyperparameter** called the **learning rate**, which controls the size of the step taken in *parameter space*.
- ▶ The computational cost of gradient descent scales linearly with the size of the dataset \mathcal{A} , as one just needs to compute the gradient for each sample and then add them up.

Aside: Gradient Descent

For sufficiently small η , the GD updates are guaranteed to decrease the training loss $\mathcal{L}_{\mathcal{A}}$:

$$\Delta\mathcal{L}_{\mathcal{A}} \equiv \mathcal{L}_{\mathcal{A}}(\theta(t+1)) - \mathcal{L}_{\mathcal{A}}(\theta(t))$$

Aside: Gradient Descent

For sufficiently small η , the GD updates are guaranteed to decrease the training loss $\mathcal{L}_{\mathcal{A}}$:

$$\begin{aligned}\Delta\mathcal{L}_{\mathcal{A}} &\equiv \mathcal{L}_{\mathcal{A}}(\theta(t+1)) - \mathcal{L}_{\mathcal{A}}(\theta(t)) \\ &= \mathcal{L}_{\mathcal{A}}\left(\theta_{\mu}(t) - \eta \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}(t)}\right) - \mathcal{L}_{\mathcal{A}}(\theta(t))\end{aligned}$$

Aside: Gradient Descent

For sufficiently small η , the GD updates are guaranteed to decrease the training loss $\mathcal{L}_{\mathcal{A}}$:

$$\begin{aligned}\Delta\mathcal{L}_{\mathcal{A}} &\equiv \mathcal{L}_{\mathcal{A}}(\theta(t+1)) - \mathcal{L}_{\mathcal{A}}(\theta(t)) \\ &= \mathcal{L}_{\mathcal{A}}\left(\theta_{\mu}(t) - \eta \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}(t)}\right) - \mathcal{L}_{\mathcal{A}}(\theta(t)) \\ &= -\eta \sum_{\mu} \left(\frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}}\right)^2 \Big|_{\theta=\theta(t)} + O(\eta^2).\end{aligned}$$

Aside: Gradient Descent

For sufficiently small η , the GD updates are guaranteed to decrease the training loss $\mathcal{L}_{\mathcal{A}}$:

$$\begin{aligned}\Delta\mathcal{L}_{\mathcal{A}} &\equiv \mathcal{L}_{\mathcal{A}}(\theta(t+1)) - \mathcal{L}_{\mathcal{A}}(\theta(t)) \\ &= \mathcal{L}_{\mathcal{A}}\left(\theta_{\mu}(t) - \eta \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}(t)}\right) - \mathcal{L}_{\mathcal{A}}(\theta(t)) \\ &= -\eta \sum_{\mu} \left(\frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}}\right)^2 \Big|_{\theta=\theta(t)} + O(\eta^2).\end{aligned}$$

In practice, small variants of gradient descent are responsible for almost all training and optimization in deep learning.

Aside of the Aside: Tensorial Gradient Descent

In one such variant, we modify the update as

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\nu}} \right|_{\theta=\theta(t)},$$

to get a more general family of learning algorithms.

Aside of the Aside: Tensorial Gradient Descent

In one such variant, we modify the update as

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\nu}} \right|_{\theta=\theta(t)},$$

to get a more general family of learning algorithms.

- ▶ Here $\lambda_{\mu\nu}$ is a **learning-rate tensor** on parameter space; *vanilla* GD is a special case with $\lambda_{\mu\nu} = \delta_{\mu\nu}$.

Aside of the Aside: Tensorial Gradient Descent

In one such variant, we modify the update as

$$\theta_{\mu}(t+1) = \theta_{\mu}(t) - \eta \sum_{\nu} \lambda_{\mu\nu} \left. \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\nu}} \right|_{\theta=\theta(t)},$$

to get a more general family of learning algorithms.

- ▶ Here $\lambda_{\mu\nu}$ is a **learning-rate tensor** on parameter space; *vanilla* GD is a special case with $\lambda_{\mu\nu} = \delta_{\mu\nu}$.
- ▶ Repeating our analysis with this generalized update, we find

$$\Delta\mathcal{L}_{\mathcal{A}} = -\eta \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\mu}} \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\nu}} + O(\eta^2),$$

meaning $\Delta\mathcal{L}_{\mathcal{A}}$ will decrease for $\eta \ll 1$, so long as the learning-rate tensor $\lambda_{\mu\nu}$ is a positive semidefinite matrix.

The Theoretical Minimum

Let's start by seeing how gradient descent solves the *linear model*:

$$\mathcal{L}_{\mathcal{A}}(W) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]^2 ,$$

The Theoretical Minimum

Let's start by seeing how gradient descent solves the *linear model*:

$$\mathcal{L}_{\mathcal{A}}(W) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]^2,$$
$$\frac{\partial \mathcal{L}_{\mathcal{A}}(W)}{\partial W_{ab}} = - \sum_{\tilde{\alpha}, i, j} \delta_{ia} \delta_{jb} \phi_j(x_{\tilde{\alpha}}) \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]$$

The Theoretical Minimum

Let's start by seeing how gradient descent solves the *linear model*:

$$\mathcal{L}_{\mathcal{A}}(W) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]^2,$$
$$\frac{\partial \mathcal{L}_{\mathcal{A}}(W)}{\partial W_{ab}} = - \sum_{\tilde{\alpha}, i, j} \delta_{ia} \delta_{jb} \phi_j(x_{\tilde{\alpha}}) \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]$$
$$= \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) (z_{a;\tilde{\alpha}} - y_{a;\tilde{\alpha}})$$

The Theoretical Minimum

Let's start by seeing how gradient descent solves the *linear model*:

$$\begin{aligned}\mathcal{L}_{\mathcal{A}}(W) &= \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]^2, \\ \frac{\partial \mathcal{L}_{\mathcal{A}}(W)}{\partial W_{ab}} &= - \sum_{\tilde{\alpha}, i, j} \delta_{ia} \delta_{jb} \phi_j(x_{\tilde{\alpha}}) \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right] \\ &= \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) (z_{a;\tilde{\alpha}} - y_{a;\tilde{\alpha}}) \\ &= \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) \epsilon_{a;\tilde{\alpha}}\end{aligned}$$

The Theoretical Minimum

Let's start by seeing how gradient descent solves the *linear model*:

$$\begin{aligned}\mathcal{L}_{\mathcal{A}}(W) &= \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right]^2, \\ \frac{\partial \mathcal{L}_{\mathcal{A}}(W)}{\partial W_{ab}} &= - \sum_{\tilde{\alpha}, i, j} \delta_{ia} \delta_{jb} \phi_j(x_{\tilde{\alpha}}) \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) \right] \\ &= \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) (z_{a;\tilde{\alpha}} - y_{a;\tilde{\alpha}}) \\ &= \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) \epsilon_{a;\tilde{\alpha}}\end{aligned}$$

In the last line, we defined the **residual training error**:

$$\epsilon_{i;\tilde{\alpha}} \equiv z_{i;\tilde{\alpha}} - y_{i;\tilde{\alpha}}$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) + \sum_{a,b} \frac{\partial z_{i;\delta}(t)}{\partial W_{ab}} \left[W_{ab}(t+1) - W_{ab}(t) \right] + \dots$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) + \sum_{a,b} \frac{\partial z_{i;\delta}(t)}{\partial W_{ab}} \left[W_{ab}(t+1) - W_{ab}(t) \right]$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) + \sum_{a,b} \frac{\partial z_{i;\delta}(t)}{\partial W_{ab}} \left[-\eta \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) \epsilon_{a;\tilde{\alpha}}(t) \right]$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) + \sum_{a,b} \delta_{ia} \phi_b(x_{\delta}) \left[-\eta \sum_{\tilde{\alpha}} \phi_b(x_{\tilde{\alpha}}) \epsilon_{a;\tilde{\alpha}}(t) \right]$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_b \phi_b(x_{\delta}) \phi_b(x_{\tilde{\alpha}}) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t)$$

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t)$$

- Fixed $k_{\delta\tilde{\alpha}}$ generates the dynamics of the model.

