



Probabilistic Approaches to Transfer Learning for Sparse and Noisy Data Environments

SIAM Conference on Uncertainty Quantification

4/15/2022

PRESENTED BY

Moe Khalil

Sandia National Laboratories, Livermore
mkhalil@sandia.gov

SAND2022-XXXX



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

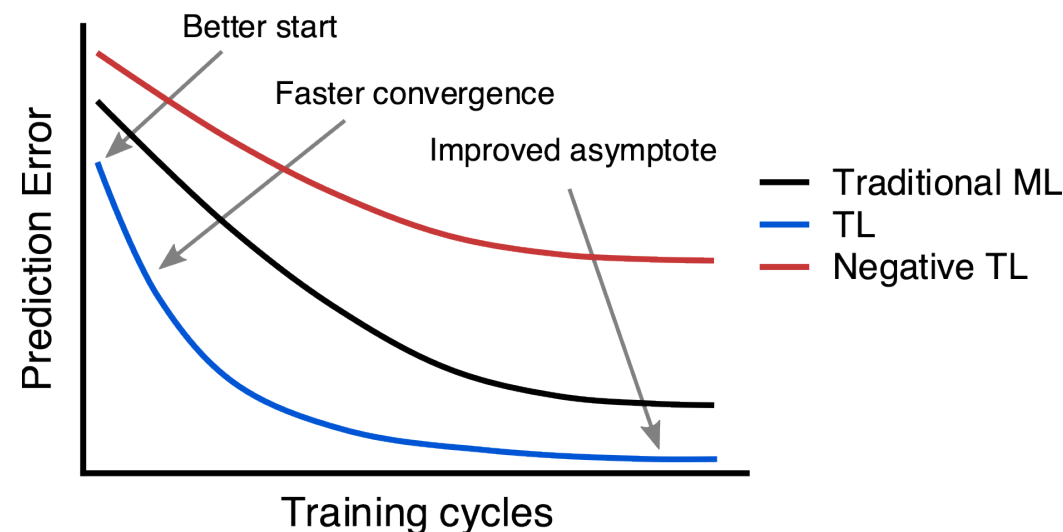
- **Sandia LDRD:**

- Moe Khalil (Sandia): principal investigator, algorithm/software development - end-to-end workflows
- Uma Balakrishnan (Sandia): numerical experiments
- Justin Jacobs (Sandia): algorithm/software development – UQ for ML
- Reese Jones (Sandia): algorithm/software development – tempering transformations
- Wyatt Bridgman (Sandia): algorithm/software development – multi-fidelity UQ
- Jackie Chen and Bruno Soriano (Sandia): TL application in turbulent combustion modeling
- Tarek Echehki (North Carolina State University): academic collaborator, TL application in turbulent combustion modeling
- Prof. Chris Pettit (US Naval Academy): TL for DoD relevant applications in atmospheric acoustics
- Erin Swansen and Naira Hovakimyan (University of Illinois at Urbana–Champaign): TL application for adaptive control of urban air mobility

- **SBIR subcontract (Cyentech Consulting LLC, Houston):**

- Moe Khalil (Sandia): principal investigator, algorithm/software development
- Cosmin Safta (Sandia): algorithm/software development
- Yueqin Huang (Cyentech Consulting LLC): customer - TL application for subsurface characterization
- Xuqing Wu (University of Houston): academic collaborator, consultant
- Jiefu Chen (University of Houston): academic collaborator, consultant
- Han Lu (University of Houston): graduate student, algorithm/software development

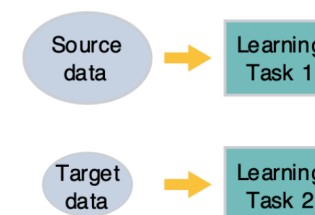
- **Challenge:** Many Sandia mission domains are defined by a lack of reliable data, effectively precluding the use of many modern deep learning/machine learning techniques for predictive modeling:
 - ✗ Excessive expense of computer simulations
 - ✗ Prohibitive experimental data acquisition cost
 - ✗ Limited access to classified and/or sensitive data
- **Goal:** *Enhance the trust in machine learning (ML) model predictions within noisy and sparse data settings*
- **Proposed Solution:** Novel probabilistic transfer learning framework.
- Transfer learning (TL): knowledge gained through similar training tasks is used to possibly improve the training process on a target domain having limited/noisy data:
 - ✓ *Improved initialization*
 - ✓ *Increased rate of convergence*
 - ✓ *Greater achievable performance*
- Proposed framework will aim to alleviate potential negative transfer: TL resulting in decreased Performance



- State-of-the-art algorithms in TL
 - ✗ tend to be ML model-specific [George et al., 2018]
 - ✗ do not consider all (if any) types of uncertainties (data, parametric, model-form/fidelity) [Colbaugh et al., 2017, Raina et al., 2006]
 - ✗ use simplified (i.e. Gaussian) probability representations of data [Karbalayghareh et al., 2018].
- Most importantly, existing methods do not address key questions relating to
 - ✗ when it is worth applying TL (as opposed to traditional ML)
 - ✗ which ML model to use in TL (out of a set of plausible ones)
 - ✗ how much knowledge is to be transferred in order to safeguard against negative learning

