

EECE 5644: Regularization

Mark Zolotas

E-mail: m.zolotas@northeastern.edu

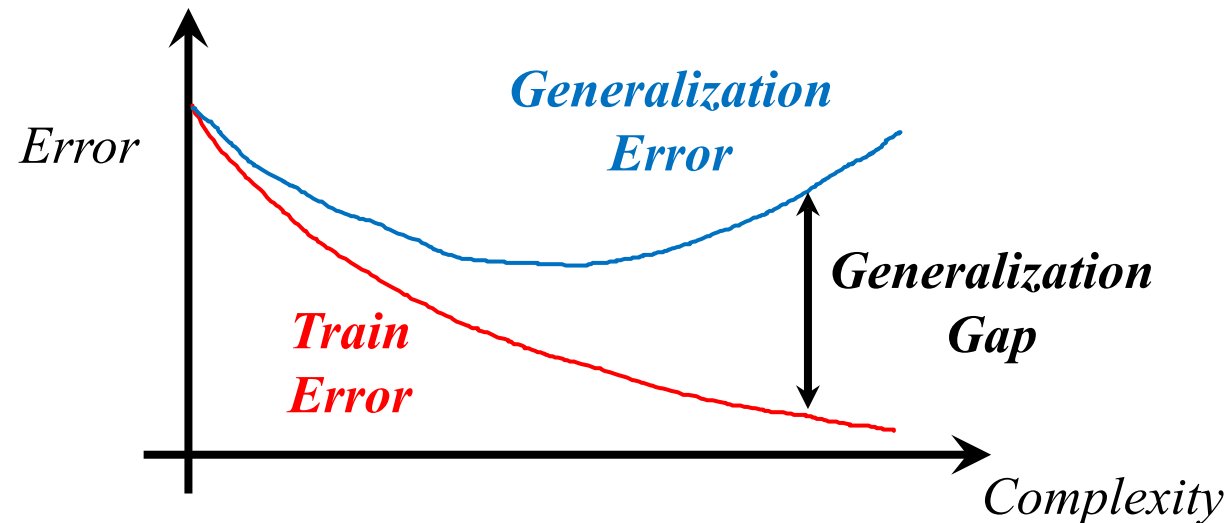
Webpage: <https://coe.northeastern.edu/people/zolotas-mark/>

Tentative Course Outline (Wks. 3-4)

Topics	Dates	Assignments	Additional Reading
Naïve Bayes Classifier & Homework 0 Practice Lab	07/18	Homework 2 released on Canvas on 07/22 Due 08/01	N/A
Model Fitting/Training: Bayesian Parameter Estimation	07/19-20		Chpts. 4.1-4.3, 8.7.2-3 Murphy 2022
Logistic Regression	07/21		Chpt. 10 Murphy 2022
Model Selection: Hyperparameter Tuning, k-fold Cross-Validation	07/26	Homework 3 released on Canvas on 07/29 Due 08/08	Chpts. 4.5, 5.2, 5.4.3 Murphy 2022
Regularization, Ridge and Lasso Regression	07/27		Chpts. 4.5, 11.1-11.4 Murphy 2022
Neural Networks: Multilayer Perceptrons & Backpropagation	07/28		Chpts. 13.1-13.5 Murphy 2022