The Theoretical Minimum

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_j(x_{\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

For the theoretical analysis, it's more convenient to understand how the output of the model updates:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t)$$

- ▶ Fixed $k_{\delta\tilde{\alpha}}$ generates the dynamics of the model.
- ▶ $\epsilon_{i;\tilde{\alpha}}(t)$ sources the updates for general inputs $\delta \in \mathcal{D}$.

The Theoretical Minimum

We have to solve a linear difference equation:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t).$$

The Theoretical Minimum

We have to solve a linear difference equation:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t).$$

Restricting to the training set, we get a first-order homogeneous linear difference equation,

$$z_{i;\tilde{\alpha}_1}(t+1) = z_{i;\tilde{\alpha}_1}(t) - \eta \sum_{\tilde{\alpha}_2} k_{\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_2}(t),$$

The Theoretical Minimum

We have to solve a linear difference equation:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t).$$

Restricting to the training set, we get a first-order homogeneous linear difference equation,

$$z_{i;\tilde{\alpha}_1}(t+1) = z_{i;\tilde{\alpha}_1}(t) - \eta \sum_{\tilde{\alpha}_2} k_{\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_2}(t),$$

for the residual training error:

$$\epsilon_{i;\tilde{\alpha}_1}(t+1) = \epsilon_{i;\tilde{\alpha}_1}(t) - \eta \sum_{\tilde{\alpha}_2} k_{\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_2}(t),$$

The Theoretical Minimum

We can rewrite these dynamics:

$$\epsilon_{i;\tilde{\alpha}_1}(t+1) = \epsilon_{i;\tilde{\alpha}_1}(t) - \eta \sum_{\tilde{\alpha}_2} k_{\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_2}(t)$$

The Theoretical Minimum

We can rewrite these dynamics:

$$\epsilon_{i;\tilde{\alpha}_1}(t+1) = \sum_{\tilde{\alpha}_2} (\delta_{\tilde{\alpha}_1\tilde{\alpha}_2} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_2}) \epsilon_{i;\tilde{\alpha}_2}(t)$$

The Theoretical Minimum

We can rewrite these dynamics:

$$\epsilon_{i;\tilde{\alpha}_1}(t+1) = \sum_{\tilde{\alpha}_2} (\delta_{\tilde{\alpha}_1\tilde{\alpha}_2} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_2}) \epsilon_{i;\tilde{\alpha}_2}(t)$$

This is a repeated multiplication by a constant matrix:

$$\begin{aligned} U_{\tilde{\alpha}_t\tilde{\alpha}_0}(t) &\equiv \left[(\delta - \eta k)^t \right]_{\tilde{\alpha}_t\tilde{\alpha}_0} \\ &= \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_{t-1}} (\delta_{\tilde{\alpha}_t\tilde{\alpha}_{t-1}} - \eta k_{\tilde{\alpha}_t\tilde{\alpha}_{t-1}}) \cdots (\delta_{\tilde{\alpha}_1\tilde{\alpha}_0} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_0}). \end{aligned}$$

The Theoretical Minimum

We can rewrite these dynamics:

$$\epsilon_{i;\tilde{\alpha}_1}(t+1) = \sum_{\tilde{\alpha}_2} (\delta_{\tilde{\alpha}_1\tilde{\alpha}_2} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_2}) \epsilon_{i;\tilde{\alpha}_2}(t)$$

This is a repeated multiplication by a constant matrix:

$$\begin{aligned} U_{\tilde{\alpha}_t\tilde{\alpha}_0}(t) &\equiv [(\delta - \eta k)^t]_{\tilde{\alpha}_t\tilde{\alpha}_0} \\ &= \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_{t-1}} (\delta_{\tilde{\alpha}_t\tilde{\alpha}_{t-1}} - \eta k_{\tilde{\alpha}_t\tilde{\alpha}_{t-1}}) \cdots (\delta_{\tilde{\alpha}_1\tilde{\alpha}_0} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_0}). \end{aligned}$$

The solution is given by

$$\epsilon_{i;\tilde{\alpha}_1}(t) = \sum_{\tilde{\alpha}_2} U_{\tilde{\alpha}_1\tilde{\alpha}_2}(t) \epsilon_{i;\tilde{\alpha}_2}(0),$$

and $U(t) \rightarrow 0$ as $t \rightarrow \infty$ so that the error vanishes: $z_{i;\tilde{\alpha}} \rightarrow y_{i;\tilde{\alpha}}$.

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$z_{i;\delta}(t+1) = z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}(t)$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$z_{i;\delta}(t) = z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{t-1} \epsilon_{i;\tilde{\alpha}}(s) \right\}$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$z_{i;\delta}(\infty) = z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}}(s) \right\}$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$\begin{aligned} z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}}(s) \right\} \\ &= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[\sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(s) \epsilon_{i;\tilde{\alpha}_1}(0) \right] \right\} \end{aligned}$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$\begin{aligned}z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}}(s) \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[\sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(s) \epsilon_{i;\tilde{\alpha}_1}(0) \right] \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[(\delta - \eta k)^s \right]_{\tilde{\alpha}\tilde{\alpha}_1} \right\} \epsilon_{i;\tilde{\alpha}_1}(0)\end{aligned}$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$\begin{aligned}z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}}(s) \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[\sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(s) \epsilon_{i;\tilde{\alpha}_1}(0) \right] \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} [(\delta - \eta k)^s]_{\tilde{\alpha}\tilde{\alpha}_1} \right\} \epsilon_{i;\tilde{\alpha}_1}(0) \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \left[\delta - (\delta - \eta k) \right]^{-1} \right\}^{\tilde{\alpha}\tilde{\alpha}_1} \epsilon_{i;\tilde{\alpha}_1}(0)\end{aligned}$$

