

Bayesian Interpolation with Deep Linear Networks

Boris Hanin (Princeton ORFE) and Alexander Zlokapa (MIT Physics)

PNAS 2023 (arXiv:2212.14457)

Motivating Questions

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Q1. How do P, L, N_ℓ jointly affect learning?

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Q1. How do P, L, N_ℓ jointly affect learning?

- $L = 0$: $(x_i^T x_j)_{1 \leq i, j \leq P}$ depends on P/N_0 (Marchenko-Pastur)

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Q1. How do P, L, N_ℓ jointly affect learning?

- $L = 0$: $(x_i^T x_j)_{1 \leq i, j \leq P}$ depends on P/N_0 (Marchenko-Pastur)
- $P < \infty$: Deviation from linear model depends on

$$\frac{1}{N_1} + \dots + \frac{1}{N_L} \simeq \frac{L}{N}$$

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Q1. How do P, L, N_ℓ jointly affect learning?

Q2. How does one analyze learning in non-linear models?

Motivating Questions

Given

Model. $z(x, \theta) = W^{(L+1)} \sigma W^{(L)} \dots \sigma W^{(1)} x, \quad W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. $\mathcal{A} : (z, \mathcal{D}) \mapsto \theta_{\text{learned}}$

Q1. How do P, L, N_ℓ jointly affect learning?

Q2. How does one analyze learning in non-linear models?

Q3. How do models make predictions on test data?

This Work: Q1 - Q3 for Linear Networks

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. Bayesian inference at zero temperature:

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. Bayesian inference at zero temperature:

- $\theta \sim \mathbb{P}_{\text{prior}} \iff W_{ij}^{(\ell)} \sim \mathcal{N}(0, \sigma^2 / N_{\ell-1})$

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. Bayesian inference at zero temperature:

- $\theta \sim \mathbb{P}_{\text{prior}} \iff W_{ij}^{(\ell)} \sim \mathcal{N}(0, \sigma^2 / N_{\ell-1})$
- $-\log \mathbb{P}_{\beta}(\mathcal{D} \mid \theta) = \beta \sum_i (z(x_i; \theta) - y_i)^2$

This Work: Q1 - Q3 for Linear Networks

Model. Set $N_{L+1} = 1$ and take identity non-linearity

$$z^{(L+1)}(x, \theta) = W^{(L+1)} W^{(L)} \dots W^{(1)} x = \theta^T x$$

Data. $\mathcal{D} = \{(x_i, y_i), \quad i = 1, \dots, P\}$

Alg. Bayesian inference at zero temperature:

- $\theta \sim \mathbb{P}_{\text{prior}} \iff W_{ij}^{(\ell)} \sim \mathcal{N}(0, \sigma^2 / N_{\ell-1})$
- $-\log \mathbb{P}_{\beta}(\mathcal{D} \mid \theta) = \beta \sum_i (z(x_i; \theta) - y_i)^2$
- $\mathbb{P}_{\text{post}}(\theta \mid \mathcal{D}) = \lim_{\beta \rightarrow \infty} Z_{\beta}^{-1}(\mathcal{D}) \mathbb{P}_{\text{prior}}(\theta) \mathbb{P}_{\beta}(\mathcal{D} \mid \theta)$

Key Theorems, Informally

Key Theorems, Informally

T1. Bayesian inference is exactly solvable

Key Theorems, Informally

T1. Bayesian inference is exactly solvable

T2. Effective depth of prior:

$$\frac{1}{N_1} + \dots + \frac{1}{N_L} \simeq \frac{L}{N}$$

Effective depth of posterior:

$$P \left(\frac{1}{N_1} + \dots + \frac{1}{N_L} \right) \simeq \frac{PL}{N}$$

Key Theorems, Informally

T1. Bayesian inference is exactly solvable

T2. Effective depth of prior:

$$\frac{1}{N_1} + \cdots + \frac{1}{N_L} \simeq \frac{L}{N}$$

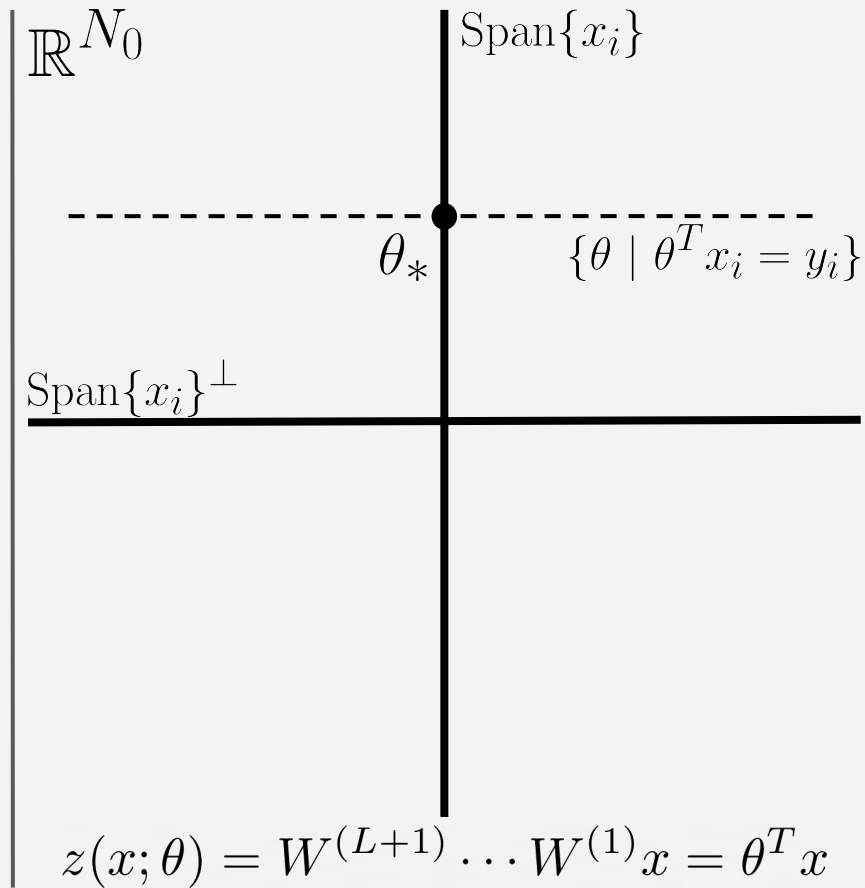
Effective depth of posterior:

$$P \left(\frac{1}{N_1} + \cdots + \frac{1}{N_L} \right) \simeq \frac{PL}{N}$$

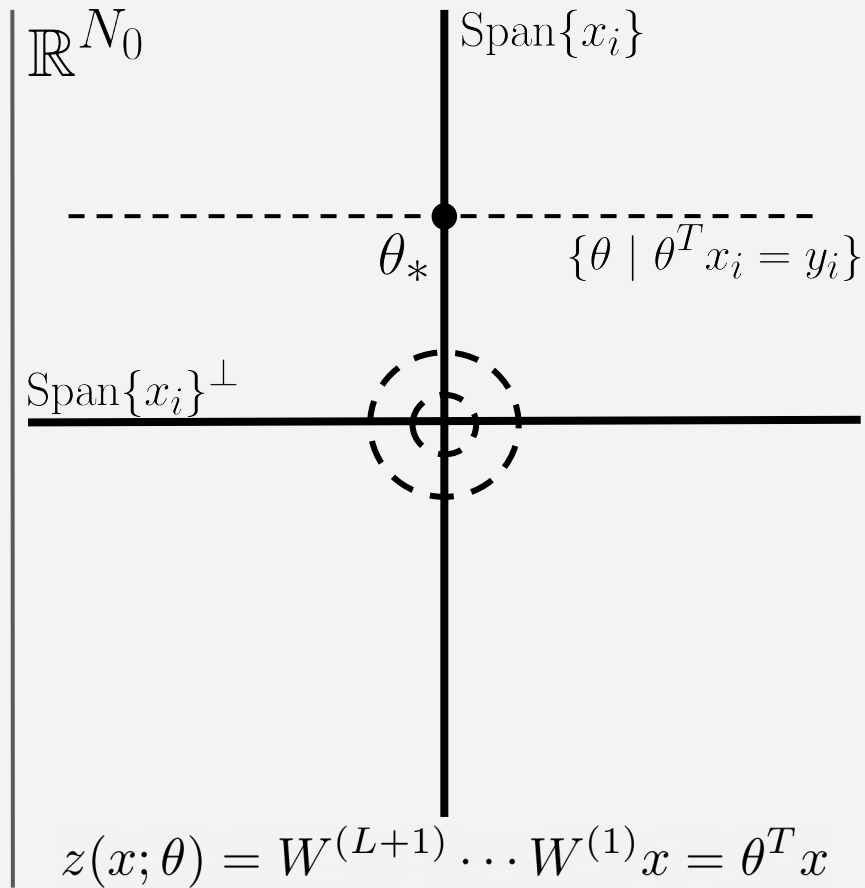
T3. Deep networks with universal priors learn same posteriors as shallow networks with optimal data-dependent priors

