

Subspace Based Meta-Learning



Yue Sun

Joint work with:



Halil Ibrahim
Gulluk



Adhyyan Narang



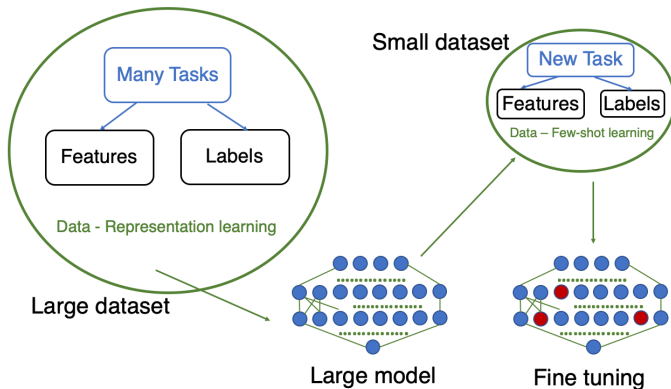
Samet Oymak



Maryam Fazel

April 21, 2021

Meta learning



Meta learning

Task, feature in \mathbb{R}^d , label in \mathbb{R} .

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, Σ_β approx low rank, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$,

Noise: $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, **Label:** $y = \mathbf{x}^\top \beta + \varepsilon$.

- **Two steps:** Representation learning, Few-shot learning

Meta learning

Task, feature in \mathbb{R}^d , label in \mathbb{R} .

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, Σ_β **approx low rank**, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$,

Noise: $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, **Label:** $y = \mathbf{x}^\top \beta + \varepsilon$.

- ▶ **Two steps:** Representation learning, Few-shot learning
- ▶ **Rep learning:** Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β , or span of Σ_β .

Meta learning

Task, feature in \mathbb{R}^d , label in \mathbb{R} .

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, Σ_β **approx low rank**, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$,

Noise: $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, **Label:** $y = \mathbf{x}^\top \beta + \varepsilon$.

- ▶ **Two steps:** Representation learning, Few-shot learning
- ▶ **Rep learning:** Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β , or span of Σ_β .
- ▶ **Few-shot learning:** Sample $\beta, \mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and estimate β in principal subspace of Σ_β .

Meta learning

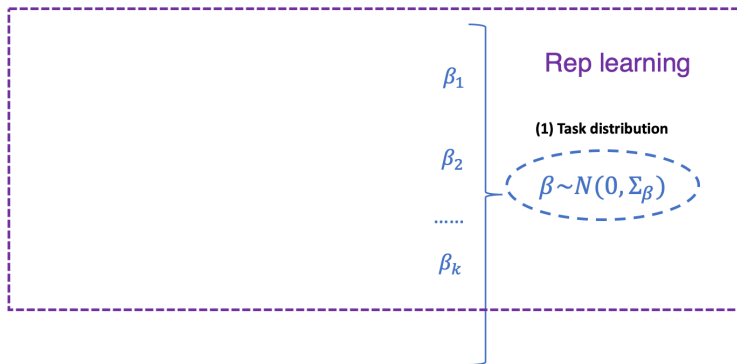
Task, feature in \mathbb{R}^d , label in \mathbb{R} .

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, Σ_β **approx low rank**, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$,

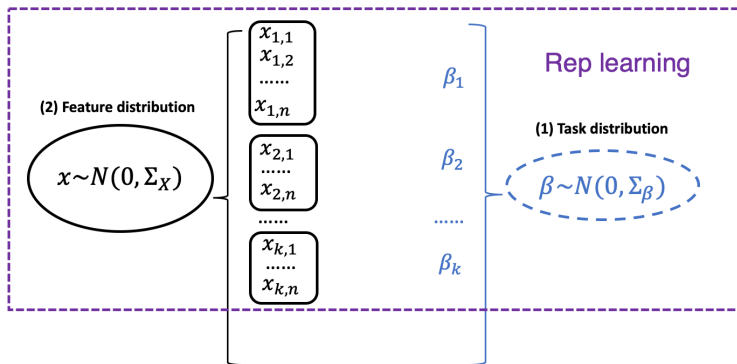
Noise: $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, **Label:** $y = \mathbf{x}^\top \beta + \varepsilon$.

- ▶ **Two steps:** Representation learning, Few-shot learning
- ▶ **Rep learning:** Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β , or span of Σ_β .
- ▶ **Few-shot learning:** Sample $\beta, \mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and estimate β in principal subspace of Σ_β .
- ▶ **Few-shot learning:** Sample $\beta, \mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