Theoretical Predictions

We *still* have to solve the difference equation for the test error:

$$\begin{aligned}z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}}(s) \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha} \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[\sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(s) \epsilon_{i;\tilde{\alpha}_1}(0) \right] \right\} \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} [(\delta - \eta k)^s]_{\tilde{\alpha}\tilde{\alpha}_1} \right\} \epsilon_{i;\tilde{\alpha}_1}(0) \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \left\{ \eta [\delta - (\delta - \eta k)]^{-1} \right\}^{\tilde{\alpha}\tilde{\alpha}_1} \epsilon_{i;\tilde{\alpha}_1}(0) \\&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} \epsilon_{i;\tilde{\alpha}_1}(0)\end{aligned}$$

Theoretical Predictions

Compare *gradient descent* vs. the *direct optimization* solution:

$$z_{i;\delta}(\infty) = z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta \tilde{\alpha}} \tilde{k}^{\tilde{\alpha} \tilde{\alpha}_1} \epsilon_{i; \tilde{\alpha}_1}(0)$$

$$z_i(x_\delta; \theta^*) = \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha} \tilde{\alpha}_1} y_{i; \tilde{\alpha}_1}.$$

Theoretical Predictions

Compare *gradient descent* vs. the *direct optimization* solution:

$$z_{i;\delta}(\infty) = z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta \tilde{\alpha}} \tilde{k}^{\tilde{\alpha} \tilde{\alpha}_1} \epsilon_{i; \tilde{\alpha}_1}(0)$$

$$z_i(x_\delta; \theta^*) = \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha} \tilde{\alpha}_1} y_{i; \tilde{\alpha}_1}.$$

- ▶ The same if $z_{i;\delta}(0) = 0$, e.g. if $W_{ij}(0) = 0$.

Theoretical Predictions

Compare *gradient descent* vs. the *direct optimization* solution:

$$z_{i;\delta}(\infty) = z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} \epsilon_{i;\tilde{\alpha}_1}(0)$$

$$z_i(x_\delta; \theta^*) = \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} y_{i;\tilde{\alpha}_1}.$$

- ▶ The same if $z_{i;\delta}(0) = 0$, e.g. if $W_{ij}(0) = 0$.
- ▶ Otherwise, linear models have **algorithm independence**.

Theoretical Predictions

Compare *gradient descent* vs. the *direct optimization* solution:

$$z_{i;\delta}(\infty) = z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} \epsilon_{i;\tilde{\alpha}_1}(0)$$

$$z_i(x_\delta; \theta^*) = \sum_{\tilde{\alpha}, \tilde{\alpha}_1 \in \mathcal{A}} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} y_{i;\tilde{\alpha}_1}.$$

- ▶ The same if $z_{i;\delta}(0) = 0$, e.g. if $W_{ij}(0) = 0$.
- ▶ Otherwise, linear models have **algorithm independence**.
- ▶ Importantly, $k_{\delta\tilde{\alpha}_1}$ is fixed, and the $\phi_i(x)$ *do not evolve*.

Quadratic Regression Returns!

Supervised learning a quadratic model doesn't have a particular name, but if it did, we'd all probably agree that its name should be **quadratic regression**:

$$\mathcal{L}_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\tilde{\alpha} \in \mathcal{A}} \sum_{i=1}^{n_{\text{out}}} \left[y_{i;\tilde{\alpha}} - \sum_{j=0}^{n_f} W_{ij} \phi_j(x_{\tilde{\alpha}}) - \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(x_{\tilde{\alpha}}) \right]^2 .$$

The loss is now *quartic* in the parameters, but we can optimize with *gradient descent*:

$$W_{ij}(t+1) = W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_{\mathcal{A}}}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} .$$

This will find a minimum in practice.

Quadratic Model Dynamics

The weights will update as

$$\begin{aligned}W_{ij}(t+1) &= W_{ij}(t) - \eta \left. \frac{d\mathcal{L}_A}{dW_{ij}} \right|_{W_{ij}=W_{ij}(t)} \\ &= W_{ij}(t) - \eta \sum_{\tilde{\alpha}} \phi_{ij;\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t).\end{aligned}$$

While the model and effective features update as

$$\begin{aligned}z_{i;\delta}(t+1) &= z_{i;\delta}(t) + \sum_j dW_{ij}(t) \phi_{ij;\delta}^E(t) \\ &\quad + \frac{\epsilon}{2} \sum_{j_1, j_2} dW_{ij_1}(t) dW_{ij_2}(t) \psi_{j_1 j_2}(\mathbf{x}_\delta), \\ \phi_{ij;\delta}^E(t+1) &= \phi_{ij;\delta}^E(t) + \epsilon \sum_{k=0}^{n_f} dW_{ik}(t) \psi_{kj}(\mathbf{x}_\delta).\end{aligned}$$

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) + \sum_j dW_{ij}(t) \phi_{ij;\delta}^E(t) + \frac{\epsilon}{2} \sum_{j_1, j_2} dW_{ij_1}(t) dW_{ij_2}(t) \psi_{j_1 j_2}(x_\delta) \end{aligned}$$

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ &= z_{i;\delta}(t) + \sum_j dW_{ij}(t) \phi_{ij;\delta}^E(t) + \frac{\epsilon}{2} \sum_{j_1, j_2} dW_{ij_1}(t) dW_{ij_2}(t) \psi_{j_1 j_2}(x_\delta) \\ &= z_{i;\delta}(t) + \sum_j \left[-\eta \sum_{\tilde{\alpha}} \phi_{ij;\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) \right] \phi_{ij;\delta}^E(t) \\ &+ \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} \left[-\eta \sum_{\tilde{\alpha}_1} \phi_{ij_1;\tilde{\alpha}_1}^E(t) \epsilon_{i;\tilde{\alpha}_1}(t) \right] \left[-\eta \sum_{\tilde{\alpha}_2} \phi_{ij_2;\tilde{\alpha}_2}^E(t) \epsilon_{i;\tilde{\alpha}_2}(t) \right] \psi_{j_1 j_2}(x_\delta) \end{aligned}$$

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ &= z_{i;\delta}(t) + \sum_j dW_{ij}(t) \phi_{ij;\delta}^E(t) + \frac{\epsilon}{2} \sum_{j_1, j_2} dW_{ij_1}(t) dW_{ij_2}(t) \psi_{j_1 j_2}(x_\delta) \\ &= z_{i;\delta}(t) + \sum_j \left[-\eta \sum_{\tilde{\alpha}} \phi_{ij;\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) \right] \phi_{ij;\delta}^E(t) \\ &\quad + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} \left[-\eta \sum_{\tilde{\alpha}_1} \phi_{ij_1;\tilde{\alpha}_1}^E(t) \epsilon_{i;\tilde{\alpha}_1}(t) \right] \left[-\eta \sum_{\tilde{\alpha}_2} \phi_{ij_2;\tilde{\alpha}_2}^E(t) \epsilon_{i;\tilde{\alpha}_2}(t) \right] \psi_{j_1 j_2}(x_\delta) \\ &= z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_j \phi_{ij;\delta}^E(t) \phi_{ij;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ &\quad + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left[\sum_{j_1, j_2} \epsilon \psi_{j_1 j_2}(x_\delta) \phi_{ij_1;\tilde{\alpha}_1}^E(t) \phi_{ij_2;\tilde{\alpha}_2}^E(t) \right] \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) \end{aligned}$$

Aside: Effective Kernel

To better understand this from the dual sample-space picture, let's analogously define an **effective kernel**

$$k_{ii;\delta_1\delta_2}^E(\theta) \equiv \sum_{j=0}^{n_f} \phi_{ij}^E(x_{\delta_1}; \theta) \phi_{ij}^E(x_{\delta_2}; \theta),$$

which measures a parameter-dependent similarity between two inputs x_{δ_1} and x_{δ_2} using our *effective features* $\phi_{ij}^E(x_\delta; \theta)$.