Structure of Prior/Posterior

Structure of Prior/Posterior



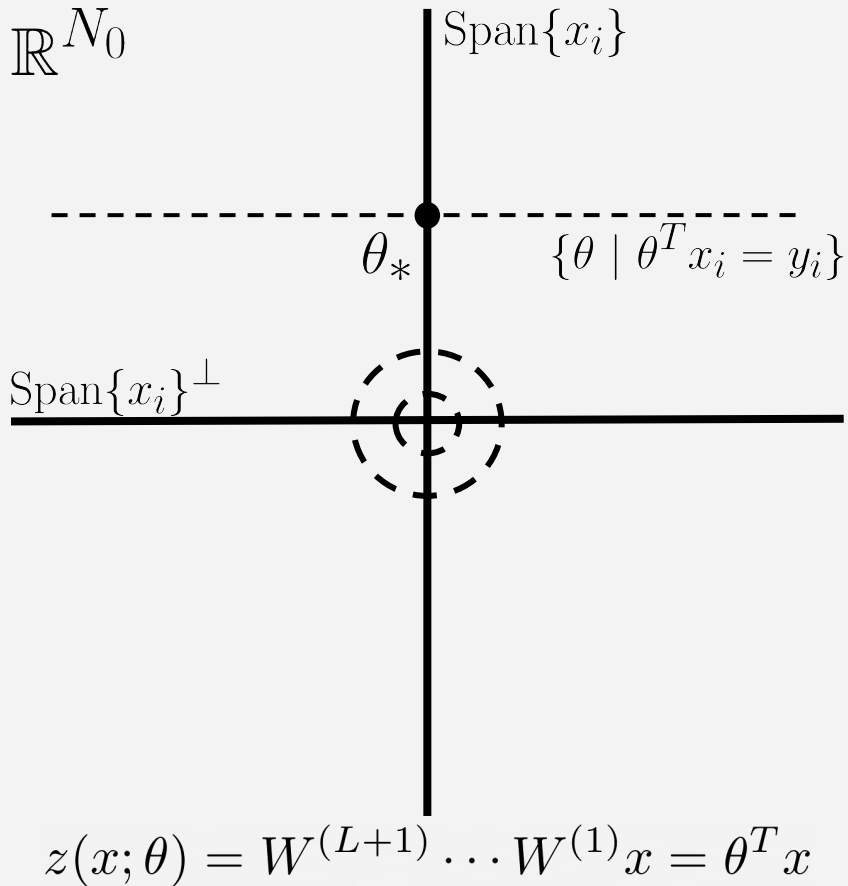
Structure of Prior/Posterior



Structure of Prior/Posterior

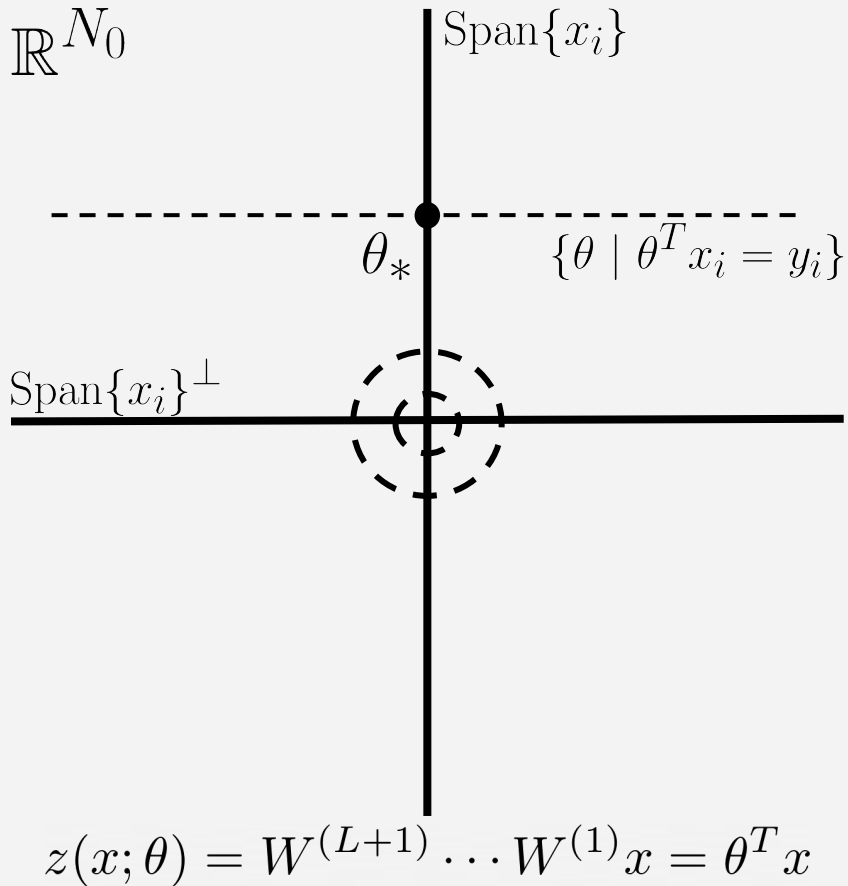
$$\boxed{1} \quad \mathbb{P}_{\text{post}}(\theta \mid L, N, \mathcal{D}, \sigma^2)$$

$$= \lim_{\beta \rightarrow \infty} \frac{\mathbb{P}_{\text{prior}}(\theta \mid L, N, \sigma^2) e^{-\beta \mathcal{L}_{\mathcal{D}}(\theta)}}{Z_{\beta}(\mathcal{D} \mid L, N, \sigma^2)}$$



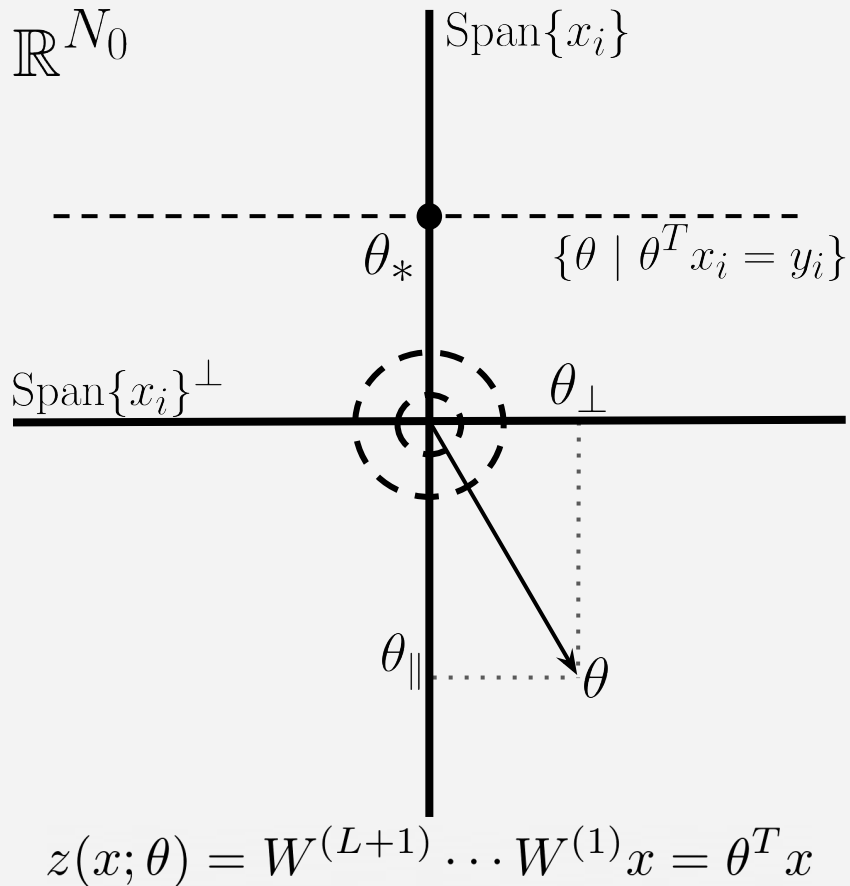
Structure of Prior/Posterior

$$\begin{aligned}
 \mathbf{1} \quad & \mathbb{P}_{\text{post}}(\theta \mid L, N, \mathcal{D}, \sigma^2) \\
 &= \lim_{\beta \rightarrow \infty} \frac{\mathbb{P}_{\text{prior}}(\theta \mid L, N, \sigma^2) e^{-\beta \mathcal{L}_{\mathcal{D}}(\theta)}}{Z_{\beta}(\mathcal{D} \mid L, N, \sigma^2)} \\
 &\propto \mathbb{P}_{\text{prior}}(\theta \mid \theta^T x_i = y_i)
 \end{aligned}$$