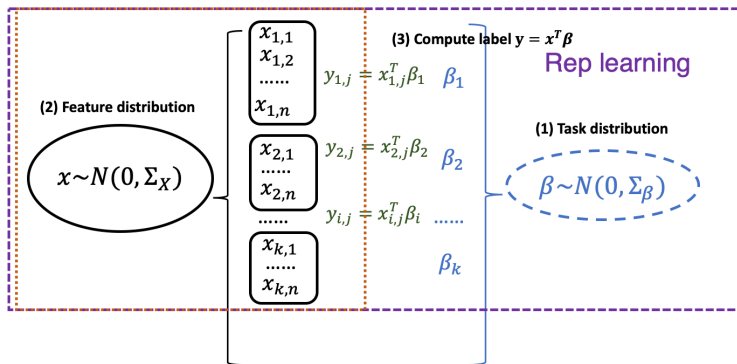
Meta learning - Linear model



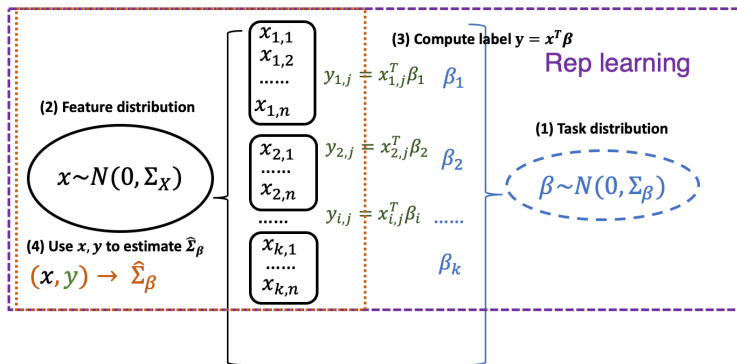
Meta learning - Linear model



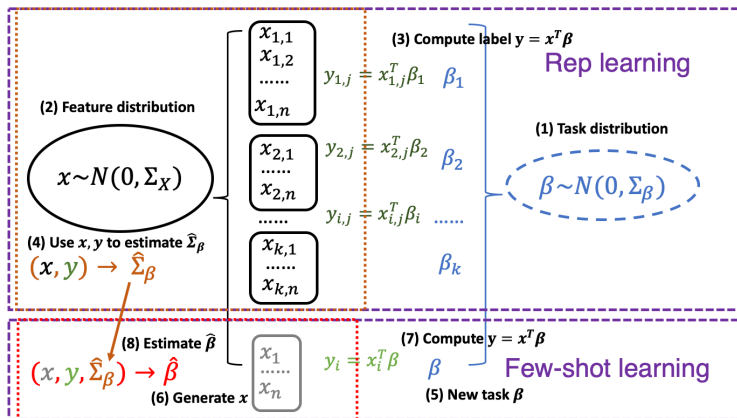
Meta learning - Linear model



Meta learning - Linear model



Meta learning - Linear model



Meta learning

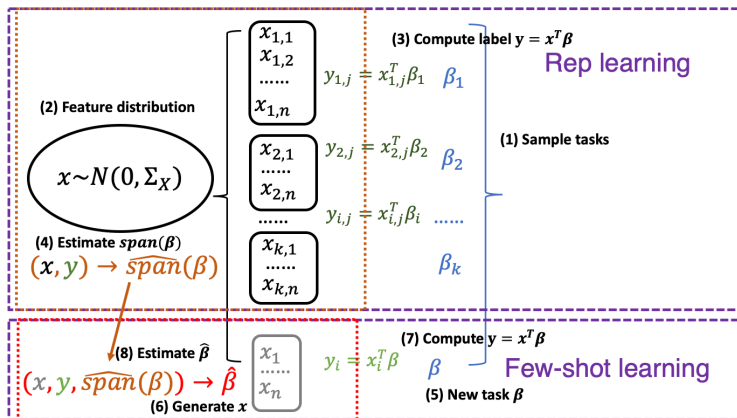
Task, feature in \mathbb{R}^d , label in \mathbb{R} .

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, Σ_β **approx low rank**, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$,

Noise: $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon)$, **Label:** $y = \mathbf{x}^\top \beta + \varepsilon$.

- ▶ **Two steps:** Representation learning, Few-shot learning
- ▶ **Rep learning:** Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β , or span of Σ_β .
- ▶ **Few-shot learning:** Sample $\beta, \mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and estimate β in principal subspace of Σ_β .
- ▶ **Few-shot learning:** Sample $\beta, \mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Meta learning - Linear model



Meta-learning - Linear - Prior works

- ▶ Mei & Montanari. Double descent.
- ▶ Du et al. Matrix factorization type. No algorithm.
- ▶ Kong et al. Method of moment (MoM) estimator. $O(dr^2)$ samples for rep learning, $O(r)$ samples for few-shot learning.
- ▶ Tripuraneni et al. MoM estimator and gradient descent. MoM: same sample complexity as above. GD: $O(dr^4)$ samples for rep learning
- ▶ Bartlett et al., Wu & Xu, Nakkiran et al. Overparameterized few-shot learning via optimal ridge regularization.

Overview

Representation learning - Linear

Few-shot learning - Linear

Meta learning - Nonlinear

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$.
Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$.
Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Question: What properties of feature & task covariance guarantees good overparameterized learning?

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Question: What properties of feature & task covariance guarantees good overparameterized learning?

Previous meta-learning work: $\Sigma_\mathbf{x} = I$.

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$.
Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Question: What properties of feature & task covariance guarantees good overparameterized learning?

Previous meta-learning work: $\Sigma_\mathbf{x} = I$.

Overparameterization: Spike feature covariance.

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Question: What properties of feature & task covariance guarantees good overparameterized learning?

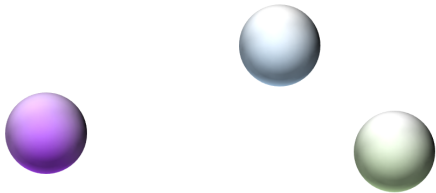
Previous meta-learning work: $\Sigma_\mathbf{x} = I$.

Overparameterization: Spike feature covariance.

Our contribution: Analysis for general cov, **feature-task alignment**.

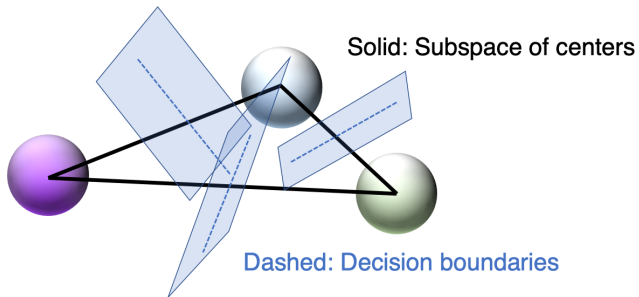
E.g., $\Sigma_\beta = \text{diag}(I_{r_t}, 0)$, $\Sigma_\mathbf{x} = \text{diag}(I_{r_f}, \iota I_{d-r_f})$.