Aside 2: Meta Kernel

This last line suggests that an important object worth defining is

$$\begin{aligned}\mu_{\delta_0 \delta_1 \delta_2} &\equiv \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{j_1}(\mathbf{x}_{\delta_1}) \phi_{j_2}(\mathbf{x}_{\delta_2}) \\ &= \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{ij_1}^E(\mathbf{x}_{\delta_1}; \theta) \phi_{ij_2}^E(\mathbf{x}_{\delta_2}; \theta) + O(\epsilon^2)\end{aligned}$$

which we will call the **meta kernel**.

Aside 2: Meta Kernel

This last line suggests that an important object worth defining is

$$\begin{aligned}\mu_{\delta_0\delta_1\delta_2} &\equiv \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{j_1}(\mathbf{x}_{\delta_1}) \phi_{j_2}(\mathbf{x}_{\delta_2}) \\ &= \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{ij_1}^E(\mathbf{x}_{\delta_1}; \theta) \phi_{ij_2}^E(\mathbf{x}_{\delta_2}; \theta) + O(\epsilon^2)\end{aligned}$$

which we will call the **meta kernel**.

- ▶ This is a *parameter-independent* tensor given entirely in terms of the fixed $\phi_j(\mathbf{x})$ and $\psi_{j_1 j_2}(\mathbf{x})$ that define the model.

Aside 2: Meta Kernel

This last line suggests that an important object worth defining is

$$\begin{aligned}\mu_{\delta_0\delta_1\delta_2} &\equiv \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{j_1}(\mathbf{x}_{\delta_1}) \phi_{j_2}(\mathbf{x}_{\delta_2}) \\ &= \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{ij_1}^E(\mathbf{x}_{\delta_1}; \theta) \phi_{ij_2}^E(\mathbf{x}_{\delta_2}; \theta) + O(\epsilon^2)\end{aligned}$$

which we will call the **meta kernel**.

- ▶ This is a *parameter-independent* tensor given entirely in terms of the fixed $\phi_j(\mathbf{x})$ and $\psi_{j_1 j_2}(\mathbf{x})$ that define the model.
- ▶ For a fixed input \mathbf{x}_{δ_0} , $\mu_{\delta_0\delta_1\delta_2}$ computes a different feature-space inner product between the two inputs, \mathbf{x}_{δ_1} & \mathbf{x}_{δ_2} .

Aside 2: Meta Kernel

This last line suggests that an important object worth defining is

$$\begin{aligned}\mu_{\delta_0\delta_1\delta_2} &\equiv \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{j_1}(\mathbf{x}_{\delta_1}) \phi_{j_2}(\mathbf{x}_{\delta_2}) \\ &= \sum_{j_1, j_2=0}^{n_f} \epsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{ij_1}^E(\mathbf{x}_{\delta_1}; \theta) \phi_{ij_2}^E(\mathbf{x}_{\delta_2}; \theta) + O(\epsilon^2)\end{aligned}$$

which we will call the **meta kernel**.

- ▶ This is a *parameter-independent* tensor given entirely in terms of the fixed $\phi_j(\mathbf{x})$ and $\psi_{j_1 j_2}(\mathbf{x})$ that define the model.
- ▶ For a fixed input \mathbf{x}_{δ_0} , $\mu_{\delta_0\delta_1\delta_2}$ computes a different feature-space inner product between the two inputs, \mathbf{x}_{δ_1} & \mathbf{x}_{δ_2} .
- ▶ Due to the inclusion of ϵ into the definition of $\mu_{\delta_0\delta_1\delta_2}$, we should think of it as being parametrically small too.

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_j \phi_{ij;\delta}^E(t) \phi_{ij;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ & + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left[\epsilon \sum_{j_1, j_2} \phi_{ij_1; \tilde{\alpha}_1}^E(t) \phi_{ij_2; \tilde{\alpha}_2}^E(t) \psi_{j_1 j_2}(x_\delta) \right] \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) \end{aligned}$$

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_j \phi_{ij;\delta}^E(t) \phi_{ij;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ & + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left[\epsilon \sum_{j_1, j_2} \phi_{ij_1; \tilde{\alpha}_1}^E(t) \phi_{ij_2; \tilde{\alpha}_2}^E(t) \psi_{j_1 j_2}(x_\delta) \right] \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{ii;\delta\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) + O(\epsilon^2) \end{aligned}$$

Model Prediction Dynamics

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_j \phi_{ij;\delta}^E(t) \phi_{ij;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ & + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left[\epsilon \sum_{j_1, j_2} \phi_{ij_1; \tilde{\alpha}_1}^E(t) \phi_{ij_2; \tilde{\alpha}_2}^E(t) \psi_{j_1 j_2}(\mathbf{x}_\delta) \right] \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{ii;\delta\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) + O(\epsilon^2) \end{aligned}$$

This is a coupled nonlinear difference equation...

Effective Kernel Dynamics

$$\phi_{ij;\delta}^{\mathbb{E}}(t+1) = \phi_{ij;\delta}^{\mathbb{E}}(t) + \epsilon \sum_{k=0}^{n_f} dW_{ik}(t) \psi_{kj}(x_{\delta})$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) + \epsilon \sum_{k=0}^{n_f} \left[-\eta \sum_{\tilde{\alpha}} \phi_{ik;\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) \right] \psi_{kj}(x_\delta)$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_\delta) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_\delta) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

\implies

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_{\delta}) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

\Rightarrow

$$\begin{aligned} & \sum_j \phi_{ij;\delta_1}^E(t+1) \phi_{ij;\delta_2}^E(t+1) \\ &= \sum_j \phi_{ij;\delta_1}^E(t) \phi_{ij;\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{j,k} \epsilon \psi_{kj}(x_{\delta_1}) \phi_{ij;\delta_2}^E(t) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ & \quad - \eta \sum_{\tilde{\alpha}} \left[\sum_{j,k} \epsilon \psi_{kj}(x_{\delta_2}) \phi_{ij;\delta_1}^E(t) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2) \end{aligned}$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_{\delta}) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

\implies

$$\begin{aligned} & \sum_j \phi_{ij;\delta_1}^E(t+1) \phi_{ij;\delta_2}^E(t+1) \\ &= \sum_j \phi_{ij;\delta_1}^E(t) \phi_{ij;\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{j,k} \epsilon \psi_{kj}(x_{\delta_1}) \phi_j(x_{\delta_2}) \phi_k(x_{\tilde{\alpha}}) \right] \epsilon_{i;\tilde{\alpha}}(t) \\ & \quad - \eta \sum_{\tilde{\alpha}} \left[\sum_{j,k} \epsilon \psi_{kj}(x_{\delta_2}) \phi_j(x_{\delta_1}) \phi_k(x_{\tilde{\alpha}}) \right] \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2) \end{aligned}$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_\delta) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

\implies

$$k_{ii;\delta_1\delta_2}^E(t+1) = k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2)$$

Effective Kernel Dynamics

$$\phi_{ij;\delta}^E(t+1) = \phi_{ij;\delta}^E(t) - \eta \sum_{\tilde{\alpha}} \left[\sum_{k=0}^{n_f} \epsilon \psi_{kj}(x_\delta) \phi_{ik;\tilde{\alpha}}^E(t) \right] \epsilon_{i;\tilde{\alpha}}(t)$$

\implies

$$k_{ii;\delta_1\delta_2}^E(t+1) = k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2)$$

Linear difference equation, w/ $\mu_{\delta_1\delta_2\tilde{\alpha}}$ playing the role of $k_{\delta\tilde{\alpha}}$...

