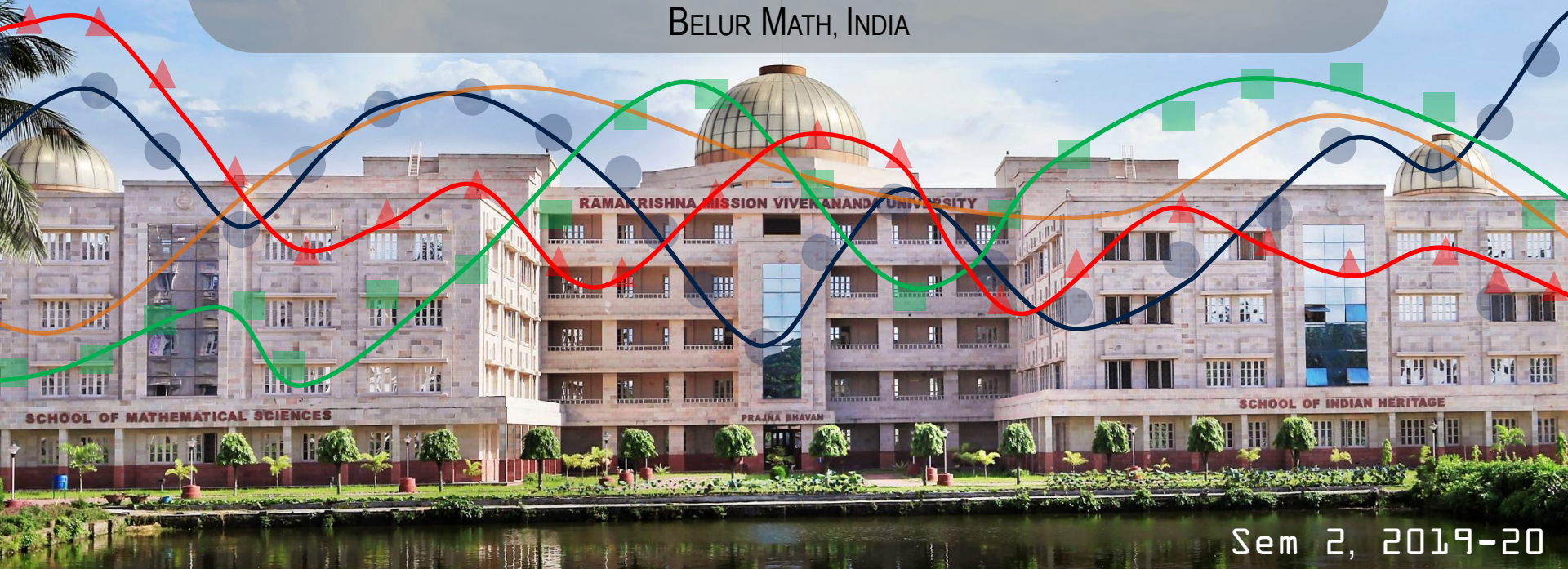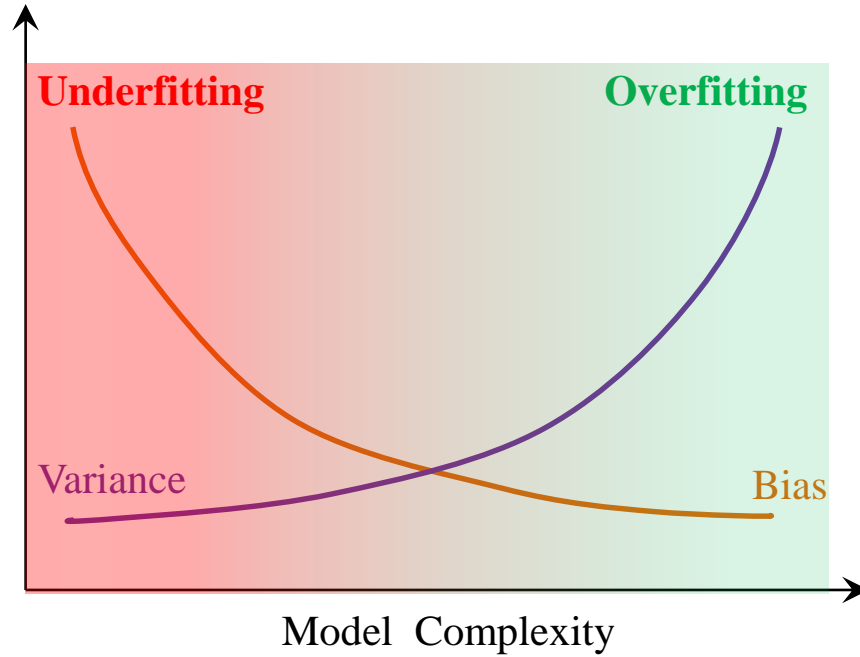# Regularization