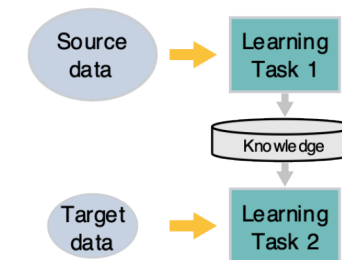
Traditional ML

- Isolated, single task learning
- Knowledge is not retained
- Learning is performed while ignoring previously obtained knowledge



Transfer Learning

- Learning a new task relies on previously learned tasks
- Learning process on new task may be more accurate with less data

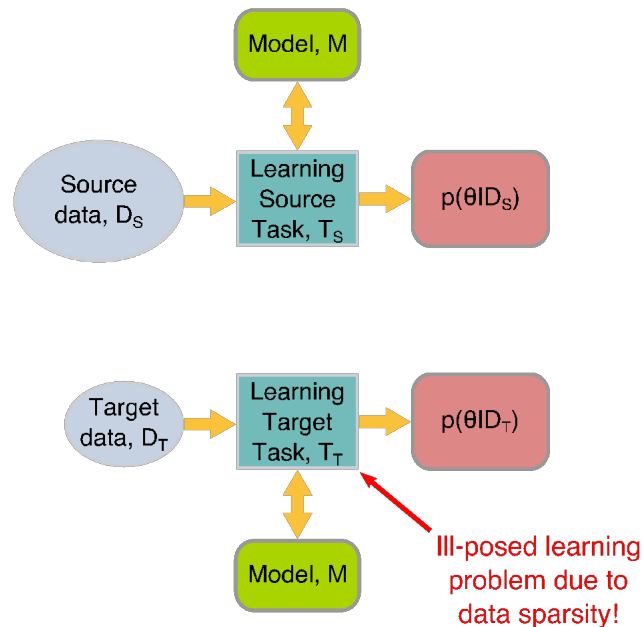


Transfer Learning Approaches	Brief Description
<i>Instance-transfer</i>	To re-weight some labeled data in the source domain for use in the target domain
<i>Feature-representation-transfer</i>	Find a “good” feature representation that reduces difference between the source and the target domains and the error of classification and regression models
<i>Parameter-transfer</i>	Discover shared parameters or priors between the source domain and target domain models, which can benefit for transfer learning
<i>Relational-knowledge-transfer</i>	Build mapping of relational knowledge between the source domain and the target domains. Both domains are relational domains and i.i.d assumption is relaxed in each domain

- Proposed TL framework aims to address the shortcomings in existing methodologies.
- It determines when to apply TL, which model to use, and how much knowledge to transfer.
- It relies on probability/measure theories to characterize and propagate uncertainties, thereby enhancing the trustworthiness of ML models in making predictions based on noisy and sparse training data.

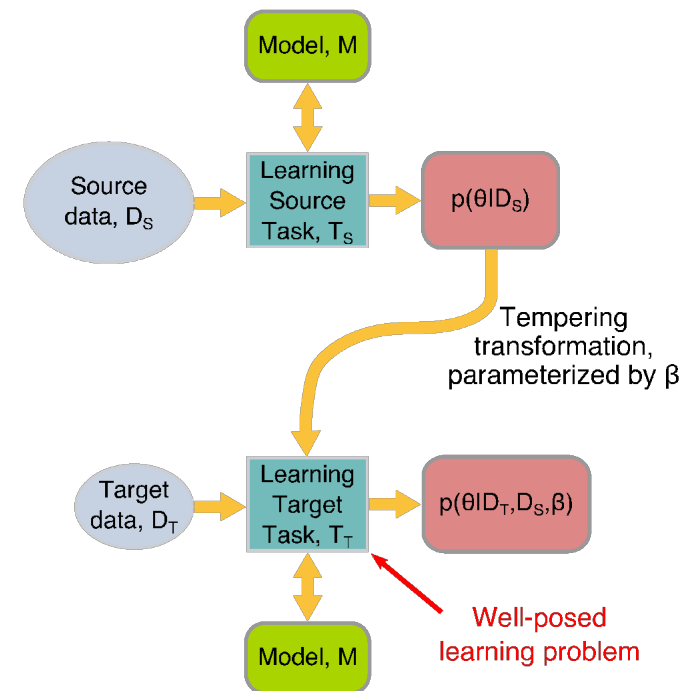
Traditional ML

- Isolated, single task learning
- Knowledge is not retained
- Learning is performed while ignoring previously obtained knowledge



Proposed TL

- Learning a new task relies on previously learned tasks
- Learning process on new task may be more accurate with sparse target data



- The proposed framework comprises of four inter-related tasks:
 - Capturing the knowledge to be transferred in training on source data
 - Provide flexibility in capturing PDFs
 - Low-fidelity Gaussian approximations (obtained using, for example, variational inference)
 - High-fidelity Gaussian mixture-models (GMM), able to characterize general non-Gaussian PDFs while enabling analytical scrutiny of the Bayesian framework.
 - Result in a spectrum of performance gains in TL
 - Propagating the knowledge to be transferred to target training tasks
 - Achieved via extensions of sequential (Bayesian) data assimilation
 - Rely on prior PDF tempering transformations (more on this later)
 - Determining how much knowledge to transfer given a choice of tempering transformation
 - Hierarchical or empirical Bayesian approaches for (joint) inference of tempering hyper-parameters
 - Information-theoretic measures; similarity and distance metrics
 - Selecting optimal ML model to use in TL
 - Probabilistic TL framework facilitates the use of Bayesian techniques for optimal model selection
 - Investigate feasibility of enhancing model complexity by leveraging Relevance Vector Machine learning techniques

$$\mathcal{M}(x, \theta) = y \approx d + \epsilon$$

Diagram illustrating the components of the ML model equation $\mathcal{M}(x, \theta) = y \approx d + \epsilon$:

- \mathcal{M} : ML model
- x : features
- θ : parameters
- y : target
- d : observation
- ϵ : noise

- Forward Problem: Given ML model, \mathcal{M} , model parameters, θ , and feature vector, x , predict “clean” targets, y
- Inverse Problem: Given a set of “noisy” observations, $D = \{d_1, \dots, d_N\}$, and feature vectors, $X = \{x_1, \dots, x_N\}$, infer parameters
 - Observations are
 - inherently noisy with unknown (or weakly known) noise model
 - sparse in space and time (insufficient resolution)
 - Problem typically ill-posed, i.e. no guarantee of solution existence nor uniqueness
- Solution: Probability density function (PDF) over the parameter space obtained using Bayes’ rule:

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

Diagram illustrating the components of Bayes’ rule equation $p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$:

- $p(\theta|D)$: posterior
- $p(D|\theta)$: likelihood
- $p(\theta)$: prior
- $p(D)$: evidence

- $p(\theta)$ is the prior PDF of θ : describes prior knowledge, inducing regularization
- $p(d|\theta)$ is the likelihood PDF of θ : describes data fit
- $p(\theta|d)$ is the posterior PDF of θ : full Bayesian solution
 - Not a single point estimate
 - Completely characterizes the uncertainty in θ
 - Subsequently used in making predictions under uncertainty

- In a transfer learning context, we have a target task of interest (regression/classification) with associated target data $\{D_T, X_T\}$. We also have access to “supplementary” source data $\{D_S, X_S\}$.
- Extending on mechanisms of propagating knowledge in sequential data assimilation (e.g. Kalman-based filters), we can take the captured knowledge from the source data in the form of the likelihood function and use it as prior knowledge in the target task:

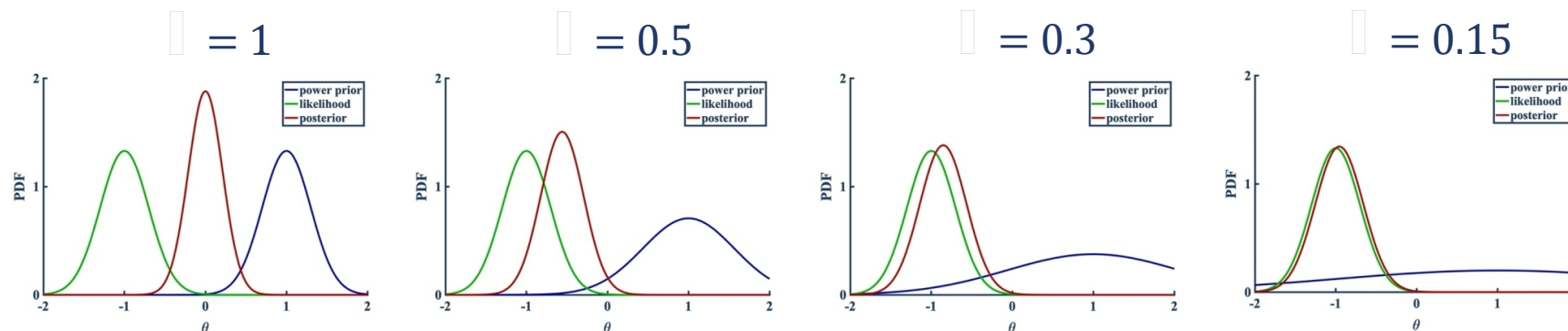
$$\text{posterior} \quad p(\theta|D_T, D_S) \propto p(D_T|\theta)p_S(\theta) \quad \text{likelihood of target data} \quad \text{prior from source data}$$

- Sequential data assimilation would dictate that the prior PDF is in fact the likelihood PDF obtained using the source data, i.e. $p_S(\theta) = p(D_S|\theta)$
- ✗ This approach does not provide flexibility in allowing the modeler to dictate how much knowledge, if any, is transferred:
 - In a traditional setting of data assimilation, all data, whether source or target, can be captured by the same model with the same parameter values (or PDFs). This assumption is not longer guaranteed to be valid in a transfer learning setting
- **Need a mechanism to control how much knowledge, if any, is transferred from source task to target task**

- How much knowledge to transfer: Tempering-based methodologies
- Tempering transformations allow us to “diffuse” or “concentrate” probabilistic knowledge (PDFs) gained through source domain learning tasks, effectively dictating how much knowledge is transferred to the target learning task
- Many PDF-tempering transformations that dictate how knowledge is transferred are envisaged
- Two proposed strategies consist of extensions/modifications of existing Bayesian priors:
 - $p_S(\theta \mid \beta) \propto p(D_S \mid \theta)^\beta$ Based on “power” priors
 - $p_S(\theta \mid \beta) = \beta p(D_S \mid \theta) + (1 - \beta) \mathcal{N}(\theta; 0, \sigma^2 I)$ Based on “mixture” priors
- For the two types of transformations above
 - **Full transfer:** $\beta \rightarrow 1$ reverts back to the full likelihood from the source training task (i.e. traditional Bayes)
 - **No transfer:** $\beta \rightarrow 0$ results in a flat prior
 - **Partial transfer:** $0 < \beta < 1$
- Optimal choice of β depends on many factors, including:
 - ML model used (can capture local vs global trends)
 - Disparity between source and target domains
 - Degree of relative data sparsity (between source and target domains)
 - Relative intensity of noise in source and target domain data