Motivating example: Multi-class classification

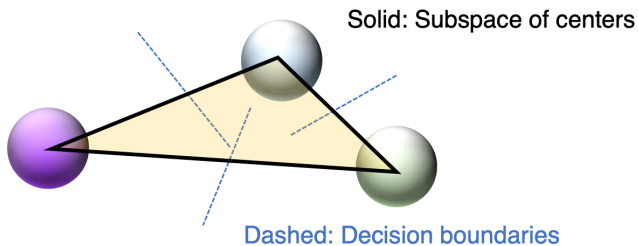


Motivation: classification of Gaussian mixture

Motivating example: Multi-class classification



Motivating example: Multi-class classification



$$\dim(\text{span}(\text{centers})) < \dim(\text{space})$$

Naive case: estimating $\Sigma_{\mathbf{X}}$ is enough

1. $\Sigma_{\beta} = \Sigma_{\mathbf{X}}$.
2. $\text{span}(\Sigma_{\beta}) = \text{span}(\Sigma_{\mathbf{X}})$.
3. $\text{span}(\Sigma_{\beta}) \subset \text{span}(\Sigma_{\mathbf{X}})$ but we are satisfied with $\text{span}(\Sigma_{\mathbf{X}})$.

Naive case: estimating $\Sigma_{\mathbf{X}}$ is enough

1. $\Sigma_{\beta} = \Sigma_{\mathbf{X}}$.
2. $\text{span}(\Sigma_{\beta}) = \text{span}(\Sigma_{\mathbf{X}})$.
3. $\text{span}(\Sigma_{\beta}) \subset \text{span}(\Sigma_{\mathbf{X}})$ but we are satisfied with $\text{span}(\Sigma_{\mathbf{X}})$.

MoM-F estimator:

$$\hat{\Sigma}_{\mathbf{X}} = \frac{1}{nk} \sum_{j=1}^n \sum_{i=1}^k \mathbf{x}_{ij} \mathbf{x}_{ij}^{\top}$$

Sample complexity: $\mathcal{O}(r_f)$. Error: $\mathcal{O}(\sqrt{r_f/(nk)})$.

General case: MoM estimator

When n is small.

$$\hat{\mathbf{Q}} = \frac{1}{k} \sum_{i=1}^k \frac{1}{n} \sum_{j=1}^n y_{ij}^2 \mathbf{x}_{ij} \mathbf{x}_{ij}^{\top}.$$

$$\hat{\mathbf{M}} = \frac{1}{k} \sum_{i=1}^k \frac{2}{n^2} \left[\sum_{j=1}^{n/2} y_{ij} y_{i(j+n/2)} \cdot (\mathbf{x}_{ij} \mathbf{x}_{i(j+n/2)}^{\top} + \mathbf{x}_{i(j+n/2)} \mathbf{x}_{ij}^{\top}) \right].$$

General case: MoM estimator

When n is small.

$$\hat{\mathbf{Q}} = \frac{1}{k} \sum_{i=1}^k \frac{1}{n} \sum_{j=1}^n y_{ij}^2 \mathbf{x}_{ij} \mathbf{x}_{ij}^\top.$$

$$\hat{\mathbf{M}} = \frac{1}{k} \sum_{i=1}^k \frac{2}{n^2} \left[\sum_{j=1}^{n/2} y_{ij} y_{i(j+n/2)} \cdot (\mathbf{x}_{ij} \mathbf{x}_{i(j+n/2)}^\top + \mathbf{x}_{i(j+n/2)} \mathbf{x}_{ij}^\top) \right].$$

Mean:

$$\mathbf{Q} = 2 \Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}} + \text{tr}(\Sigma_{\beta} \Sigma_{\mathbf{X}}) \Sigma_{\mathbf{X}}.$$

$$\mathbf{M} = \Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}}.$$

Sample complexity: $\mathcal{O}(r_f r_t^2)$. Error: $\mathcal{O}(\sqrt{r_f r_t^2 / (nk)} + \sqrt{r_t / k})$.

General case: MoM-TA estimator

We first define $\hat{\mathbf{b}}_i = \sum_{j=1}^n y_{ij} \mathbf{x}_{ij}$, for every $i = 1, \dots, k$.

$$\begin{aligned}\hat{\mathbf{B}} &= [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_k], \\ \hat{\mathbf{G}} &= k^{-1} \hat{\mathbf{B}} \hat{\mathbf{B}}^\top.\end{aligned}$$

General case: MoM-TA estimator

We first define $\hat{\mathbf{b}}_i = \sum_{j=1}^n y_{ij} \mathbf{x}_{ij}$, for every $i = 1, \dots, k$.

$$\begin{aligned}\hat{\mathbf{B}} &= [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_k], \\ \hat{\mathbf{G}} &= k^{-1} \hat{\mathbf{B}} \hat{\mathbf{B}}^\top.\end{aligned}$$

Mean:

$$\mathbf{G} = \Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}} + n^{-1} (\Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}} + \text{tr}(\Sigma_{\beta} \Sigma_{\mathbf{X}}) \Sigma_{\mathbf{X}})$$

Sample complexity:

1. Generally $\mathcal{O}(r_f r_t^2)$.
2. $\mathcal{O}(r_f r_t)$ when $n \geq r_t$.

Rep learning - learn Σ_β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_\mathbf{x})$, **Label:** $y = \mathbf{x}^\top \beta$.

Suppose $\text{rank}(\Sigma_\mathbf{x}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$. Suppose the principal subspaces of $\Sigma_\mathbf{x}$ and Σ_β align.

Rep learning: Sample β_1, \dots, β_k . For $i \in [k]$, Sample $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,n}$. Evaluate y . Use \mathbf{x}, y to estimate Σ_β .

Estimators: MoM, MoM-TA, MoM-F.

MoM: $\sum_{i,j} y_{ij}^2 \mathbf{x}_{ij} \mathbf{x}_{ij}^\top$.

MoM-TA: Let $\hat{\mathbf{b}}_i = \sum_{j=1}^n y_{ij} \mathbf{x}_{ij}$. $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_k]$. Need $n \geq r_t$.