Quadratic Model Dynamics: Dual Sample Space

The *model predictions* will update as

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{ii;\delta\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) + O(\epsilon^2) \end{aligned}$$

While the *effective kernel* will update as

$$k_{ii;\delta_1\delta_2}^E(t+1) = k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} (\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}}) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2)$$

Quadratic Model Dynamics: Dual Sample Space

The *model predictions* will update as

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{ii;\delta\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) + O(\epsilon^2) \end{aligned}$$

While the *effective kernel* will update as

$$k_{ii;\delta_1\delta_2}^E(t+1) = k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2)$$

- ▶ These joint updates are coupled *difference equations*, and the first is *nonlinear* in the training error.

Quadratic Model Dynamics: Dual Sample Space

The *model predictions* will update as

$$\begin{aligned} & z_{i;\delta}(t+1) \\ = & z_{i;\delta}(t) - \eta \sum_{\tilde{\alpha}} k_{ii;\delta\tilde{\alpha}}^E(t) \epsilon_{i;\tilde{\alpha}}(t) + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}(t) \epsilon_{i;\tilde{\alpha}_2}(t) + O(\epsilon^2) \end{aligned}$$

While the *effective kernel* will update as

$$k_{ii;\delta_1\delta_2}^E(t+1) = k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2)$$

- ▶ These joint updates are coupled *difference equations*, and the first is *nonlinear* in the training error.
- ▶ We are now going to solve these equations in a closed form to leading order in ϵ using **perturbation theory**.

Another Theoretical Minimum

Decompose the model prediction into *free* and *interacting* parts:

$$z_{i;\delta}(t) \equiv z_{i;\delta}^F(t) + z_{i;\delta}^I(t).$$

Another Theoretical Minimum

Decompose the model prediction into *free* and *interacting* parts:

$$z_{i;\delta}(t) \equiv z_{i;\delta}^F(t) + z_{i;\delta}^I(t).$$

- ▶ $z_{i;\delta}^F(t)$ solves the linear model dynamics from before:

$$\begin{aligned} z_{i;\delta}^F(t+1) &= z_{i;\delta}^F(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}^F(t) \\ &= z_{i;\delta}^F(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \left[z_{i;\tilde{\alpha}}^F(t) - y_{i;\tilde{\alpha}} \right]. \end{aligned}$$

Another Theoretical Minimum

Decompose the model prediction into *free* and *interacting* parts:

$$z_{i;\delta}(t) \equiv z_{i;\delta}^F(t) + z_{i;\delta}^I(t).$$

- ▶ $z_{i;\delta}^F(t)$ solves the linear model dynamics from before:

$$\begin{aligned} z_{i;\delta}^F(t+1) &= z_{i;\delta}^F(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \epsilon_{i;\tilde{\alpha}}^F(t) \\ &= z_{i;\delta}^F(t) - \eta \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \left[z_{i;\tilde{\alpha}}^F(t) - y_{i;\tilde{\alpha}} \right]. \end{aligned}$$

- ▶ $z_{i;\delta}^I(t) = O(\epsilon)$, since $\epsilon \rightarrow 0$ gives back the linear model.

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t+1) \\ &= k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \epsilon_{i;\tilde{\alpha}}(t) + O(\epsilon^2) \end{aligned}$$

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t+1) \\ &= k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \left[z_{i;\tilde{\alpha}}^F(t) + z_{i;\tilde{\alpha}}^I(t) - y_{i;\tilde{\alpha}} \right] + O(\epsilon^2) \end{aligned}$$

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t+1) \\ &= k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \left[z_{i;\tilde{\alpha}}^F(t) + z_{i;\tilde{\alpha}}^I(t) - y_{i;\tilde{\alpha}} \right] \end{aligned}$$

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t+1) \\ &= k_{ii;\delta_1\delta_2}^E(t) - \eta \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \left[z_{i;\tilde{\alpha}}^F(t) - y_{i;\tilde{\alpha}} \right] \end{aligned}$$

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t) \\ &= k_{ii;\delta_1\delta_2}^E(0) - \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) \left\{ \eta \sum_{s=0}^{t-1} \left[z_{i;\tilde{\alpha}}^F(s) - y_{i;\tilde{\alpha}} \right] \right\} \end{aligned}$$

Another Theoretical Minimum

Let's solve the effective kernel equation first:

$$\begin{aligned} & k_{ii;\delta_1\delta_2}^E(t) \\ &= k_{ii;\delta_1\delta_2}^E(0) - \sum_{\tilde{\alpha}} \left(\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}} \right) a_{i;\tilde{\alpha}}(t) \end{aligned}$$

“ $a_{i;\tilde{\alpha}}$ ” Dynamical Helper Function

The explicit representation of $a_{i;\tilde{\alpha}}(t)$ is given by

$$\begin{aligned} & a_{i;\tilde{\alpha}}(t) \\ & \equiv \eta \sum_{s=0}^{t-1} \left[z_{i;\tilde{\alpha}}^F(s) - y_{i;\tilde{\alpha}} \right] \\ & = \eta \sum_{\tilde{\alpha}_1} \left\{ \sum_{s=0}^{t-1} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}\tilde{\alpha}_1} \right\} (z_{i;\tilde{\alpha}_1} - y_{i;\tilde{\alpha}_1}) \\ & = \eta \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left\{ \left[\delta - (I - \eta \tilde{k}) \right]^{-1} \right\}^{\tilde{\alpha}\tilde{\alpha}_2} \left[\delta - (\delta - \eta \tilde{k})^t \right]_{\tilde{\alpha}_2\tilde{\alpha}_1} (z_{i;\tilde{\alpha}_1} - y_{i;\tilde{\alpha}_1}) \\ & = \sum_{\tilde{\alpha}_2} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_2} \left\{ z_{i;\tilde{\alpha}_2} - y_{i;\tilde{\alpha}_2} - \left[z_{i;\tilde{\alpha}_2}^F(t) - y_{i;\tilde{\alpha}_2} \right] \right\} . \end{aligned}$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^l(t)$:

$$z_{i;\delta}^l(t+1) = z_{i;\delta}^l(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^l(t) + \eta \mathbb{F}_{i;\delta}(t).$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^1(t)$:

$$z_{i;\delta}^1(t+1) = z_{i;\delta}^1(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^1(t) + \eta \mathbb{F}_{i;\delta}(t).$$

- ▶ We used the fact that the *free part* satisfies the free dynamics.