Structure of Prior/Posterior

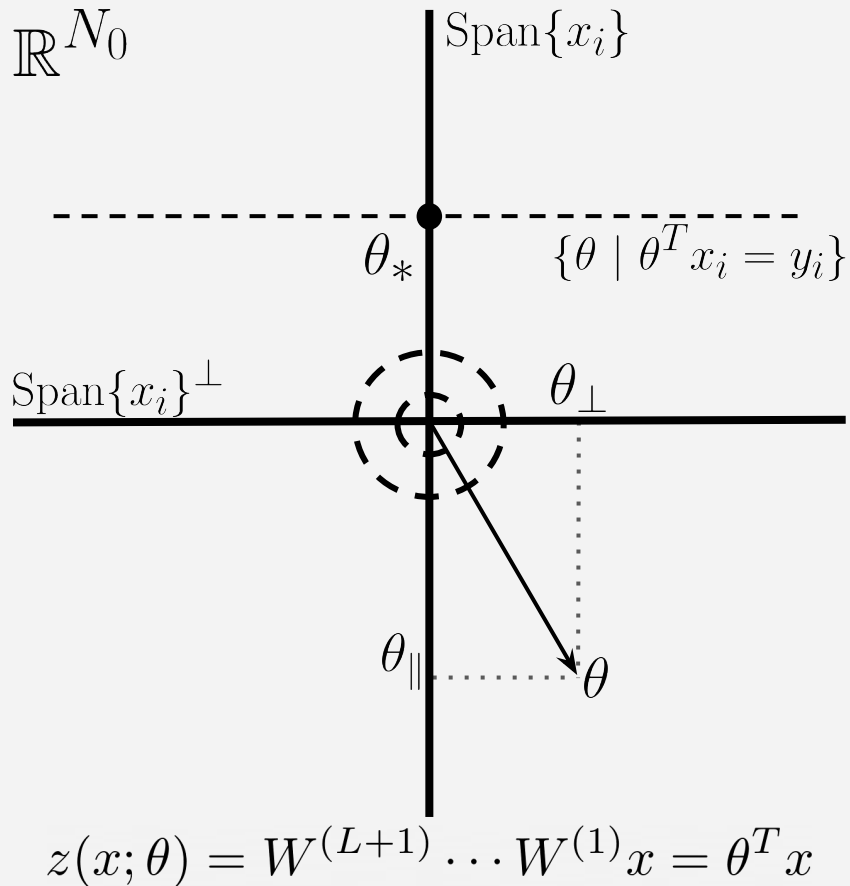
$$\begin{aligned}
 \mathbf{1} \quad & \mathbb{P}_{\text{post}}(\theta \mid L, N, \mathcal{D}, \sigma^2) \\
 &= \lim_{\beta \rightarrow \infty} \frac{\mathbb{P}_{\text{prior}}(\theta \mid L, N, \sigma^2) e^{-\beta \mathcal{L}_{\mathcal{D}}(\theta)}}{Z_{\beta}(\mathcal{D} \mid L, N, \sigma^2)} \\
 &\propto \mathbb{P}_{\text{prior}}(\theta \mid \theta^T x_i = y_i)
 \end{aligned}$$



Structure of Prior/Posterior

$$\begin{aligned}
 \mathbf{1} \quad & \mathbb{P}_{\text{post}}(\theta \mid L, N, \mathcal{D}, \sigma^2) \\
 &= \lim_{\beta \rightarrow \infty} \frac{\mathbb{P}_{\text{prior}}(\theta \mid L, N, \sigma^2) e^{-\beta \mathcal{L}_{\mathcal{D}}(\theta)}}{Z_{\beta}(\mathcal{D} \mid L, N, \sigma^2)} \\
 &\propto \mathbb{P}_{\text{prior}}(\theta \mid \theta^T x_i = y_i)
 \end{aligned}$$

$$\mathbf{2} \quad \Rightarrow \theta = \theta_* + \|\theta_{\perp}\| \cdot \text{unif}$$

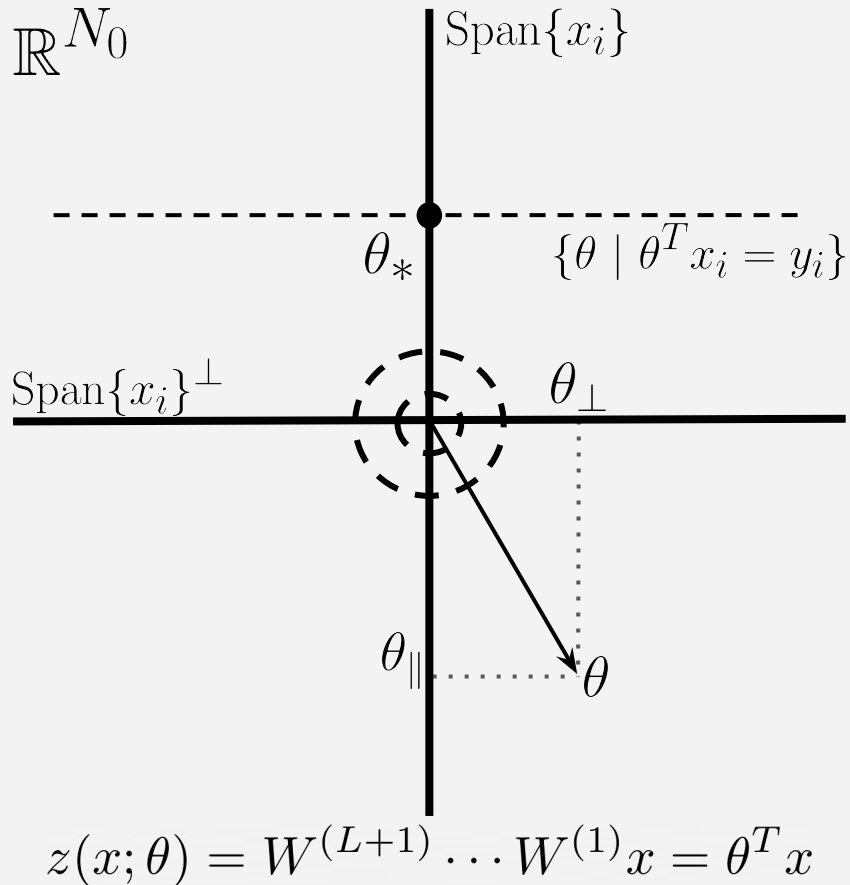


Structure of Prior/Posterior

$$\begin{aligned}
 \text{1} \quad & \mathbb{P}_{\text{post}}(\theta \mid L, N, \mathcal{D}, \sigma^2) \\
 &= \lim_{\beta \rightarrow \infty} \frac{\mathbb{P}_{\text{prior}}(\theta \mid L, N, \sigma^2) e^{-\beta \mathcal{L}_{\mathcal{D}}(\theta)}}{Z_{\beta}(\mathcal{D} \mid L, N, \sigma^2)} \\
 &\propto \mathbb{P}_{\text{prior}}(\theta \mid \theta^T x_i = y_i)
 \end{aligned}$$

$$\text{2} \quad \Rightarrow \theta = \theta_* + \|\theta_{\perp}\| \cdot \text{unif}$$

3 $\|\theta_{\perp}\|$ is a learnable data-dependent scale
 for predictions in new directions