Overfitting

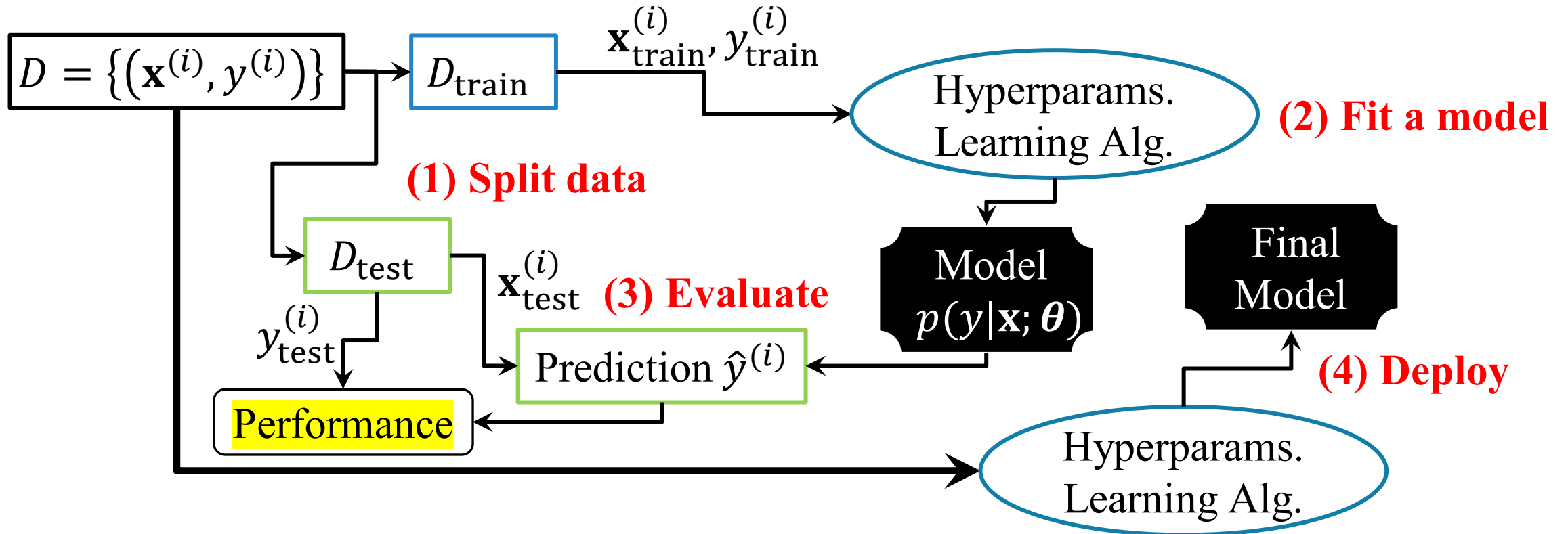
- Training models using MLE, or $L(\theta)$ in general, risks perfectly fitting the training data D_{train} and **not generalizing** well to unseen, future data...



- **Generalization gap:** $L(\theta; p^*) - L(\theta; D_{\text{train}})$
- Large generalization gap (low empirical, high theoretical loss) = **Overfitting**

Performance Estimation – Test Set

- If we are only interested in evaluating **generalization performance**:

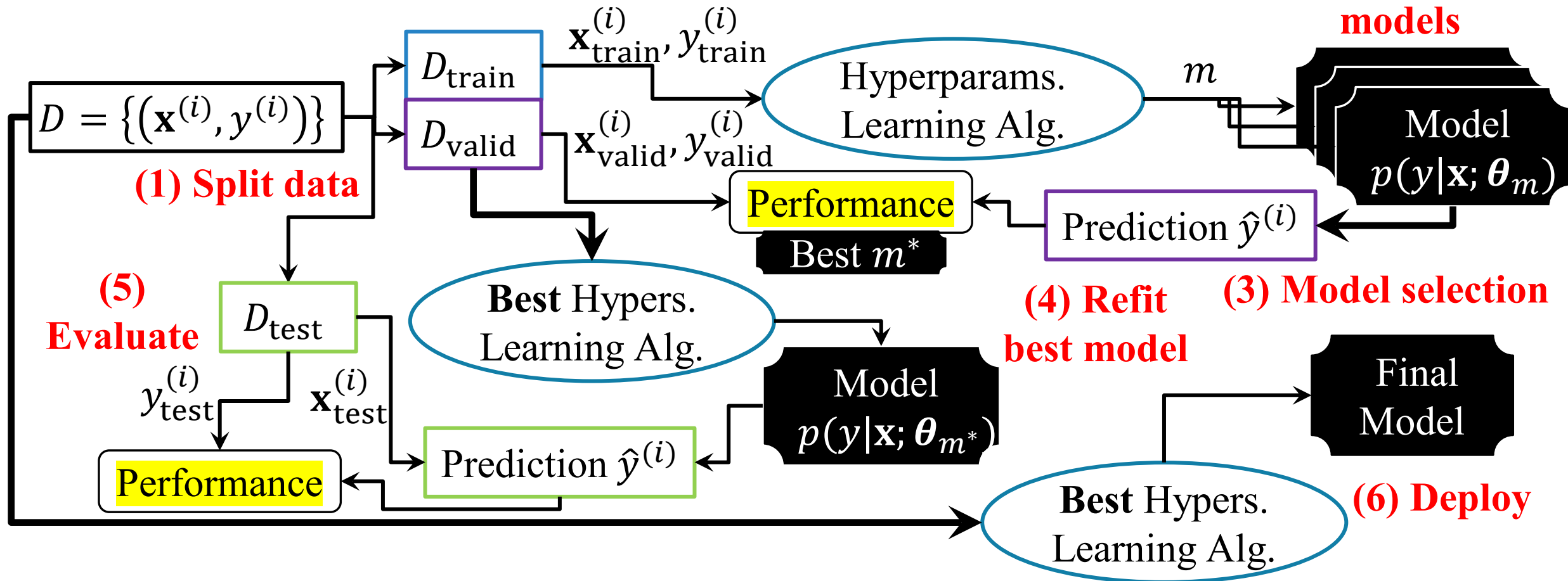


- If **large N** \rightarrow train/test split (80/20); **small N** \rightarrow K-fold CV on D_{train}

Model Selection – Validation Set

- Wish to estimate performance AND perform **model selection**:

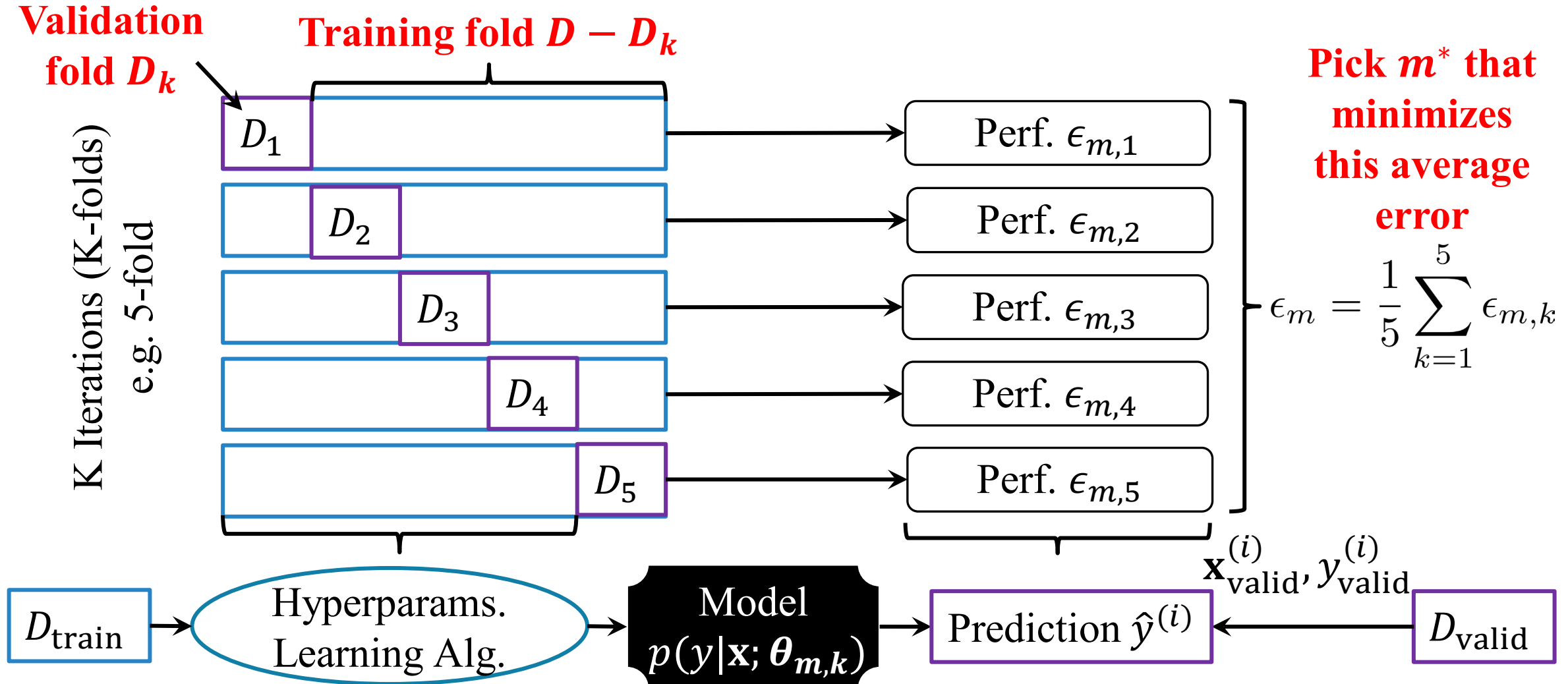
(2) Fit multiple models



- If **large N** \rightarrow train/valid/test (60/20/20); **small N** \rightarrow K-fold CV on D_{train} !

Model Selection – K-fold CV

- Inner loop of K-fold CV for one tested model m :



K-fold CV Algorithm

Algorithm 1 K-fold Cross-validation to estimate (hyper)parameters

Partition D into D_1, D_2, \dots, D_K with equal partition sizes (or close)

for $m \in \{1, \dots, M\}$ **do**

for $k \in \{1, \dots, K\}$ **do**

$$D_{\text{valid}} = D_k; D_{\text{train}} = D - D_k$$

$$\theta_{m,k}^* = \arg \min_{\theta_m} L_{\text{train}}(\theta_m; D_{\text{train}})$$

$$\epsilon_{m,k} = L_{\text{valid}}(\theta_{m,k}^*; D_{\text{valid}})$$

end for

$$\epsilon_m = \frac{1}{K} \sum_{k=1}^K \epsilon_{m,k} \quad \leftarrow \text{Avg. CV error for model } m$$