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^l(t)$:

$$z_{i;\delta}^l(t+1) = z_{i;\delta}^l(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^l(t) + \eta \mathbb{F}_{i;\delta}(t).$$

- ▶ We used the fact that the *free part* satisfies the free dynamics.
- ▶ We've defined a *damping force*:

$$\begin{aligned} \mathbb{F}_{i;\delta}(t) \equiv & - \sum_{\tilde{\alpha}} \left[k_{\delta\tilde{\alpha}}^E(t) - k_{\delta\tilde{\alpha}} \right] \epsilon_{i;\tilde{\alpha}}^F(t) \\ & + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \end{aligned}$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^l(t)$:

$$z_{i;\delta}^l(t+1) = z_{i;\delta}^l(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^l(t) + \eta \mathbb{F}_{i;\delta}(t).$$

- ▶ We used the fact that the *free part* satisfies the free dynamics.
- ▶ We've defined a *damping force*:

$$\begin{aligned} \mathbb{F}_{i;\delta}(t) \equiv & - \sum_{\tilde{\alpha}} \left[k_{\delta\tilde{\alpha}}^E(t) - k_{\delta\tilde{\alpha}} \right] \epsilon_{i;\tilde{\alpha}}^F(t) \\ & + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \end{aligned}$$

- ▶ Importantly, we've explicit expressions for all in $\mathbb{F}_{i;\delta}(t)$.

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^l(t)$:

$$z_{i;\delta}^l(t+1) = z_{i;\delta}^l(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^l(t) + \eta \mathbb{F}_{i;\delta}(t).$$

Restricting to $\tilde{\alpha} \in \mathcal{A}$, we get a first-order inhomogeneous linear difference equation, with the damping force ruining the homogeneity:

$$z_{i;\tilde{\alpha}}^l(t+1) = \sum_{\tilde{\alpha}_1} (\delta_{\tilde{\alpha}\tilde{\alpha}_1} - \eta k_{\tilde{\alpha}\tilde{\alpha}_1}) z_{i;\tilde{\alpha}_1}^l(t) + \eta \mathbb{F}_{i;\tilde{\alpha}}(t).$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^!(t)$:

$$z_{i;\delta}^!(t+1) = z_{i;\delta}^!(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^!(t) + \eta \mathbb{F}_{i;\delta}(t).$$

Restricting to $\tilde{\alpha} \in \mathcal{A}$, we get a first-order inhomogeneous linear difference equation, with the damping force ruining the homogeneity:

$$z_{i;\tilde{\alpha}}^!(t+1) = \sum_{\tilde{\alpha}_1} (\delta_{\tilde{\alpha}\tilde{\alpha}_1} - \eta k_{\tilde{\alpha}\tilde{\alpha}_1}) z_{i;\tilde{\alpha}_1}^!(t) + \eta \mathbb{F}_{i;\tilde{\alpha}}(t).$$

We can formally solve this by convoluting $\mathbb{F}_{i;\tilde{\alpha}}(t)$ with $U(t)$:

$$z_{i;\tilde{\alpha}}^!(t) = \eta \sum_{s=0}^{t-1} \sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(t-1-s) \mathbb{F}_{i;\tilde{\alpha}_1}(s).$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^1(t)$:

$$z_{i;\delta}^1(t+1) = z_{i;\delta}^1(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^1(t) + \eta \mathbb{F}_{i;\delta}(t).$$

For general inputs $\delta \in \mathcal{D}$, we can iterate to find:

$$z_{i;\delta}^1(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} z_{i;\tilde{\alpha}}^1(s) \right].$$

Another Theoretical Minimum

Let's solve the *interacting part* of the model prediction, $z_{i;\delta}^1(t)$:

$$z_{i;\delta}^1(t+1) = z_{i;\delta}^1(t) - \sum_{j,\tilde{\alpha}} \eta k_{\delta\tilde{\alpha}} z_{j;\tilde{\alpha}}^1(t) + \eta \mathbb{F}_{i;\delta}(t).$$

For general inputs $\delta \in \mathcal{D}$, we can iterate to find:

$$z_{i;\delta}^1(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} z_{i;\tilde{\alpha}}^1(s) \right].$$

After some manipulations of this sum to be done in the privacy of your own notebook, we get a slightly less formal solution:

$$z_{i;\delta}^1(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right] + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^1(t).$$

Aside: Convergence on Training Set

Ultimately, we want $z_{i;\delta}^!(t)$ at the end of training, $t \rightarrow \infty$.

Aside: Convergence on Training Set

Ultimately, we want $z_{i;\delta}^!(t)$ at the end of training, $t \rightarrow \infty$.

- ▶ Before we showed that for small enough η we have:

$$\lim_{t \rightarrow \infty} U(t) \propto \exp(-\eta \tilde{k} t) .$$

Aside: Convergence on Training Set

Ultimately, we want $z_{i;\delta}^l(t)$ at the end of training, $t \rightarrow \infty$.

- ▶ Before we showed that for small enough η we have:

$$\lim_{t \rightarrow \infty} U(t) \propto \exp(-\eta \tilde{k} t) .$$

Then, for training inputs, $z_{i;\tilde{\alpha}}^l(t)$, will exponentially converge

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^l(t) = 0 .$$

Aside: Convergence on Training Set

Ultimately, we want $z_{i;\delta}^l(t)$ at the end of training, $t \rightarrow \infty$.

- ▶ Before we showed that for small enough η we have:

$$\lim_{t \rightarrow \infty} U(t) \propto \exp(-\eta \tilde{k} t) .$$

Then, for training inputs, $z_{i;\tilde{\alpha}}^l(t)$, will exponentially converge

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^l(t) = 0 .$$

- ▶ To see why this holds, note that we have

$$\lim_{s \rightarrow \infty} \mathbb{F}(s) \propto \exp(-\eta \tilde{k} s) ,$$

since its leading behavior is linearly proportional to $\epsilon_{j;\tilde{\alpha}}^F(t)$.

Aside: Convergence on Training Set

Ultimately, we want $z_{i;\delta}^l(t)$ at the end of training, $t \rightarrow \infty$.

- ▶ Before we showed that for small enough η we have:

$$\lim_{t \rightarrow \infty} U(t) \propto \exp(-\eta \tilde{k} t) .$$

Then, for training inputs, $z_{i;\tilde{\alpha}}^l(t)$, will exponentially converge

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^l(t) = 0 .$$

- ▶ To see why this holds, note that we have

$$\lim_{s \rightarrow \infty} \mathbb{F}(s) \propto \exp(-\eta \tilde{k} s) ,$$

since its leading behavior is linearly proportional to $\epsilon_{j;\tilde{\alpha}}^F(t)$.