MoM-F: $\sum_{i,j} \mathbf{x}_{ij} \mathbf{x}_{ij}^\top$.

$\Sigma_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$. Extra $(r_t/k)^{1/2}$ term in MoM and MoM-TA ignored.

feature cov	$\Sigma_\mathbf{x} = \mathbf{I}$		$\Sigma_\mathbf{x} = \text{diag}(\mathbf{I}_{r_f}, 0)$	
estimator	min sample	error	min sample	error
MoM	dr_t^2	$(dr_t^2/(nk))^{1/2}$	$r_f r_t^2$	$(r_f r_t^2/(nk))^{1/2}$
MoM-TA	dr_t	$(r_t/n)^{1/2}$	$r_f r_t$	$(r_t/n)^{1/2}$
MoM-F	-	-	r_f	$(r_f/(nk))^{1/2}$

Tradeoff between n and k

n : Sample per task

k : Number of tasks need enough tasks to estimate Σ_β

n_{tot} : Total samples $= nk$.

Question: Fix n_{tot} and change n, k .

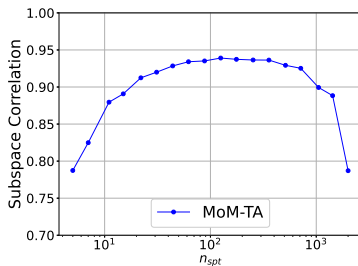
Tradeoff between n and k

n : Sample per task

k : Number of tasks need enough tasks to estimate Σ_{β}

n_{tot} : Total samples = nk .

Question: Fix n_{tot} and change n, k .



$$\Sigma_{\beta} = (I_{10}, 0_{90}), \Sigma_{\mathbf{x}} = I_{100}, \sigma_{\varepsilon} = 0.5, n_{\text{tot}} = 20000.$$

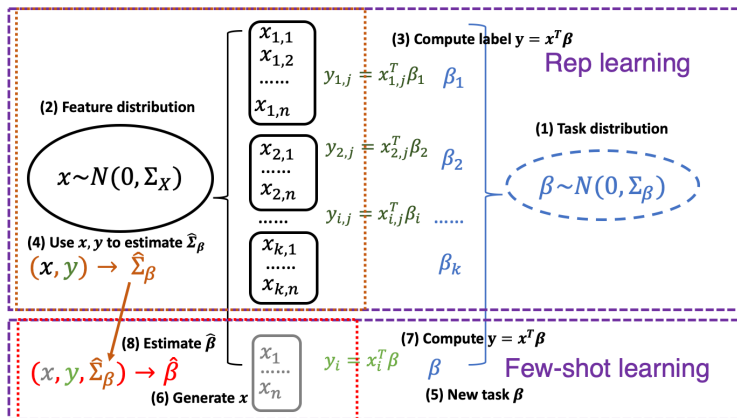
Overview

Representation learning - Linear

Few-shot learning - Linear

Meta learning - Nonlinear

Meta learning



Few-shot learning - learn β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_{\mathbf{x}}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$.

Few-shot learning: Sample β , $\mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use x, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Few-shot learning - learn β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_{\mathbf{x}}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$.

Few-shot learning: Sample β , $\mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Prior work: Restrict $\hat{\beta}$ in principal subspace of $\hat{\Sigma}_\beta$. Dimension $< n$.

Few-shot learning - learn β

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$, **Label:** $y = \mathbf{x}^\top \beta$.
Suppose $\text{rank}(\Sigma_{\mathbf{x}}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$.

Few-shot learning: Sample β , $\mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Prior work: Restrict $\hat{\beta}$ in principal subspace of $\hat{\Sigma}_\beta$. Dimension $< n$.

Our work: An arbitrary dimension R , and set a shaping matrix $\Lambda \in \mathbb{R}^{R \times d}$ as a function of $\hat{\Sigma}_\beta$ that helps with few-shot learning.

How does shaping matrix work

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$, **Label:** $y = \mathbf{x}^\top \beta$.

Suppose $\text{rank}(\Sigma_{\mathbf{x}}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$.

Few-shot learning: Sample β , $\mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Min norm solution with Λ .

$$\hat{\alpha}_\Lambda = \arg \min_{\alpha} \|\alpha\|_{\ell_2} \text{ s.t. } \mathbf{y} = \mathbf{X}\Lambda\alpha$$

$$\hat{\beta}_\Lambda = \Lambda \hat{\alpha}_\Lambda = \Lambda(\mathbf{X}\Lambda)^\dagger \mathbf{y}.$$

How does shaping matrix work

Task: $\beta \sim \mathcal{N}(0, \Sigma_\beta)$, **Feature:** $\mathbf{x} \sim \mathcal{N}(0, \Sigma_{\mathbf{x}})$, **Label:** $y = \mathbf{x}^\top \beta$.

Suppose $\text{rank}(\Sigma_{\mathbf{x}}) = r_f$, $\text{rank}(\Sigma_\beta) = r_t$.

Few-shot learning: Sample β , $\mathbf{x}_1, \dots, \mathbf{x}_n$, evaluate y . Use \mathbf{x}, y and a shaping matrix as a function of $\hat{\Sigma}_\beta$ to estimate β .