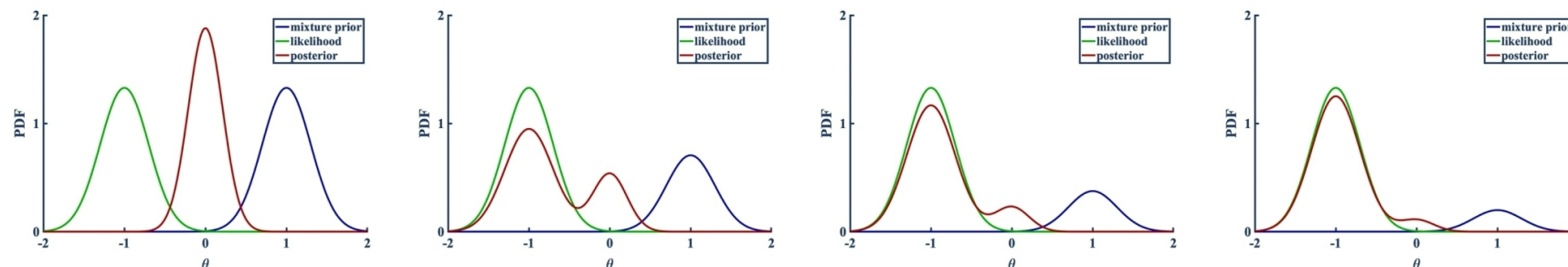
- The following is an example of the extension of power-based and mixture-based prior tempering transformation to “diffuse” knowledge in the prior PDF
- The prior and likelihood PDFs are chosen to be Gaussian
- Note: one can show that for a Gaussian PDF, raising it to a power β is equivalent to scaling the associated covariance matrix by the same β (mean vector unaffected)
- Power prior: Gaussian posterior

$$p(\theta|D, \beta) \propto p(D|\theta)p(\theta)^\beta$$



- Mixture prior: Gaussian-mixture posterior

$$p(\theta|D, \beta) \propto p(D|\theta)[\beta p(\theta) + (1 - \beta) \mathcal{N}(\theta; 0, 1 \times 10^4)]$$



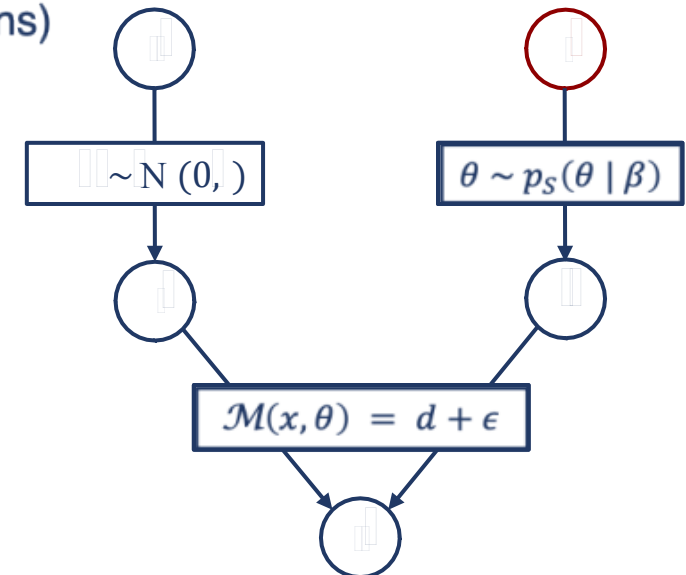
- The “tempering” hyper-parameter(s) β allow us to control the degree to which learning is transferred from the source task, characterized by the prior PDF, to the target task
- There are two approaches within a Bayesian context to determining β :

1. Hierarchical Bayes

- A fully Bayesian treatment of model parameters, θ , and noise and prior hyper-parameters, e.g. γ and β
- Proceed with joint inference of all unknowns according to joint posterior:

$$\begin{aligned} p(\theta, \beta | D) &= p(\theta | \beta, D) p(\beta | D) && \text{(probability chain rule)} \\ &\propto p(D | \theta, \beta) p(\theta | \beta) p(\beta | D) && \text{(Bayes' rule)} \\ &= p(D | \theta) p(\theta | \beta) p(\beta) && \text{(independence assumptions)} \end{aligned}$$

- Posterior distribution over the ML model parameters, θ , can be obtained by marginalizing over the hyper-parameters
- ✓ Propagates uncertainty in hyper-parameters through to parameter posterior
- ✗ Added complexity associated with inference of “less relevant” parameters and propagation of uncertainty associated with it