Evidence and Posterior Via G-Functions

Thm. Write $Z(t) = \mathbb{E}_{\text{post}} [\exp \{-it \cdot \theta\}]$. Then,

$$Z(t) = \left(\frac{4\pi}{\|\theta_*\|^2} \right)^{\frac{P}{2}} \times e^{-i\langle \theta_*, t_{||} \rangle} \times \prod_{\ell=1}^L \Gamma \left(\frac{N_\ell}{2} \right)^{-1} \\ \times \sum_{k=0}^{\infty} \frac{(-\|t^\perp\|^2 M)^k}{k!} G_{L+1} \left(\frac{\|\theta_*\|^2}{4M} \left| \frac{P}{2}, \frac{N_1}{2} + k, \dots, \frac{N_L}{2} + k \right. \right)$$

where $M = \prod_{\ell=1}^L 2\sigma^2/N_\ell$ and

$$G_\ell(z \mid b_1, \dots, b_\ell) = \frac{1}{2\pi i} \int_{-i\infty}^{i\infty} ds \, z^s \prod_{j=1}^{\ell} \Gamma(1 + b_j - s)$$

Phase Diagram for Deep Linear Posteriors

Phase Diagram for Deep Linear Posteriors

Setup

$$\left\{ \begin{array}{l} \text{model: } W^{(L+1)} \dots W^{(1)} x \\ \text{prior: } W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1}) \\ \text{likelihood: MSE on } \mathcal{D} = \{(x_i, y_i)\}_{i=1}^P \\ \text{regime: } P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1 \\ \text{prior depth: } L/N = \gamma \\ \text{post. depth: } PL/N = \lambda \end{array} \right.$$

Phase Diagram for Deep Linear Posteriors

Setup

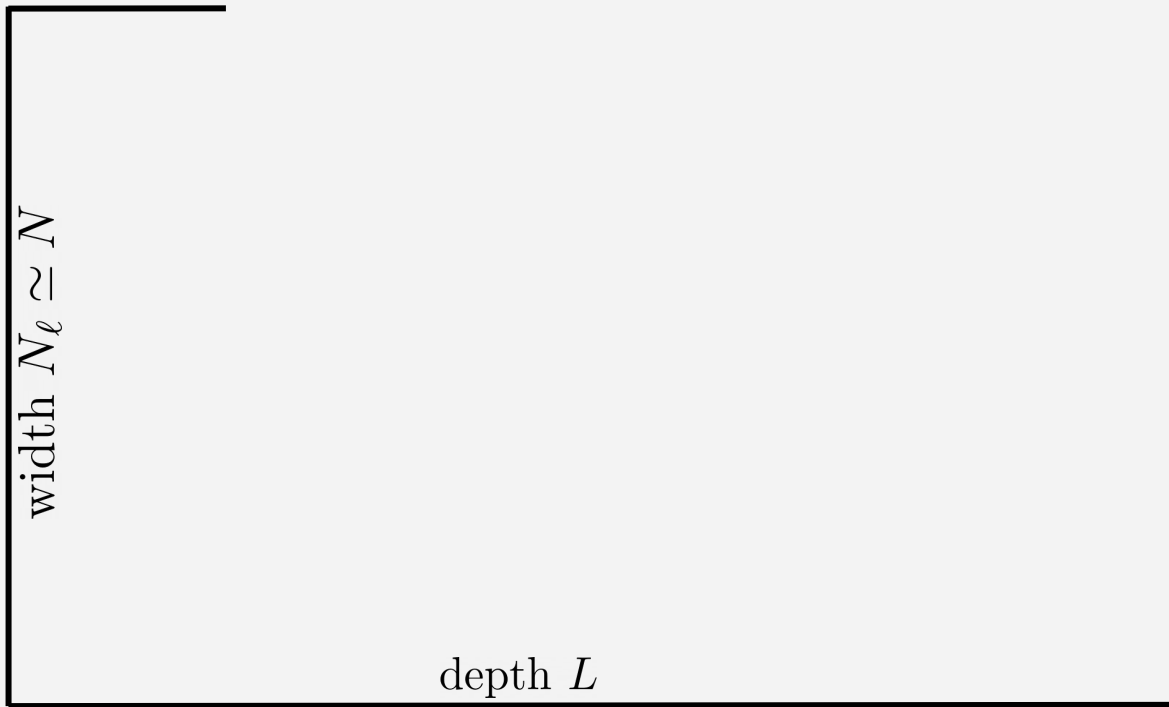
$$\left\{ \begin{array}{l} \text{model: } W^{(L+1)} \dots W^{(1)} x \\ \text{prior: } W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1}) \\ \text{likelihood: MSE on } \mathcal{D} = \{(x_i, y_i)\}_{i=1}^P \\ \text{regime: } P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1 \\ \text{prior depth: } L/N = \gamma \\ \text{post. depth: } PL/N = \lambda \end{array} \right.$$

width $N_\ell \simeq N$

depth L

Phase Diagram for Deep Linear Posteriors

$$\underline{\lambda, \gamma = 0}$$



Setup

$$\left\{ \begin{array}{l} \text{model: } W^{(L+1)} \dots W^{(1)} x \\ \text{prior: } W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1}) \\ \text{likelihood: MSE on } \mathcal{D} = \{(x_i, y_i)\}_{i=1}^P \\ \text{regime: } P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1 \\ \text{prior depth: } L/N = \gamma \\ \text{post. depth: } PL/N = \lambda \end{array} \right.$$