Min norm solution with Λ .

$$\hat{\alpha}_\Lambda = \arg \min_{\alpha} \|\alpha\|_{\ell_2} \text{ s.t. } \mathbf{y} = \mathbf{X}\Lambda\alpha$$

$$\hat{\beta}_\Lambda = \Lambda \hat{\alpha}_\Lambda = \Lambda(\mathbf{X}\Lambda)^\dagger \mathbf{y}.$$

$$\hat{\beta}_\Lambda = \lim_{t \rightarrow 0} \arg \min_{\beta} \|\mathbf{X}^\top \beta - \mathbf{y}\|^2 + t \beta^\top \Lambda^{-2} \beta$$

Error metric for asymptotic case

Risk function

$$\begin{aligned}\text{risk}(\Lambda, \Sigma_{\beta}) &= \mathbf{E}(y - \mathbf{x}^{\top} \hat{\beta}_{\Lambda})^2 \\ &= \mathbf{E}(\hat{\beta}_{\Lambda} - \beta)^{\top} \Sigma_{\mathbf{x}} (\hat{\beta}_{\Lambda} - \beta).\end{aligned}$$

Error metric for asymptotic case

Risk function

$$\begin{aligned}\text{risk}(\Lambda, \Sigma_\beta) &= \mathbf{E}(y - \mathbf{x}^\top \hat{\beta}_\Lambda)^2 \\ &= \mathbf{E}(\hat{\beta}_\Lambda - \beta)^\top \Sigma_\mathbf{x} (\hat{\beta}_\Lambda - \beta).\end{aligned}$$

Optimal shaping matrix

$$\Lambda^* = \arg \min_{\Lambda' \in \mathbf{S}_{++}^d} \text{risk}(\Lambda', \Sigma_\beta)$$

Error metric for asymptotic case

Risk function

$$\begin{aligned}\text{risk}(\Lambda, \Sigma_\beta) &= \mathbf{E}(y - \mathbf{x}^\top \hat{\beta}_\Lambda)^2 \\ &= \mathbf{E}(\hat{\beta}_\Lambda - \beta)^\top \Sigma_\mathbf{x} (\hat{\beta}_\Lambda - \beta).\end{aligned}$$

Optimal shaping matrix

$$\Lambda^* = \arg \min_{\Lambda' \in \mathbf{S}_{++}^d} \text{risk}(\Lambda', \Sigma_\beta)$$

As we do not know Σ_β , define

$$\Lambda = \arg \min_{\Lambda' \in \mathbf{S}_{++}^d} \text{risk}(\Lambda', \hat{\Sigma}_\beta)$$

Error metric for asymptotic case

Risk function

$$\begin{aligned}\text{risk}(\Lambda, \Sigma_\beta) &= \mathbf{E}(y - \mathbf{x}^\top \hat{\beta}_\Lambda)^2 \\ &= \mathbf{E}(\hat{\beta}_\Lambda - \beta)^\top \Sigma_\mathbf{x} (\hat{\beta}_\Lambda - \beta).\end{aligned}$$

Optimal shaping matrix

$$\Lambda^* = \arg \min_{\Lambda' \in \mathbf{S}_{++}^d} \text{risk}(\Lambda', \Sigma_\beta)$$

As we do not know Σ_β , define

$$\Lambda = \arg \min_{\Lambda' \in \mathbf{S}_{++}^d} \text{risk}(\Lambda', \hat{\Sigma}_\beta)$$

Asymptotic: Let $n, d \rightarrow \infty$ and n/d be fixed.

Computing shaping matrix

Asymptotic: Let $n, d \rightarrow \infty$ and n/d be fixed.

Let $\Sigma_{\mathbf{X}} = \mathbf{I}$ and Σ_{β} be diagonal. Let ξ solve

$$n = \sum_{i=1}^d (1 + (\xi \Sigma_{\mathbf{X}i})^{-1})^{-1}.$$

Define $\boldsymbol{\theta} \in \mathbb{R}^d$ to be $\theta_i = \frac{\xi \Lambda_i^2}{1 + \xi \Lambda_i^2}$, and the risk is

$$\text{risk}(\Lambda, \hat{\Sigma}_{\beta}) = \frac{1}{n - \|\boldsymbol{\theta}\|^2} \left(\frac{n}{d} \sum_{i=1}^d (1 - \theta_i)^2 \hat{\Sigma}_{\beta i} + \|\boldsymbol{\theta}\|^2 \sigma_{\varepsilon}^2 \right).$$

We denote the right hand side as $f(\boldsymbol{\theta}; \hat{\Sigma}_{\beta})$.

Computing shaping matrix

Asymptotic: Let $n, d \rightarrow \infty$ and n/d be fixed.

Let $\Sigma_X = I$ and Σ_β be diagonal. Let ξ solve

$$n = \sum_{i=1}^d (1 + (\xi \Sigma_X)_i)^{-1}.$$

Define $\theta \in \mathbb{R}^d$ to be $\theta_i = \frac{\xi \Lambda_i^2}{1 + \xi \Lambda_i^2}$, and the risk is

$$\text{risk}(\Lambda, \hat{\Sigma}_\beta) = \frac{1}{n - \|\theta\|^2} \left(\frac{n}{d} \sum_{i=1}^d (1 - \theta_i)^2 \hat{\Sigma}_{\beta i} + \|\theta\|^2 \sigma_\varepsilon^2 \right).$$

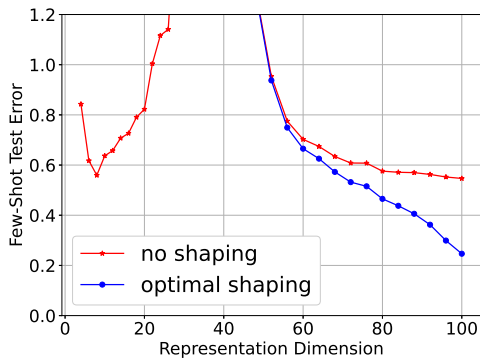
We denote the right hand side as $f(\theta; \hat{\Sigma}_\beta)$.