- There are two approaches within a Bayesian context to determining β :
 2. Empirical Bayes
 - A pseudo-Bayesian treatment of prior hyper-parameter(s), β
 - ✓ Instead of inferring and subsequently propagating uncertainties in the hyper-parameters, point estimates are obtained by maximizing *some* objective function
 - ✗ **What objective function?** Although empirical Bayes has been applied in numerous contexts for various purposes, there is not precedent for its use in transfer learning in determining such hyper-parameters
 - ✓ **Our solution:** Follow an *information-theoretic* approach, focusing on the **relative entropy** that measures the information geometry in moving from the prior to posterior:

$$\underbrace{D_{\text{KL}} [p(\theta | D, \beta), p(\theta | \beta)]}_{\text{Relative entropy (prior \& posterior)}} = \underbrace{-H[p(\theta | D, \beta), p(D | \theta)]}_{\text{Cross entropy (likelihood \& posterior)}} = \underbrace{\int \log p(D | \theta) p(\theta | D, \beta)}_{\text{Expected data-fit (posterior-averaged log-likelihood)}}$$

Relative entropy
(prior & posterior)

Cross entropy
(likelihood & posterior)

Expected data-fit
(posterior-averaged log-likelihood)

- This objective function (usually employed for Bayesian experimental design) has multiple interpretations, including an “average” log-likelihood one, resulting in a maximum “expected” likelihood estimate for hyper-parameter(s) β

- Current strategy: Empirical Bayesian treatment of the tempering hyper-parameter(s) with information-theoretic objective function in the form of the expected data-fit
- This is a data-driven approach: Optimal tempering hyper-parameters depend on the extent to which the source and target data provide “similar” information on the unknown model parameters:

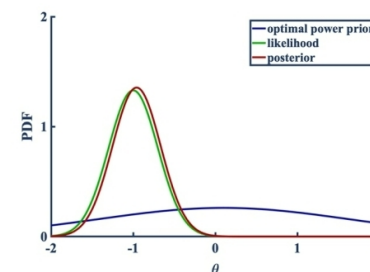
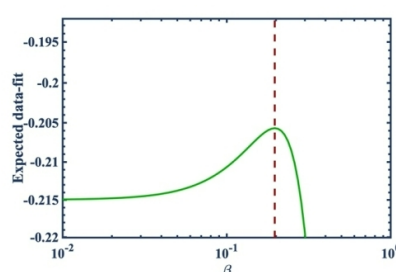
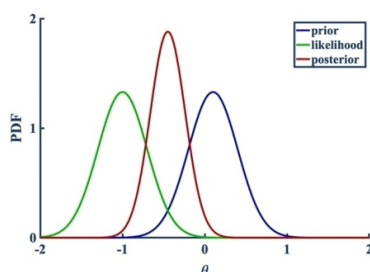
$$p(\theta|D, \beta) \propto p(D|\theta)p(\theta|\beta) \quad \beta_{\text{opt}} = \max_{\beta} E[\log p(D | \theta)] = \max_{\beta} \underbrace{\int \log p(D | \theta) p(\theta | D, \beta) d\theta}_{\text{Integral over parameter-space!}}$$

- Desirable properties:
 - Utilizes whatever available sparse target data (no need to split target data into training and testing sets)
 - Full Bayesian treatment of parametric uncertainties
 - Generally applicable to any ML model and tempering transformation
 - For Gaussian distributions (linear-in-parameter ML models) and power prior tempering transformation, **closed form expression available for the objective function and gradient/Hessian**, thus facilitating optimization step

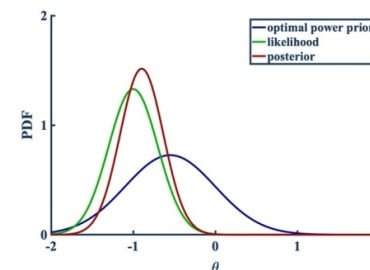
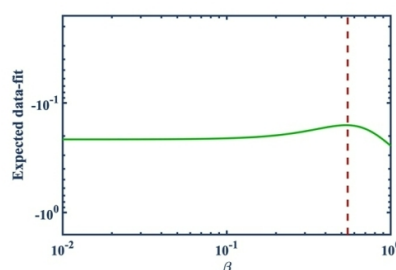
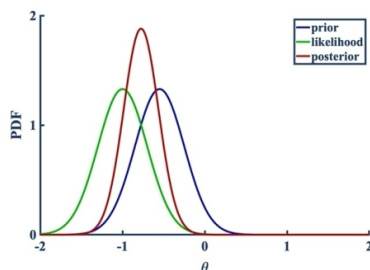
- Revisiting the 1-D example of dealing with Gaussian PDFs, let's examine the result of utilizing the expected data-fit to drive the choice of tempering hyper-parameter(s)
- Note: for Gaussian PDFs and the use power-prior based tempering transformations, the expected data-fit is available analytically

$$p(\theta|D,\beta) \propto p(D|\theta)p(\theta)^\beta$$

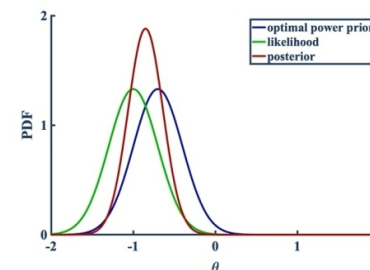
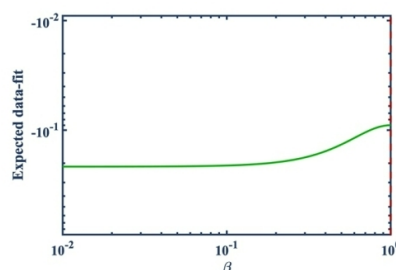
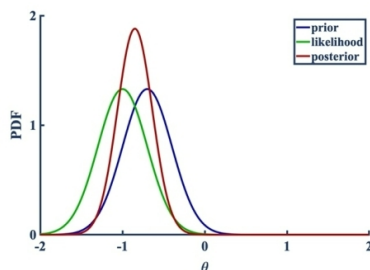
$$\beta_{\text{opt}} = \max_{\beta} E[\log p(D | \theta)] = \max_{\beta} \int \log p(D | \theta) p(\theta | \beta, D) d\theta$$