Phase Diagram for Deep Linear Posteriors

$$\underline{\lambda, \gamma = 0}$$

- $L = 0$

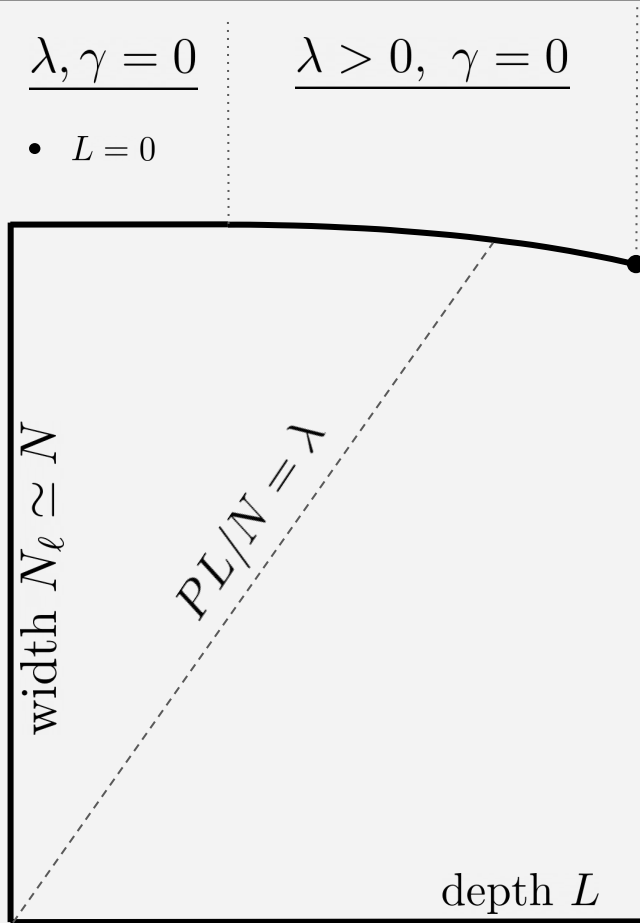
width $N_\ell \asymp N$

depth L

Setup

$$\left\{ \begin{array}{l} \text{model: } W^{(L+1)} \dots W^{(1)} x \\ \text{prior: } W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1}) \\ \text{likelihood: MSE on } \mathcal{D} = \{(x_i, y_i)\}_{i=1}^P \\ \text{regime: } P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1 \\ \text{prior depth: } L/N = \gamma \\ \text{post. depth: } PL/N = \lambda \end{array} \right.$$

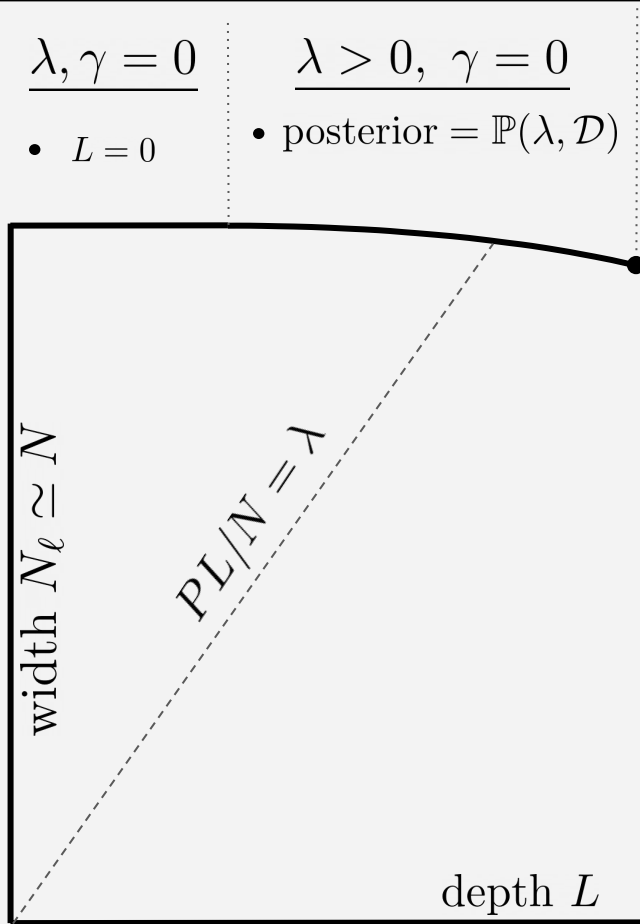
Phase Diagram for Deep Linear Posteriors



Setup

model: $W^{(L+1)} \dots W^{(1)}x$
 prior: $W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1})$
 likelihood: MSE on $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^P$
 regime: $P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1$
 prior depth: $L/N = \gamma$
 post. depth: $PL/N = \lambda$

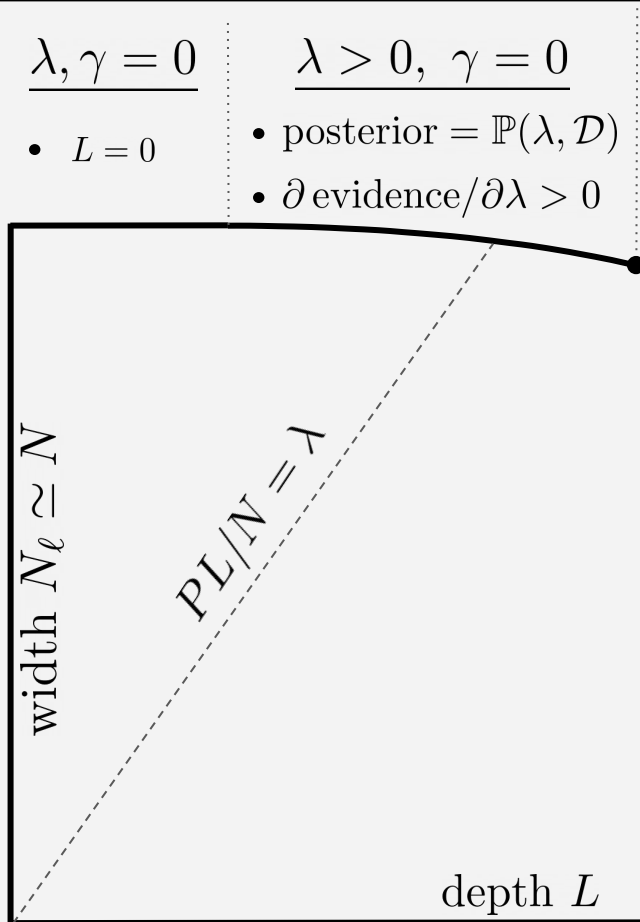
Phase Diagram for Deep Linear Posteriors



Setup

- model:** $W^{(L+1)} \dots W^{(1)}x$
- prior:** $W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1})$
- likelihood:** MSE on $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^P$
- regime:** $P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1$
- prior depth:** $L/N = \gamma$
- post. depth:** $PL/N = \lambda$