Computation of optimal representation:

$$\theta^* = \arg \min_{\theta} f(\theta; \hat{\Sigma}_\beta), \text{ s.t. } \underline{\theta} \leq \theta < 1, \sum_{i=1}^d \theta_i = n.$$

$$\Lambda_i^* = ((1/\theta_i^* - 1)\xi)^{-2}$$

Double descent



$$\Sigma_{\beta} = (25 \cdot \mathbf{I}_{10}, \mathbf{I}_{90}), \Sigma_{\mathbf{X}} = \mathbf{I}_{100} \sigma_{\varepsilon} = 0.5$$

Error of meta-learning

Suppose \mathcal{E} is the error of representation learning.

$$\text{risk}(\Lambda, \Sigma_{\beta}) \leq \text{risk}(\Lambda^*, \Sigma_{\beta}) + \mathcal{O} \left(\frac{n^2 \cdot \mathcal{E}}{(d - n)(2n - d\underline{\theta})\underline{\theta}} \right)$$

Error of meta-learning

Suppose \mathcal{E} is the error of representation learning.

$$\text{risk}(\Lambda, \Sigma_\beta) \leq \text{risk}(\Lambda^*, \Sigma_\beta) + \mathcal{O} \left(\frac{n^2 \cdot \mathcal{E}}{(d-n)(2n-d\underline{\theta})\underline{\theta}} \right)$$

We have presented $\Lambda \in \mathbb{R}^{d \times d}$. We can similarly define a $\mathbb{R}^{R \times d}$ representation Λ_R for arbitrary $R > n$ by projecting onto a subspace.

Error of meta-learning

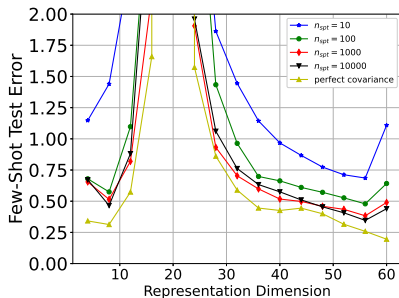
Suppose \mathcal{E} is the error of representation learning.

$$\text{risk}(\Lambda, \Sigma_\beta) \leq \text{risk}(\Lambda^*, \Sigma_\beta) + \mathcal{O} \left(\frac{n^2 \cdot \mathcal{E}}{(d-n)(2n-d\underline{\theta})\underline{\theta}} \right)$$

We have presented $\Lambda \in \mathbb{R}^{d \times d}$. We can similarly define a $\mathbb{R}^{R \times d}$ representation Λ_R for arbitrary $R > n$ by projecting onto a subspace.

Tradeoff: when R increases, $\text{risk}(\Lambda_R^*, \Sigma_\beta)$ decreases, \mathcal{E} increases.

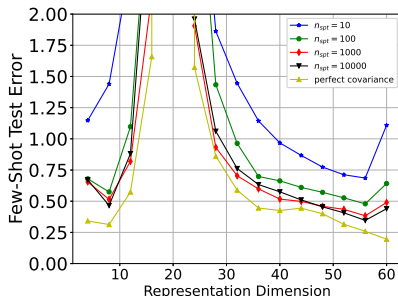
Empirical observation



$$\Sigma_{\beta} = (25 \cdot I_6, I_{54}), \Sigma_{\mathbf{X}} = I_{60}$$

We plot the error of few-shot learning versus varying dimension of Λ . Different curves correspond to different sample size for rep learning.

Empirical observation

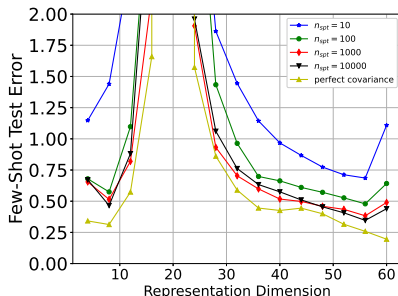


$$\Sigma_{\beta} = (25 \cdot \mathbf{I}_6, \mathbf{I}_{54}), \Sigma_{\mathbf{X}} = \mathbf{I}_{60}$$

We plot the error of few-shot learning versus varying dimension of Λ . Different curves correspond to different sample size for rep learning.

When rep learning sample size $= \infty$, $\hat{\Sigma}_{\beta} = \Sigma_{\beta}$, smallest error at $R = d$.

Empirical observation



$$\Sigma_{\beta} = (25 \cdot I_6, I_{54}), \Sigma_X = I_{60}$$

We plot the error of few-shot learning versus varying dimension of Λ . Different curves correspond to different sample size for rep learning.

When rep learning sample size $= \infty$, $\hat{\Sigma}_{\beta} = \Sigma_{\beta}$, smallest error at $R = d$. Finite sample, $\hat{\Sigma}_{\beta} \neq \Sigma_{\beta}$, smallest error when R is slightly smaller than d .