end for

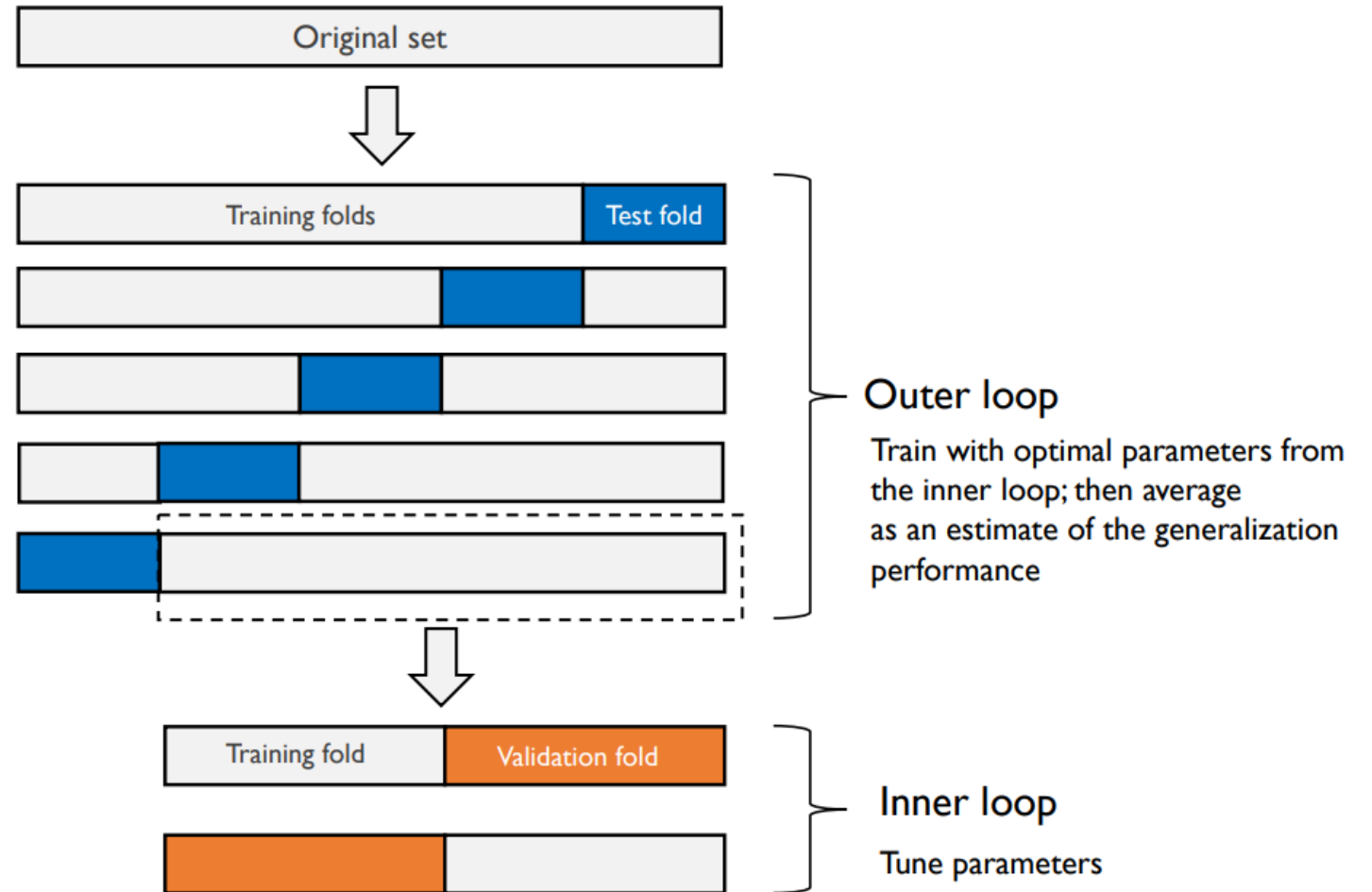
“Best” model has smallest **average** error: $m^* = \arg \min_{m \in \{1, \dots, M\}} \epsilon_m$

Given m^* , train on entire dataset: **Final trained model** $\theta_{m^*}^* = \arg \min_{\theta_{m^*}} L_{\text{train}}(\theta_{m^*}; D)$

Evaluate model m
on each K fold
of D and record
error $\epsilon_{m,k}$

Nested CV

- **Inner loop** for model selection
- **Outer loop** to estimate generalization accuracy
- Provides an **unbiased** estimate of true error
- Can be **slow**, e.g., $5 \cdot 2 \cdot M$ iterations
- Useful for **algorithm comparison**



Source: <https://androidkt.com/pytorch-k-fold-cross-validation-using-dataloader-and-sklearn/>

Techniques to Handle Overfitting

- **Hold out data** → Split into data subsets D_{train} and D_{test} (80/20)
- **CV** → K-folds of D_{valid} or D_{test} ; all data used for training eventually

Explicit • **Regularization** → Automatically controls the complexity of the model

• And others...

← Tightly coupled with MAP estimation of Θ

- Implicit* {
- ❖ **Feature selection** → Eliminate unwanted features
 - ❖ **Data augmentation** → Transform inputs for more data (e.g. rotate, flip, etc.)
 - ❖ **Early stopping** → Monitor validation performance and stop training early
 - ❖ **Feature scaling** → Expect all inputs to be similar in magnitude

Regularization

- **Key Idea** – Add a term to loss function that penalizes complexity:

$$L(\boldsymbol{\theta}; \lambda) = \frac{1}{N} \sum_{i=1}^N l\left(y^{(i)}, f(\mathbf{x}^{(i)}; \boldsymbol{\theta})\right) + \lambda C(\boldsymbol{\theta})$$