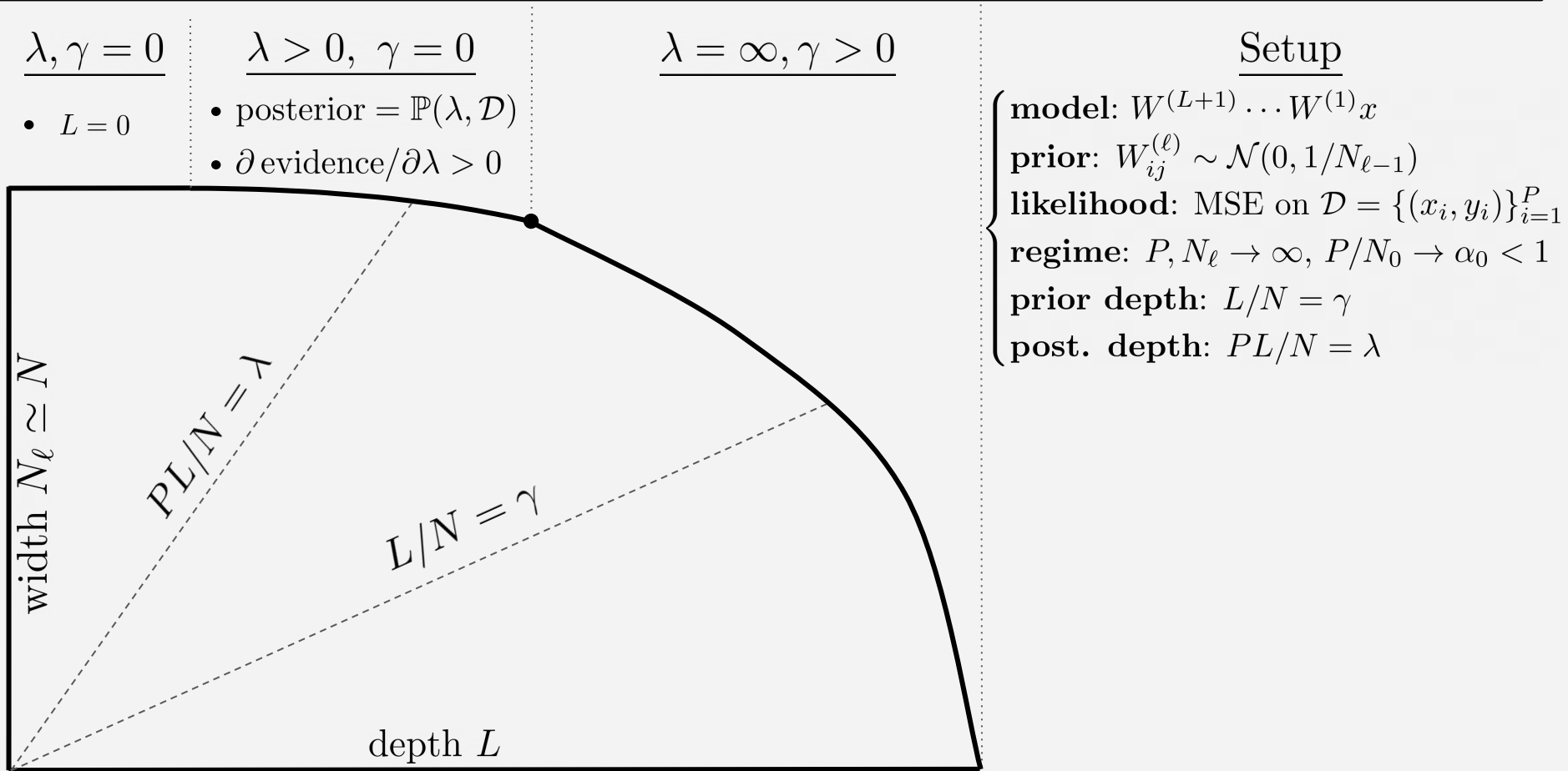
Phase Diagram for Deep Linear Posteriors