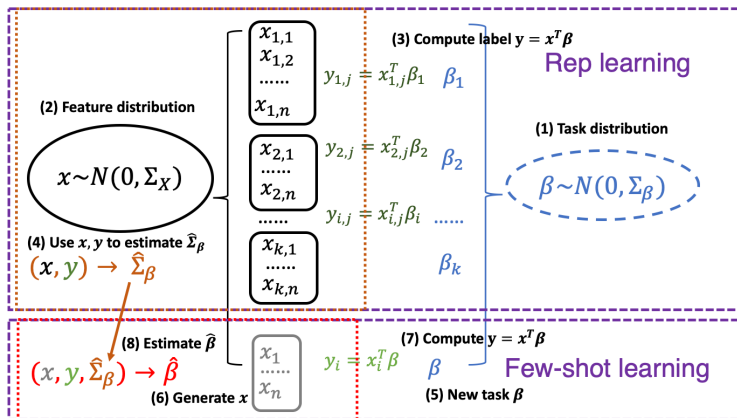
Overview

Representation learning - Linear

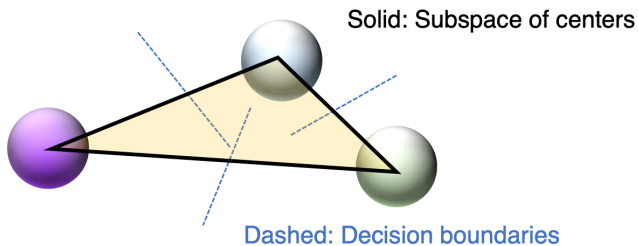
Few-shot learning - Linear

Meta learning - Nonlinear

Meta-learning - Linear



Motivating example: Multi-class classification



$$\dim(\text{span}(\text{centers})) < \dim(\text{space})$$

Meta-learning - Nonlinear - Dataset

- ▶ Fix a matrix $\mathbf{W} \in \mathbb{R}^{r \times d}$ satisfying $\mathbf{W}\mathbf{W}^\top = \mathbf{I}$.
- ▶ The i -th task is associated with function $f^i : \mathbb{R}^r \rightarrow \mathbb{R}$.
- ▶ Given input $\mathbf{x} \in \mathbb{R}^d$, the label y is distributed as $p_i(y|\mathbf{x}) = p_i(y|\mathbf{W}\mathbf{x})$ and the expectation satisfies $\mathbf{E}(y) = f^i(\mathbf{W}\mathbf{x})$.

In words, the label depends on the relevant features induced by \mathbf{W} .

Meta-learning - Nonlinear - Dataset

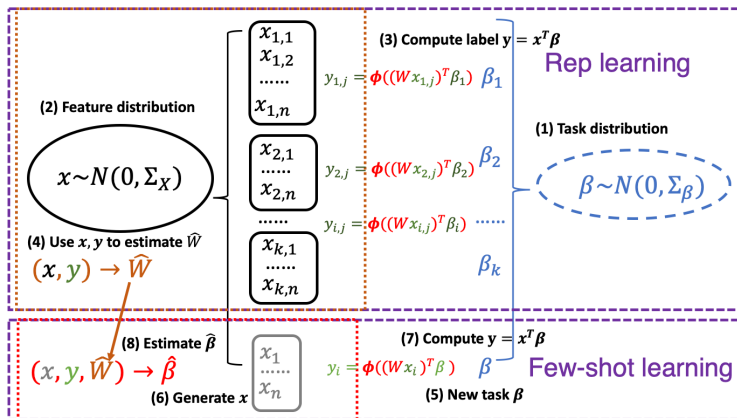
- ▶ Fix a matrix $\mathbf{W} \in \mathbb{R}^{r \times d}$ satisfying $\mathbf{W}\mathbf{W}^\top = \mathbf{I}$.
- ▶ The i -th task is associated with function $f^i : \mathbb{R}^r \rightarrow \mathbb{R}$.
- ▶ Given input $\mathbf{x} \in \mathbb{R}^d$, the label y is distributed as $p_i(y|\mathbf{x}) = p_i(y|\mathbf{W}\mathbf{x})$ and the expectation satisfies $\mathbf{E}(y) = f^i(\mathbf{W}\mathbf{x})$.

In words, the label depends on the relevant features induced by \mathbf{W} .

Example: Generalized linear models (GLM), which include logistic/linear regression, can be modeled by choosing f^i to be parameterized by a vector $\beta_i \in \mathbb{R}^r$ and a link function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ as $f^i(\mathbf{W}\mathbf{x}_{ij}) := \phi((\mathbf{W}\mathbf{x}_{ij})^\top \beta_i)$.

– logistic regression, multi-class classification, etc.

Meta-learning - Nonlinear



Representation learning

Moment estimator of covariance.

Define

$$\mathbf{v}_1 = \sum_{j=1}^{n_i/2} y_{ij} \mathbf{x}_{ij}, \quad \mathbf{v}_{-1} = \sum_{j=n_i/2+1}^{n_i} y_{ij} \mathbf{x}_{ij},$$
$$\hat{\mathbf{M}} = \sum_{i=1}^k \frac{2}{n_i^2} \left[\mathbf{v}_1 \mathbf{v}_{-1}^\top + \mathbf{v}_{-1} \mathbf{v}_1^\top \right]$$

Representation learning

Moment estimator of covariance.

Define

$$\mathbf{v}_1 = \sum_{j=1}^{n_i/2} y_{ij} \mathbf{x}_{ij}, \quad \mathbf{v}_{-1} = \sum_{j=n_i/2+1}^{n_i} y_{ij} \mathbf{x}_{ij},$$

$$\hat{\mathbf{M}} = \sum_{i=1}^k \frac{2}{n_i^2} \left[\mathbf{v}_1 \mathbf{v}_{-1}^\top + \mathbf{v}_{-1} \mathbf{v}_1^\top \right]$$

$$\mathbf{h}^i(\mathbf{W}) : \mathbb{R}^{r \times d} \rightarrow \mathbb{R}^d = \mathbf{E}_{\mathbf{x}}[f^i(\mathbf{W}\mathbf{x})\mathbf{x}]$$

$$\mathbf{M} := \mathbf{W}^\top \mathbf{W} \left(\frac{1}{k} \sum_{i=1}^k \mathbf{h}^i(\mathbf{W})(\mathbf{h}^i(\mathbf{W}))^\top \right) \mathbf{W}^\top \mathbf{W}.$$

\mathbf{M} is the mean of $\hat{\mathbf{M}}$, which is low rank.

Representation learning - Result

k tasks, each task contains n samples.

Suppose $y\mathbf{x}$ is subGaussian, $\|\mathbf{Cov}(y\mathbf{x})\| \leq \sigma^2$. (These conditions hold when $|f^i(\mathbf{x})| < \sigma$.) Let $\epsilon \in (0, 1)$.

$$kn \gtrsim \frac{d}{\epsilon^2} \Rightarrow \|\hat{\mathbf{M}} - \mathbf{M}\| \leq \epsilon\sigma^2$$

Representation learning - Result

k tasks, each task contains n samples.