- $\lambda \geq 0$ is the **regularization** parameter (also a hyperparameter)
- $C(\boldsymbol{\theta})$ is the **complexity penalty**, e.g., negative log prior $C(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta})$

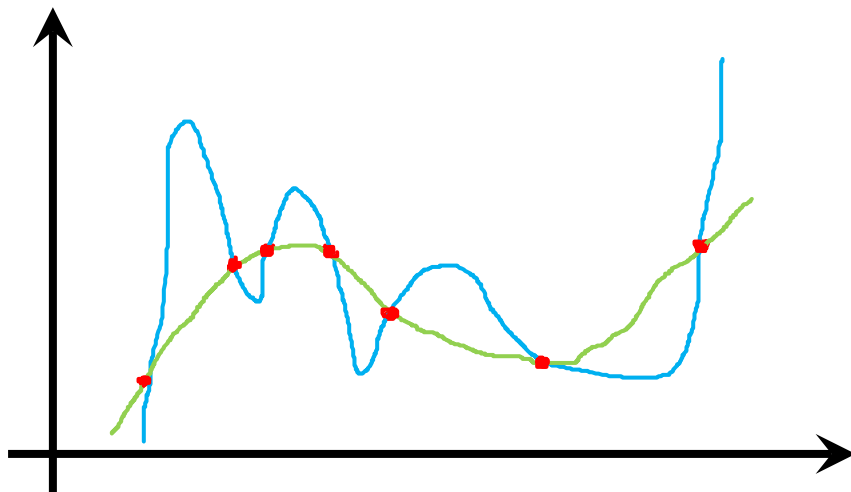
MAP
parameter
estimation
log loss

$$L(\boldsymbol{\theta}; \lambda) = -[\log p(D | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})] = -\frac{1}{N} \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta})$$