All together, this means that $z_{i;\tilde{\alpha}}^l(t)$ converges as

$$\lim_{t \rightarrow \infty} z_{i;\tilde{\alpha}}^l(t) \propto \lim_{t \rightarrow \infty} \left\{ \eta t \exp[-(t-1)\eta \tilde{k}] \right\} = 0 ,$$

slightly slower than the free solution $\epsilon_{i;\tilde{\alpha}}^F(t) \propto \exp(-\tilde{k}t)$.

Another Theoretical Minimum

$$z_{i;\delta}^l(t) = \eta \sum_{s=0}^{t-1} \left[\mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right] + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^l(t).$$

Another Theoretical Minimum

$$\lim_{t \rightarrow \infty} z_{i;\delta}^l(t) = \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right].$$

Another Theoretical Minimum

$$\lim_{t \rightarrow \infty} z_{i;\delta}^l(t) = \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right].$$

So we still have one more infinite sum to perform ...

Another Theoretical Minimum

$$\lim_{t \rightarrow \infty} z_{i;\delta}^l(t) = \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right].$$

So we still have one more infinite sum to perform ...

$$\begin{aligned} \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) &= - \sum_{\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[k_{\delta\tilde{\alpha}}^E(t) - k_{\delta\tilde{\alpha}} \right] \epsilon_{i;\tilde{\alpha}}^F(t) \right\} \\ &\quad + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \right\} \end{aligned}$$

Another Theoretical Minimum

$$\lim_{t \rightarrow \infty} z_{i;\delta}^l(t) = \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left[\eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right].$$

So we still have one more infinite sum to perform ...

$$\begin{aligned} \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) &= \sum_{\tilde{\alpha}_0, \tilde{\alpha}} \left(\mu_{\delta_1 \delta_2 \tilde{\alpha}_0} + \mu_{\delta_2 \delta_1 \tilde{\alpha}_0} \right) \left\{ \eta \sum_{s=0}^{\infty} a_{i;\tilde{\alpha}_0}(s) \epsilon_{i;\tilde{\alpha}}^F(s) \right\} \\ &+ \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2} \left\{ \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}_1}^F(s) \epsilon_{i;\tilde{\alpha}_2}^F(s) \right\} \end{aligned}$$

Another *Aside*: Geometric Sums

$$\eta \sum_{s=0}^{\infty} \epsilon_{i; \tilde{\alpha}_1}^F(t) \epsilon_{i; \tilde{\alpha}_2}^F(t)$$

Another *Aside*: Geometric Sums

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \\ &= \eta \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} \sum_{s=0}^{\infty} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_1 \tilde{\alpha}_3} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{i;\tilde{\alpha}_3} - y_{i;\tilde{\alpha}_3}) (z_{i;\tilde{\alpha}_4} - y_{i;\tilde{\alpha}_4}) \end{aligned}$$

Another *Aside*: Geometric Sums

$$\begin{aligned} & \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \\ &= \eta \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} \sum_{s=0}^{\infty} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_1 \tilde{\alpha}_3} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{i;\tilde{\alpha}_3} - y_{i;\tilde{\alpha}_3}) (z_{i;\tilde{\alpha}_4} - y_{i;\tilde{\alpha}_4}) \\ &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{i;\tilde{\alpha}_3} - y_{i;\tilde{\alpha}_3}) (z_{i;\tilde{\alpha}_4} - y_{i;\tilde{\alpha}_4}) \end{aligned}$$

Another *Aside*: Geometric Sums

$$\begin{aligned}
 & \eta \sum_{s=0}^{\infty} \epsilon_{i;\tilde{\alpha}_1}^F(t) \epsilon_{i;\tilde{\alpha}_2}^F(t) \\
 &= \eta \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} \sum_{s=0}^{\infty} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_1 \tilde{\alpha}_3} \left[(\delta - \eta \tilde{k})^s \right]_{\tilde{\alpha}_2 \tilde{\alpha}_4} (z_{i;\tilde{\alpha}_3} - y_{i;\tilde{\alpha}_3}) (z_{i;\tilde{\alpha}_4} - y_{i;\tilde{\alpha}_4}) \\
 &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} (z_{i;\tilde{\alpha}_3} - y_{i;\tilde{\alpha}_3}) (z_{i;\tilde{\alpha}_4} - y_{i;\tilde{\alpha}_4})
 \end{aligned}$$

This last operation yielded an **inverting tensor** implicitly defined:

$$\begin{aligned}
 & \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2} \\
 &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \frac{1}{\eta} \left[\delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - (\delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} - \eta \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5}) (\delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6}) \right] \\
 &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(\tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right).
 \end{aligned}$$

Another Theoretical Minimum

$$\begin{aligned}
 & z_{i;\dot{\beta}}(\infty) \\
 = & z_{i;\dot{\beta}}(0) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} k_{\dot{\beta}\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \epsilon_{i;\tilde{\alpha}_2}^F(0) \\
 & + \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4 \in \mathcal{A}} \left[\mu_{\tilde{\alpha}_1\dot{\beta}\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} k_{\dot{\beta}\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_5\tilde{\alpha}_6} \mu_{\tilde{\alpha}_1\tilde{\alpha}_6\tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} \epsilon_{i;\tilde{\alpha}_3}^F(0) \epsilon_{i;\tilde{\alpha}_4}^F(0) \\
 & + \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4 \in \mathcal{A}} \left[\mu_{\dot{\beta}\tilde{\alpha}_1\tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} k_{\dot{\beta}\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_5\tilde{\alpha}_6} \mu_{\tilde{\alpha}_6\tilde{\alpha}_1\tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} \epsilon_{i;\tilde{\alpha}_3}^F(0) \epsilon_{i;\tilde{\alpha}_4}^F(0)
 \end{aligned}$$

where the **algorithm projectors** are given by

$$\begin{aligned}
 Z_A^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} & \equiv \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2\tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2\tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1\tilde{\alpha}_5\tilde{\alpha}_3\tilde{\alpha}_4}, \\
 Z_B^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4} & \equiv \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2\tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2\tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1\tilde{\alpha}_5\tilde{\alpha}_3\tilde{\alpha}_4} + \frac{\eta}{2} X_{\parallel}^{\tilde{\alpha}_1\tilde{\alpha}_2\tilde{\alpha}_3\tilde{\alpha}_4}.
 \end{aligned}$$

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

Unlike *kernel methods*, this depends on the *learning algorithm*.

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

Unlike *kernel methods*, this depends on the *learning algorithm*.

- ▶ If we'd optimized by *direct optimization*, we'd have found:

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = 0, \quad Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = \frac{1}{2} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4}.$$

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

Unlike *kernel methods*, this depends on the *learning algorithm*.

- ▶ If we'd optimized by *direct optimization*, we'd have found:

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = 0, \quad Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = \frac{1}{2} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4}.$$

- ▶ In the ODE limit, we get different predictions

$$Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4},$$

$$\sum_{\tilde{\alpha}_3, \tilde{\alpha}_4 \in \mathcal{A}} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left(\tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right) = \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2},$$

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

We again have two ways of thinking about the solution:

- ▶ we can use the optimal parameters to make predictions, or
- ▶ we can make *nearly-kernel predictions* in which the features, the meta features, and the model parameters do not appear.