Little transfer
 $\beta_{\text{opt}} = 0.2$



Moderate transfer
 $\beta_{\text{opt}} = 0.55$

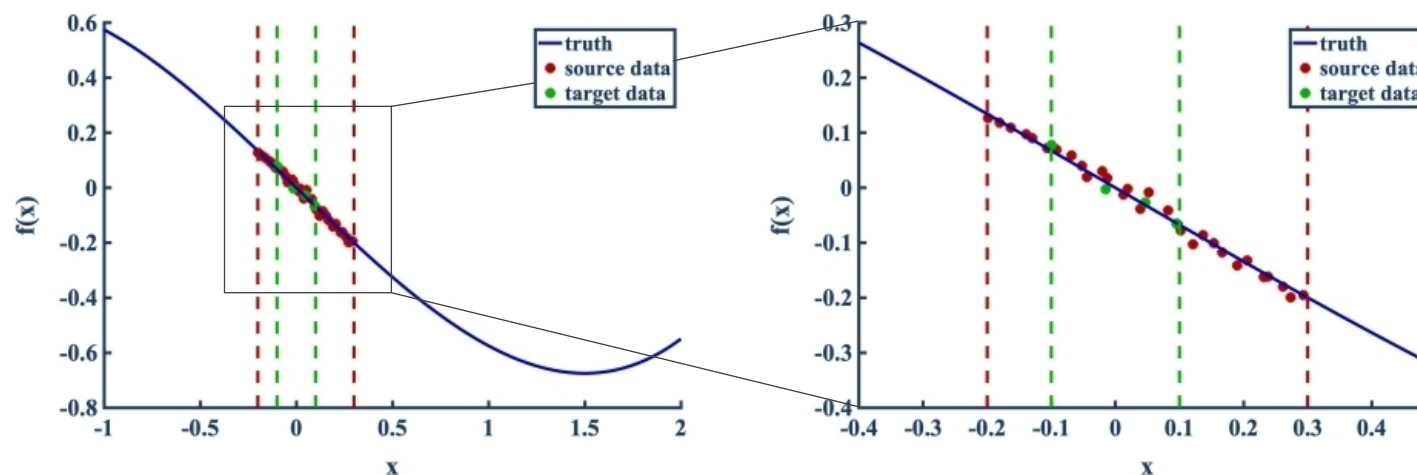


Full transfer
 $\beta_{\text{opt}} = 1$

- Assuming we're dealing with a “true” model given by

$$y_i = f(x) + \epsilon_i = 0.1x^3 - 0.75x + \epsilon_i \quad \epsilon_i \sim \mathcal{N}(0, 1 \times 10^4)$$

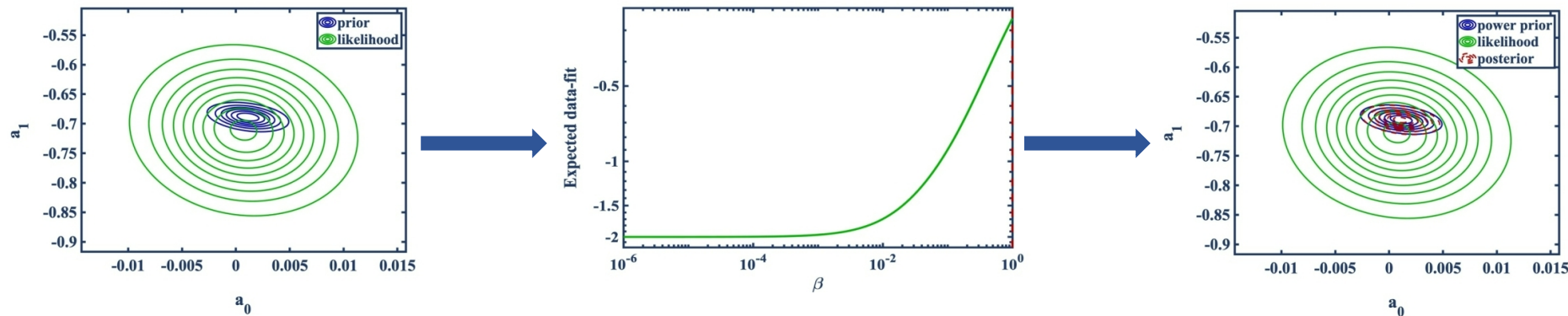
- We have 30 data points from the source domain and 4 from the target domain