$C(\boldsymbol{\theta})$ with
 $\lambda = 1$

- Works well with **high** n -dimensional inputs

Regularization for Linear Regression

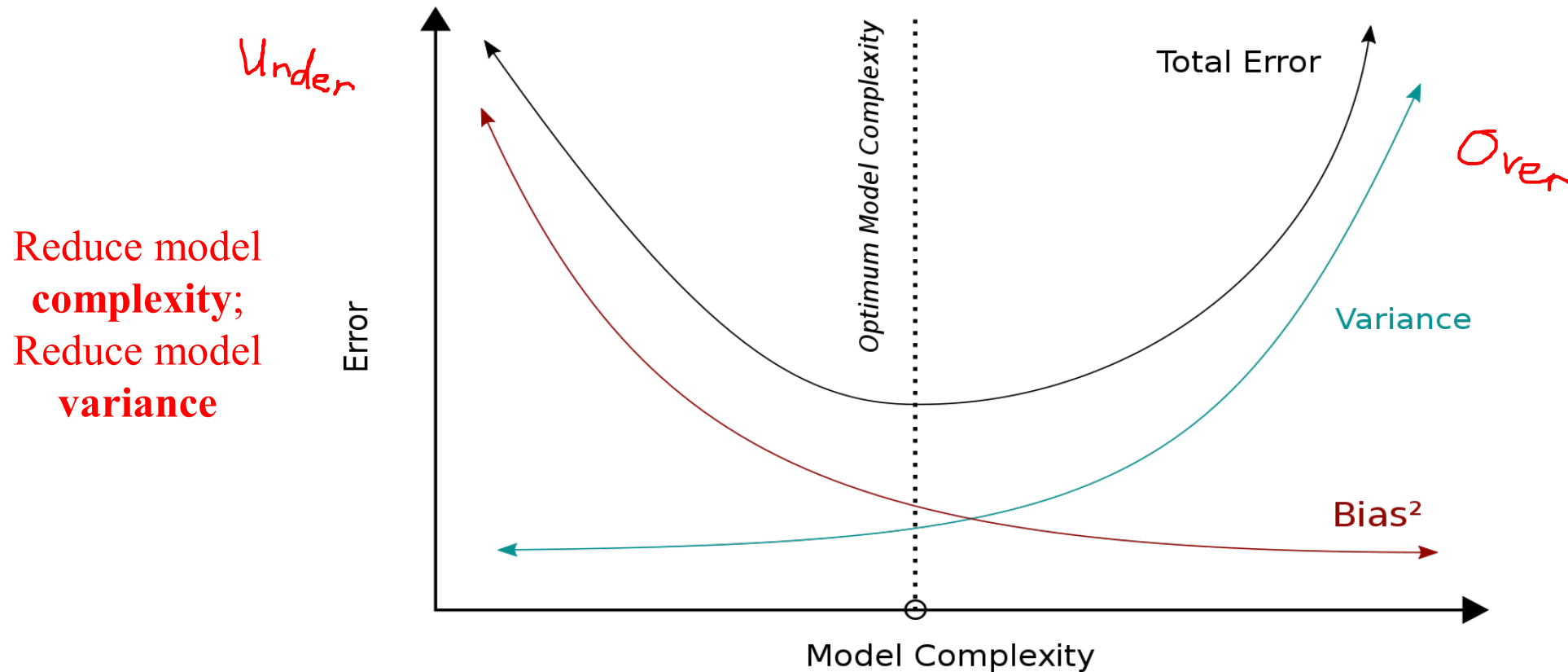


Red: Training set

Green: True target function

Blue: What we have learned (overfit)

Bias-Variance Tradeoff



Bias = Difference between estimated and true models

Variance = Model sensitivity/fluctuations to different training sets

Mean Squared Error (MSE) proportional to Variance and Bias

Least Squares Linear Regression

- Recall linear regression model as a conditional Gaussian:

$$p(y | \mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma^2)$$

- Optimize using MLE, i.e., minimize NLL:

$$\text{NLL}(\boldsymbol{\theta}) = \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2}_{\text{RSS}} + \underbrace{\frac{N}{2} \log(2\pi\sigma^2)}_{\text{const}}$$

- Assume fixed variance, e.g. $\sigma = 1$, and focus on $\mu = \mathbf{w}^\top \mathbf{x}$:

$$\text{NLL}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Ridge Regression – Formulation

- Penalize large weight magnitudes of $\text{NLL}(\boldsymbol{\theta})$ to avoid overfitting

Also called l_2 regularization or weight decay

$$\begin{aligned}\text{PNLL}(\boldsymbol{\theta}) &= \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \tilde{\mathbf{x}}^{(i)})^2 + \lambda \sum_{j=1}^n w_j^2 \\ &= \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \frac{\lambda \mathbf{w}^\top \mathbf{w}}{\lambda \|\mathbf{w}\|_2^2}\end{aligned}$$

Not bias w_0
 L_2 -norm of \mathbf{w}

- $\lambda \geq 0$ is the regularization parameter
- As if computing the $\hat{\boldsymbol{\theta}}_{\text{MAP}}$ estimate with a **zero-mean Gaussian prior** on the parameters or in our cases, weights: $p(\mathbf{w}) = N(\mathbf{w} \mid \mathbf{0}, \lambda^{-1}\mathbf{I})$

$$\boldsymbol{\theta} = [w_1 \cdots w_n]^\top$$

Lasso Regression – Formulation

- Uses a **Laplace prior** instead, with a resulting form:

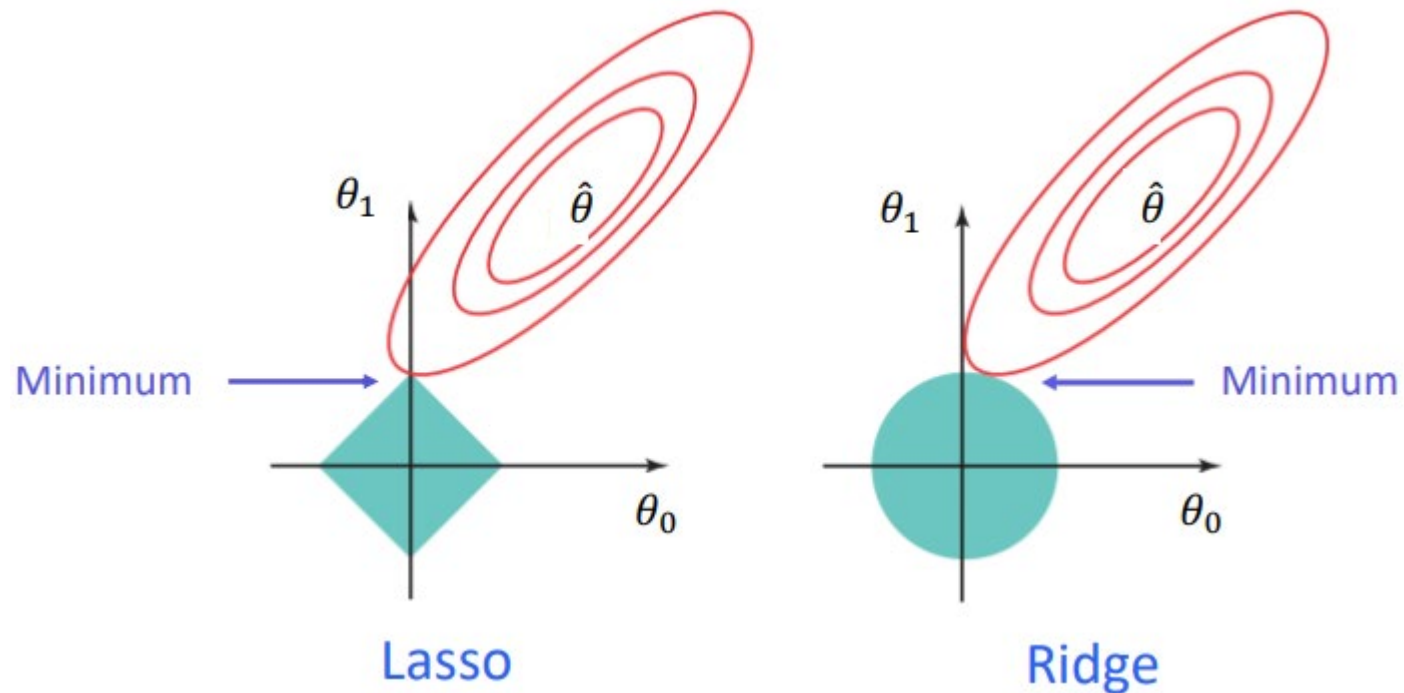
$$\begin{aligned}\text{PNLL}(\boldsymbol{\theta}) &= \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \tilde{\mathbf{x}}^{(i)})^2 + \lambda \sum_{j=1}^n |w_j| \\ &= \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_1\end{aligned}$$

L₁-norm
of \mathbf{w}
Also called l_1
regularization

- **Lasso:** “least absolute shrinkage and selection operator”
- **Warning:** Gradients cannot be computed around zero (absolute value)

Geometric Interpretation

- **Ridge** regression **shrinks** all coefficients (parameters)
- **Lasso** regression sets coefficients to **zero** (sparse solution); **feature selection**



Regularization for Logistic Regression

Same idea as **Ridge** regression:

$$\text{PNLL}(\boldsymbol{\theta}) = \text{NLL}(\boldsymbol{\theta}) + \lambda \mathbf{w}^\top \mathbf{w}$$

$$\nabla_{\boldsymbol{\theta}} \text{PNLL}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \text{NLL}(\boldsymbol{\theta}) + 2\lambda \mathbf{w}$$

Concluding Remarks

- **Explicit** regularization adds a term to optimization problem to penalize complexity, e.g., **prior** term for **MAP** parameter estimation
- **Implicit** regularization refers to other techniques like early stopping
- Some early CV-related slides adapted from:

<https://arxiv.org/pdf/1811.12808.pdf>

- Questions?