Nearly-Kernel Methods

When the prediction of a *quadratic model* is computed in this way, call it a *nearly-kernel machine* or **nearly-kernel methods**.

We again have two ways of thinking about the solution:

- ▶ we can use the optimal parameters to make predictions, or
- ▶ we can make *nearly-kernel predictions* in which the features, the meta features, and the model parameters do not appear.

Predictions are made by direct comparison with the training set:

- ▶ It has the kernel linear piece $\propto y_{i;\tilde{\alpha}_2}$, and
- ▶ it also has a new quadratic piece $\propto y_{i;\tilde{\alpha}_1} y_{i;\tilde{\alpha}_2}$.

Algorithm Projectors

The prediction takes a **universal form**, with of the *algorithm dependence* of θ^* is encoded in the *algorithm projectors*.

Algorithm Projectors

The prediction takes a **universal form**, with of the *algorithm dependence* of θ^* is encoded in the *algorithm projectors*.

- ▶ Gives theoretical means of isolating the *inductive bias of the training algorithm* from all the other inductive biases.
- ▶ Essentially this is the sample-space *dual* of a learning algorithm, cf. the relationship between $\phi_j(x)$ and $k(x_{\delta_1}, x_{\delta_2})$.

Algorithm Projectors

The prediction takes a **universal form**, with of the *algorithm dependence* of θ^* is encoded in the *algorithm projectors*.

- ▶ Gives theoretical means of isolating the *inductive bias of the training algorithm* from all the other inductive biases.
- ▶ Essentially this is the sample-space *dual* of a learning algorithm, cf. the relationship between $\phi_j(x)$ and $k(x_{\delta_1}, x_{\delta_2})$.

In principle, this gives a means of **inverse algorithm design**:

- (i) engineer a desired functional form for the projectors
- (ii) determine the parameter-space training algorithm that leads to the projectors having those desired properties or specific functional form.

Representation Learning

For simplicity, let's pick the **direct optimization** solution:

$$k_{ii;\delta_1\delta_2}^E(\theta^*) = k_{\delta_1\delta_2} + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} (\mu_{\delta_1\delta_2\tilde{\alpha}_1} + \mu_{\delta_2\delta_1\tilde{\alpha}_1}) \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2} + \mathcal{O}(\epsilon^2).$$

Representation Learning

For simplicity, let's pick the **direct optimization** solution:

$$k_{ii;\delta_1\delta_2}^E(\theta^*) = k_{\delta_1\delta_2} + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} (\mu_{\delta_1\delta_2\tilde{\alpha}_1} + \mu_{\delta_2\delta_1\tilde{\alpha}_1}) \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2} + O(\epsilon^2).$$

Then, we can define a **trained kernel** that averages between the *fixed* kernel and *dynamical* effective kernels:

$$k_{ii;\delta_1\delta_2}^\# \equiv \frac{1}{2} \left[k_{\delta_1\delta_2} + k_{ii;\delta_1\delta_2}^E(\theta^*) \right].$$

Representation Learning

For simplicity, let's pick the **direct optimization** solution:

$$k_{ii;\delta_1\delta_2}^E(\theta^*) = k_{\delta_1\delta_2} + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} (\mu_{\delta_1\delta_2\tilde{\alpha}_1} + \mu_{\delta_2\delta_1\tilde{\alpha}_1}) \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2} + O(\epsilon^2).$$

Then, we can define a **trained kernel** that averages between the *fixed* kernel and *dynamical* effective kernels:

$$k_{ii;\delta_1\delta_2}^\# \equiv \frac{1}{2} \left[k_{\delta_1\delta_2} + k_{ii;\delta_1\delta_2}^E(\theta^*) \right].$$

Now the nearly-kernel prediction formula can be compressed,

$$z_i(x_{\tilde{\beta}}; \theta^*) = \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2 \in \mathcal{A}} k_{ii;\tilde{\beta}\tilde{\alpha}_1}^\# \tilde{k}_{ii}^{\tilde{\alpha}_1\tilde{\alpha}_2} y_{i;\tilde{\alpha}_2} + O(\epsilon^2),$$

taking the form of a *kernel prediction*, but with the benefit of nontrivial feature evolution incorporated into the trained kernel.

Representation Learning as Regularization

The *direct optimization* solution in parameter space is

$$z_i(x_{\beta}; \theta^*) = \sum_{j=0}^{n_f} W_{ij}^* \phi_j(x_{\beta}) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1}^* W_{ij_2}^* \psi_{j_1 j_2}(x_{\beta})$$

and the optimal parameters can decompose as

$$W_{ij}^* \equiv W_{ij}^F + W_{ij}^I,$$

where W_{ij}^F are the optimal parameters from the linear model.

Representation Learning as Regularization

The *direct optimization* solution in parameter space is

$$z_i(x_{\beta}; \theta^*) = \sum_{j=0}^{n_f} W_{ij}^* \phi_j(x_{\beta}) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1}^* W_{ij_2}^* \psi_{j_1 j_2}(x_{\beta})$$

and the optimal parameters can decompose as

$$W_{ij}^* \equiv W_{ij}^F + W_{ij}^I,$$

where W_{ij}^F are the optimal parameters from the linear model.

- ▶ The $O(\epsilon)$ tunings W_{ij}^I ruin the **fine tuning** of the W_{ij}^F , as they are constrained by the $\psi_{kj}(x)$ defined before training.

Representation Learning as Regularization

The *direct optimization* solution in parameter space is

$$z_i(x_{\hat{\beta}}; \theta^*) = \sum_{j=0}^{n_f} W_{ij}^* \phi_j(x_{\hat{\beta}}) + \frac{\epsilon}{2} \sum_{j_1, j_2=0}^{n_f} W_{ij_1}^* W_{ij_2}^* \psi_{j_1 j_2}(x_{\hat{\beta}})$$

and the optimal parameters can decompose as

$$W_{ij}^* \equiv W_{ij}^F + W_{ij}^I,$$

where W_{ij}^F are the optimal parameters from the linear model.

- ▶ The $O(\epsilon)$ tunings W_{ij}^I ruin the **fine tuning** of the W_{ij}^F , as they are constrained by the $\psi_{kj}(x)$ defined before training.
- ▶ Assuming these $\psi_{kj}(x)$ are *useful*, we might expect that the quadratic model will overfit less and generalize better.

Course Plan

~~Lecture 1~~ **Initialization, Linear Models**

- ▶ ~~§0 + §7.1 + §10.4~~

~~Lecture 2~~ **Quadratic Models & Nearly-Kernel Methods**

- ▶ ~~§11.4 (+ §7.2) + §∞.2.2~~

Lecture 3 The Principle of Sparsity (Recurring)

- ▶ §4, §8, §11.2, §∞.3

Lecture 4 The Principle of Criticality

- ▶ §5, §9, §11.3, §∞.1, +§10.3

Lecture 5 The End of Training & More

- ▶ §∞.2.3 + Maybe { §A, §B, ... }