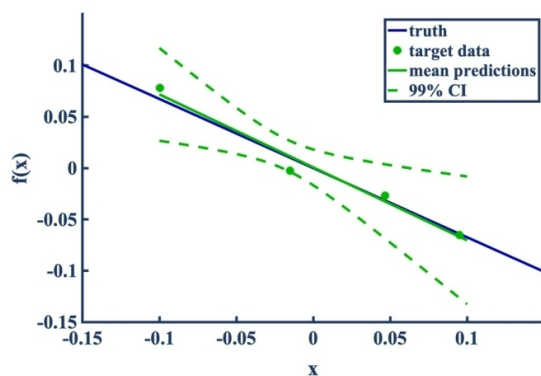
- Transfer learning task: Leverage the available source data to enhance accuracy of predictive model for target task, trained using the scarce target data
- We start with an approximate ML model to train. Let's assume a linear model:

$$y_i \approx a_0 + a_1x + \epsilon_i$$

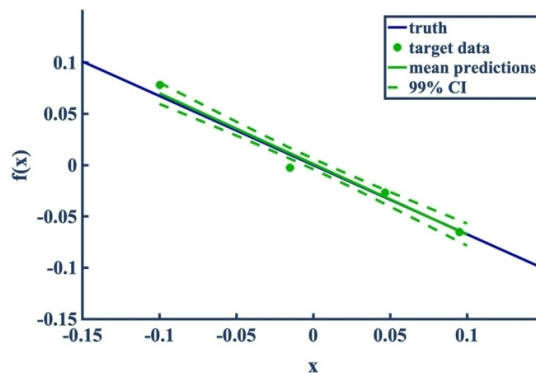
- The source data, once assimilated, provide the prior PDF for subsequent use in the target training task
- Similarly, the target data provide the likelihood PDF
- We maximize the expected data-fit to arrive at an optimal power prior



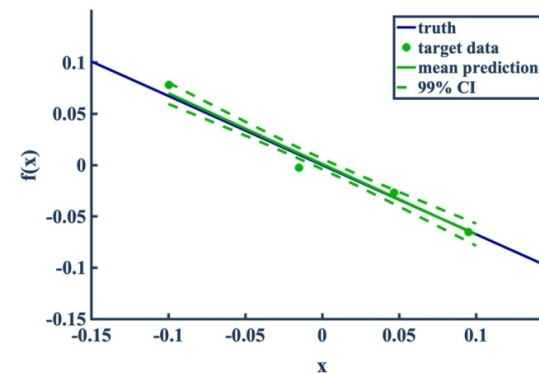
- The following are posterior predictive mean estimates and confidence intervals



no transfer



full transfer

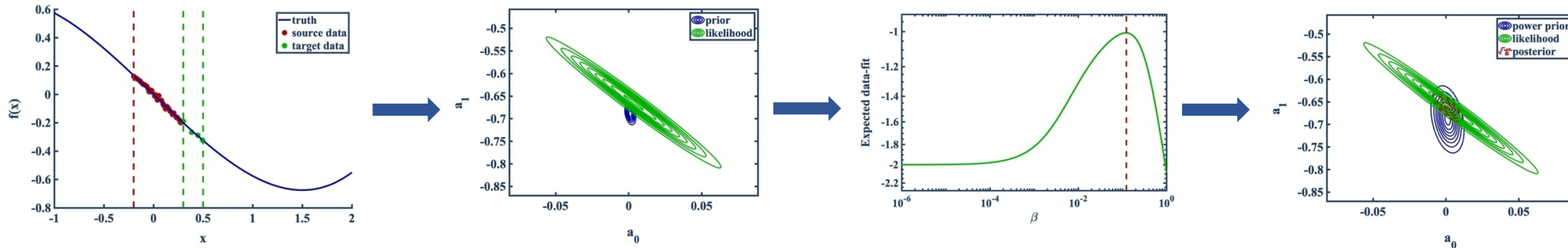


optimal transfer, $\beta_{\text{opt}} = 1$

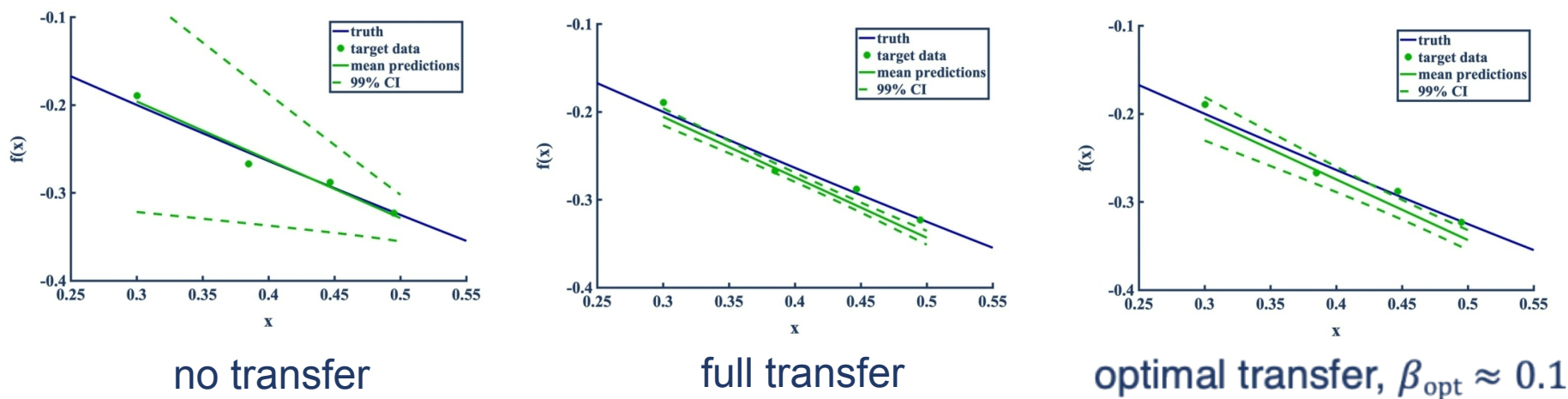
- Observation: overlapping source and target domains result in full transfer of learning