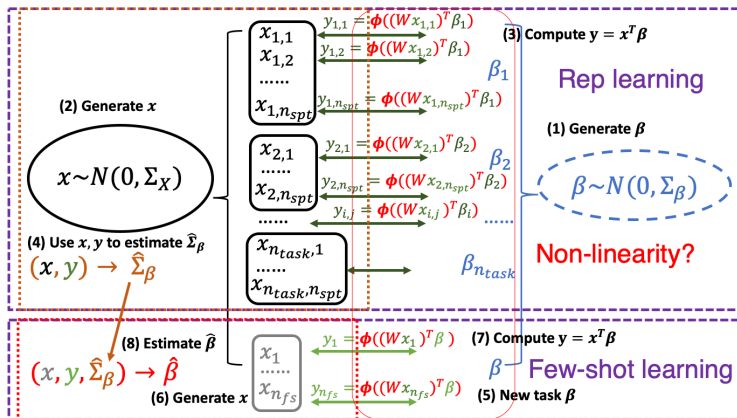
Suppose $y\mathbf{x}$ is subGaussian, $\|\mathbf{Cov}(y\mathbf{x})\| \leq \sigma^2$. (These conditions hold when $|f^i(\mathbf{x})| < \sigma$.) Let $\epsilon \in (0, 1)$.

$$kn \gtrsim \frac{d}{\epsilon^2} \Rightarrow \|\hat{\mathbf{M}} - \mathbf{M}\| \leq \epsilon\sigma^2$$

If $\lambda_r(\mathbf{M}) > \epsilon\sigma^2$, then for some orthonormal matrix $\mathbf{Q} \in \mathbb{R}^{r \times r}$,

$$\|\hat{\mathbf{W}} - \mathbf{Q}\mathbf{W}\| \leq \epsilon\sigma^2(\lambda_r(\mathbf{M}) - \epsilon\sigma^2)^{-1}.$$

Meta-learning - Nonlinear



Few-shot learning - Metric

Let $\mathcal{P}_{\mathbf{x},y}$ be the joint distribution of \mathbf{x}, y . We introduce population risk \mathcal{L} and empirical risk \mathcal{L}_e based on any single loss function between model prediction and true label.

$$\begin{aligned}\mathcal{L}(f; \mathbf{P}) &= \mathbf{E}_{\mathcal{P}_{\mathbf{x},y}} \text{loss}(f(\mathbf{P}\mathbf{x}), y) \\ \mathcal{L}_e(f; \mathbf{P}) &= \frac{1}{n} \sum_{i=1}^n \text{loss}(f(\mathbf{P}\mathbf{x}_i), y_i).\end{aligned}$$

We make the following assumption on the population risk.

1. \mathcal{L} is L Lipschitz in $\mathbf{P}\mathbf{x}$.
2. $\min_{\mathbf{P}} \mathcal{L}(f; \mathbf{P}) = \mathcal{L}(f; \mathbf{W})$.

Few-shot learning - Metric

Let $\mathcal{P}_{\mathbf{x},y}$ be the joint distribution of \mathbf{x}, y . We introduce population risk \mathcal{L} and empirical risk \mathcal{L}_e based on any single loss function between model prediction and true label.

$$\begin{aligned}\mathcal{L}(f; \mathbf{P}) &= \mathbf{E}_{\mathcal{P}_{\mathbf{x},y}} \text{loss}(f(\mathbf{P}\mathbf{x}), y) \\ \mathcal{L}_e(f; \mathbf{P}) &= \frac{1}{n} \sum_{i=1}^n \text{loss}(f(\mathbf{P}\mathbf{x}_i), y_i).\end{aligned}$$

We make the following assumption on the population risk.

1. \mathcal{L} is L Lipschitz in $\mathbf{P}\mathbf{x}$.
2. $\min_{\mathbf{P}} \mathcal{L}(f; \mathbf{P}) = \mathcal{L}(f; \mathbf{W})$.

Example: Cross entropy

$$\mathcal{L}(f; \mathbf{P}) = -\mathbf{E}_{\mathcal{P}_{\mathbf{x},y}} (y \log f(\mathbf{P}\mathbf{x}) + (1 - y) \log(1 - f(\mathbf{P}\mathbf{x}))).$$

Few-shot learning

In the few-shot learning phase, suppose $\mathbf{x}, y \sim \mathcal{P}_{\mathbf{x}, y}$ satisfy $\mathbf{E}[y \mid \mathbf{x}] = f^*(\mathbf{W}\mathbf{x})$. Let \mathcal{F} be a family of functions as the search space for few-shot learning model. We search for the solution

$$\hat{f}_e = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{L}_e(f; \hat{\mathbf{W}})$$

Few-shot learning - Result

Suppose we have n i.i.d. examples with ground-truth model $f^*(\mathbf{x}) = \phi((\mathbf{W}\mathbf{x})^\top \theta^*)$ where¹ $\|\theta^*\| \leq a$. Let \mathcal{F} be the family of functions of \mathbf{x} expressed as $\{\phi((\hat{\mathbf{W}}\mathbf{x})^\top \theta) : \|\theta\| \leq a\}$, we solve for

$$\hat{f}_e = \operatorname{argmin}_{f \in \mathcal{F}} \mathcal{L}_e(f; \hat{\mathbf{W}})$$

There exist constants $c > 1$, $\delta \in (0, 1)$, with probability at least $1 - n^{-c+1} - \delta$,

$$\begin{aligned} & \mathcal{L}(\hat{f}_e; \hat{\mathbf{W}}) - \mathcal{L}(f^*; \mathbf{W}) \\ & \leq \underbrace{\frac{caL(\sqrt{r} + \log(n))(1 + \sqrt{\log(1/\delta)})}{\sqrt{n}}}_{\text{estimation err}} + \underbrace{L\sqrt{r}\|\hat{\mathbf{W}} - \mathbf{W}\|}_{\text{rep learning err}}. \end{aligned}$$

¹bounded norm for Rademacher complexity analysis.