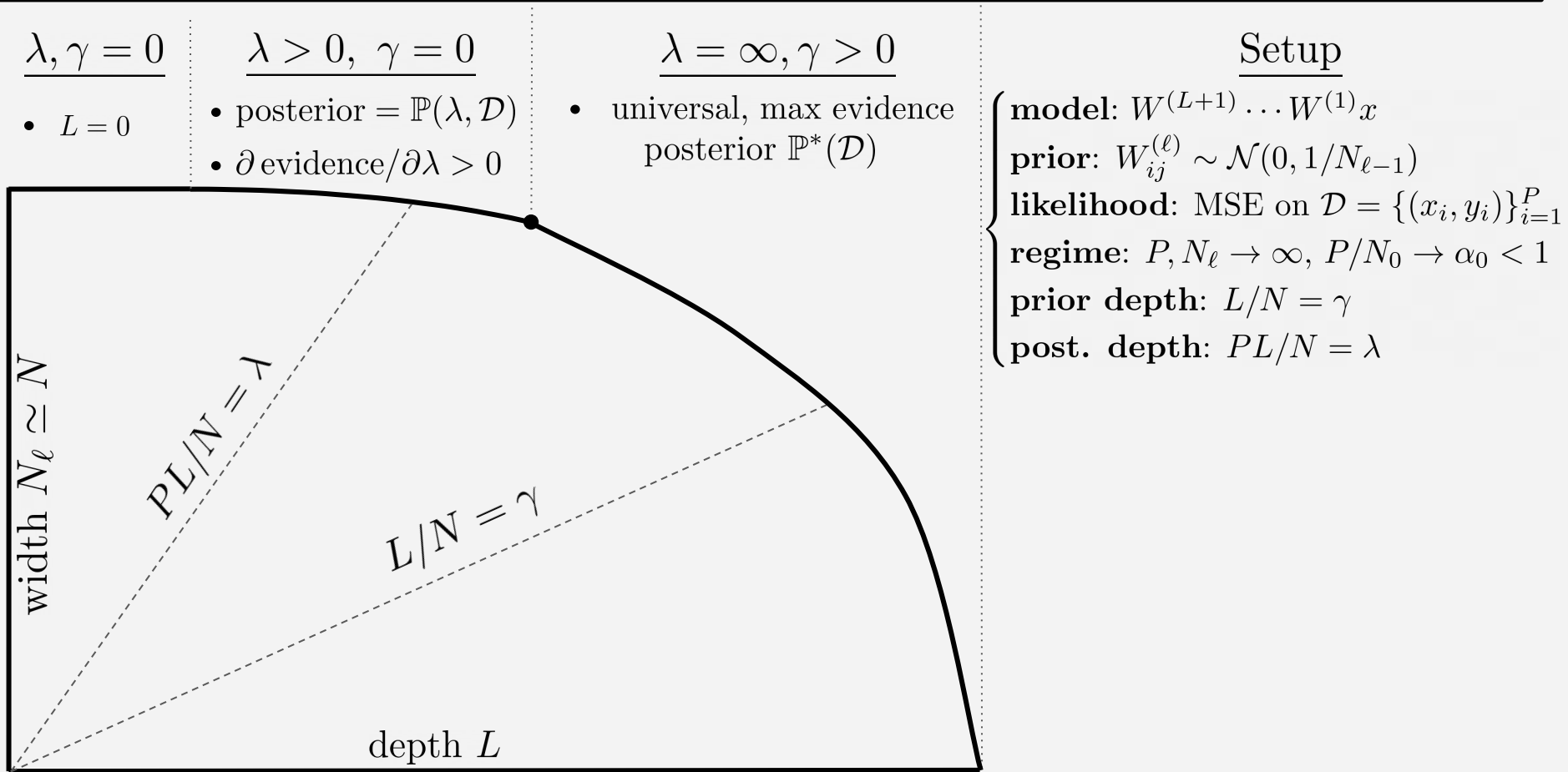
Setup

- model:** $W^{(L+1)} \dots W^{(1)} x$
- prior:** $W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1/N_{\ell-1})$
- likelihood:** MSE on $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^P$
- regime:** $P, N_\ell \rightarrow \infty, P/N_0 \rightarrow \alpha_0 < 1$
- prior depth:** $L/N = \gamma$
- post. depth:** $PL/N = \lambda$

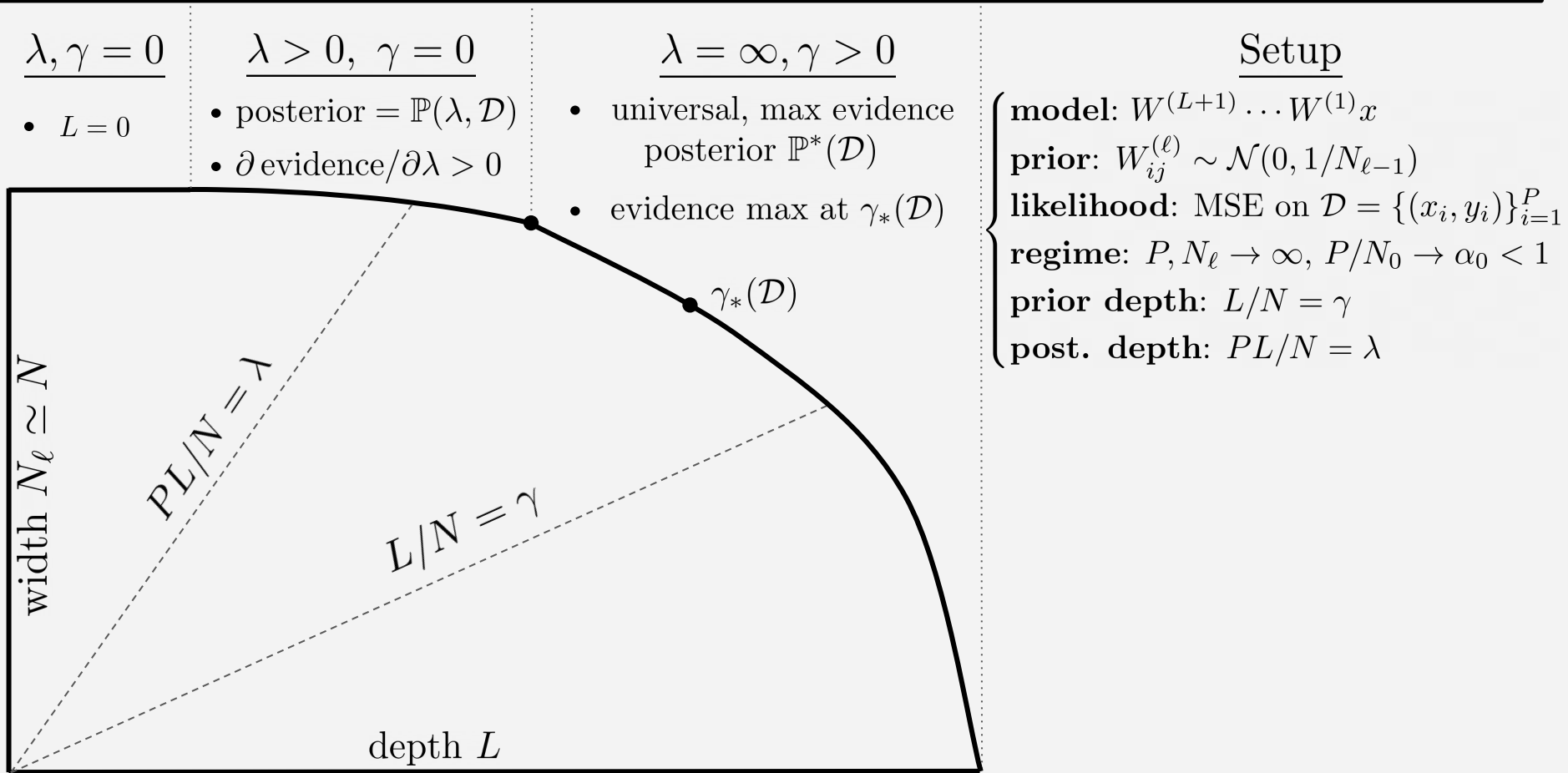
Phase Diagram for Deep Linear Posteriors



Phase Diagram for Deep Linear Posteriors



Phase Diagram for Deep Linear Posteriors



Deep Linear Networks Learn Optimal Features

Thm. Define $\sigma_* := \operatorname{argmax}_{\sigma} \lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0}} Z_{\infty}(\mathcal{D} \mid L, N_{\ell}, \sigma)$

Then,
$$\lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0}} \mathbb{P}_{\text{post}}(||\theta_{\perp}|| \mid L, N_{\ell}, \sigma = \sigma_*) = \delta_{\kappa^*}$$

for $\kappa^* = \kappa^*(\alpha_0, ||\theta_*||)$. Moreover,

$$\lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \\ PL/N \rightarrow \lambda}} \mathbb{P}_{\text{post}}(||\theta_{\perp}|| \mid L, N_{\ell}, \sigma = 1) = \delta_{\kappa(\lambda)}$$

where
$$\kappa(\lambda) \rightarrow \begin{cases} \kappa^*, & \lambda \rightarrow \infty \\ L = 0 \text{ post}, & \lambda \rightarrow 0 \end{cases}.$$