Probabilistic Transfer Learning at Work – Polynomial Surrogates

- Let's repeat the procedure with a target domain that's adjacent to the source domain (extrapolation)
- Again, we maximize the expected data-fit to arrive at an optimal power prior

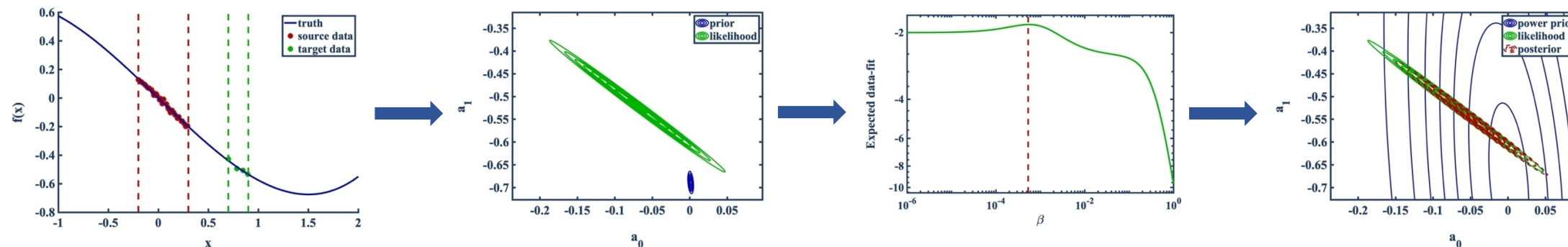


- The following are posterior predictive mean estimates and confidence intervals

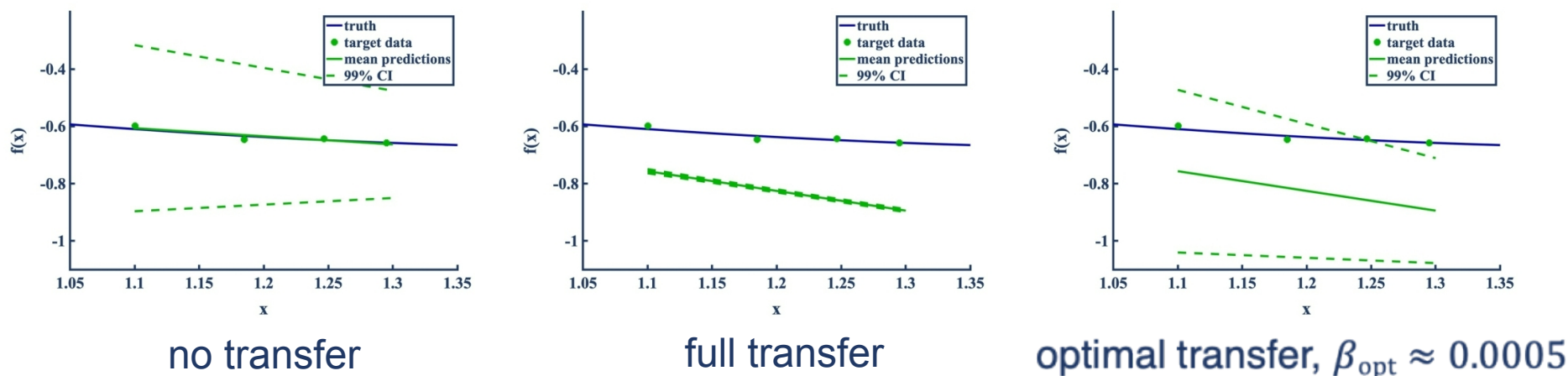


- Observation: Optimal transfer results in more accurate/precise predictions

- Lastly, we examine dissimilar target and source domains
- Again, we maximize the expected data-fit to arrive at an optimal power prior



- The following are posterior predictive mean estimates and confidence intervals



- Observation: Negligible transfer takes place with dissimilar tasks

- Developing a novel probabilistic framework for transfer learning:
 1. Reducing the validation/testing errors of such models by leveraging data from similar domains
 2. Propagating parametric, model-form, and data uncertainties towards predictions
 3. Allowing for optimal ML model selection within the TL paradigm
 4. Exhibiting moderate/strong computational scalability with increasing data volume and model complexity
 5. Safeguarding against negative learning (decreased accuracy due to task disparity, w.r.t. baseline)
 6. Consisting of strictly non-intrusive methods, applicable to most ML models without needing to modify the model (architecture) or implementation.
- Key technical steps:
 1. Capture the knowledge to be transferred in training on source data
 - Probability density functions on the calibrated ML model parameters/hyperparameters using Bayesian inversion
 - Captures data and modeling errors/uncertainties
 2. Propagate the knowledge to be transferred to target training tasks
 - Novel mechanisms for knowledge transfer that extends the traditional Bayesian approach via the application of prior PDF tempering transformations
 3. Determine how much knowledge to transfer given a choice of tempering transformation
 - Explore hierarchical or empirical Bayes approaches, based on information-theoretic measures and distance metrics
 4. Determine optimal ML model to use in TL
 - Determine the optimal ML model for TL using Bayesian model selection