**DRIPTA MJ**

Department of Mathematics

RAMAKRISHNA MISSION VIVEKANANDA EDUCATIONAL AND RESEARCH INSTITUTE

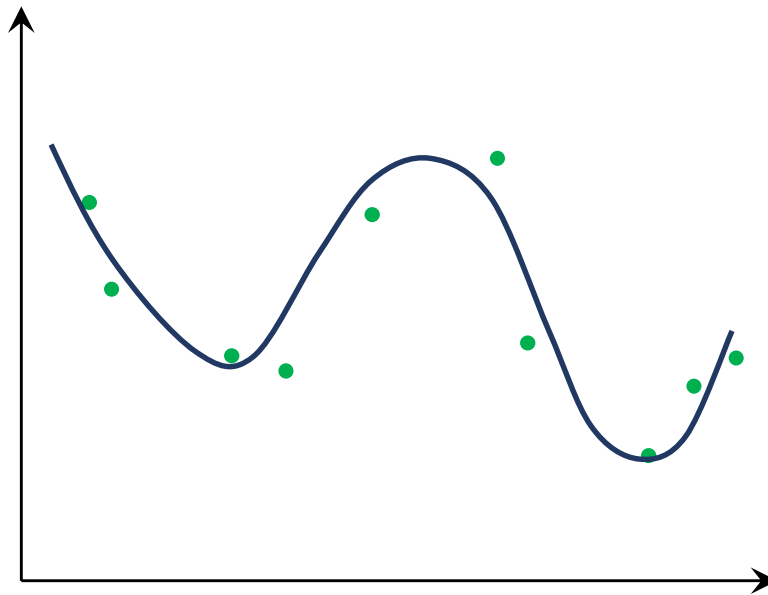BELUR MATH, INDIA

Sem 2, 2019-20
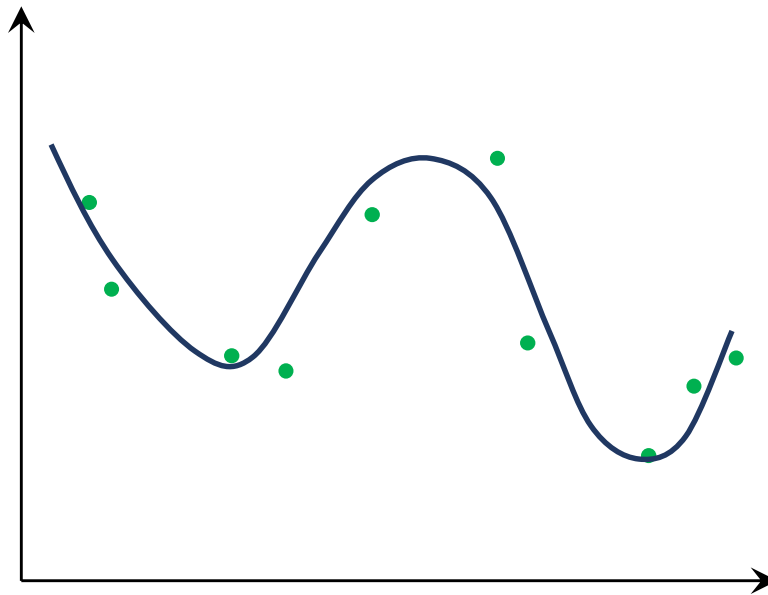
# Bias-Variance trade-off

- Data generated by a $Q$th order polynomial + some noise



- Consider fitting data with a polynomial of order $M$.

- Data preprocessing:
  - Standardize the inputs.
  - Center the outputs.

- Model can be trained using linear regression with $\left[x^1, x^2, ...., x^M\right]$ as features.

- The intercept $w_0$ can then be ignored.

# Polynomial curve fitting

- Data generated by a $Q$th order polynomial + some noise



- Predictor model:

$$f(x, \mathbf{w}) = w_1 x + w_2 x^2 + w_3 x^3 + \ldots\ldots + w_M x^M$$

$$= \sum_{i=1}^{M} w_i x^i$$

$$= \mathbf{w}^{\mathrm{T}} \phi$$

where $\mathbf{w} = \begin{bmatrix} w_1, \ldots, w_M \end{bmatrix}^{\mathrm{T}}$ and $\phi = \begin{bmatrix} x, \ldots, x^M \end{bmatrix}^{\mathrm{T}}$.

**Regularization**

# Polynomial curve fitting



- Complex hypotheses (richer class of models) lead to overfitting.

- A higher degree polynomial has more degrees of freedom which can lead to overfitting of the training data.

- Need to penalize the complexity in some way in the cost function.

*Figures are just for illustration.

**Regularization**

- Observations:

  − Weights $\mathbf{w}$ are unconstrained, and as such can lead to high variance.

  − Need to control the magnitude of the weights in order to control the variance.

- Modified objective:

$$\text{minimize} \quad \sum_{i=1}^{N} \left(y^{(i)} - \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}^{(i)})\right)^2 \quad \text{such that} \quad \sum_{j=1}^{M} w_j^2 \leq p$$

  − In vector form:

$$\text{minimize} \quad \left(\mathbf{y} - \Phi\mathbf{w}\right)^{\mathrm{T}}\left(\mathbf{y} - \Phi\mathbf{w}\right) \quad \text{such that} \quad ||\mathbf{w}||_2^2 \leq p$$

- Assumptions:

  $\Phi$ is standardized (zero mean and unit variance), and $\mathbf{y}$ is centered.

- Can show that the problem is equivalent to:

$$\text{minimize} \quad (\mathbf{y} - \Phi\mathbf{w})^{\mathrm{T}}(\mathbf{y} - \Phi\mathbf{w}) + \lambda||\mathbf{w}||_2^2$$

where $\lambda$ is the regularization coefficient.

- $\lambda$ tries to balance between the fit to the training data and the model complexity.

- Modified loss function:

$$L(\mathbf{w}) = L_E(\mathbf{w}) + \lambda L_R(\mathbf{w})$$

where

$$L_E(\mathbf{w}) = \sum_{i=1}^{N} \left(y^{(i)} - \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}^{(i)})\right)^2$$
$$= (\mathbf{y} - \Phi\mathbf{w})^{\mathrm{T}}(\mathbf{y} - \Phi\mathbf{w})$$

$$L_R(\mathbf{w}) = \sum_{j=1}^{M} w_j^2$$
$$= ||\mathbf{w}||_2^2$$
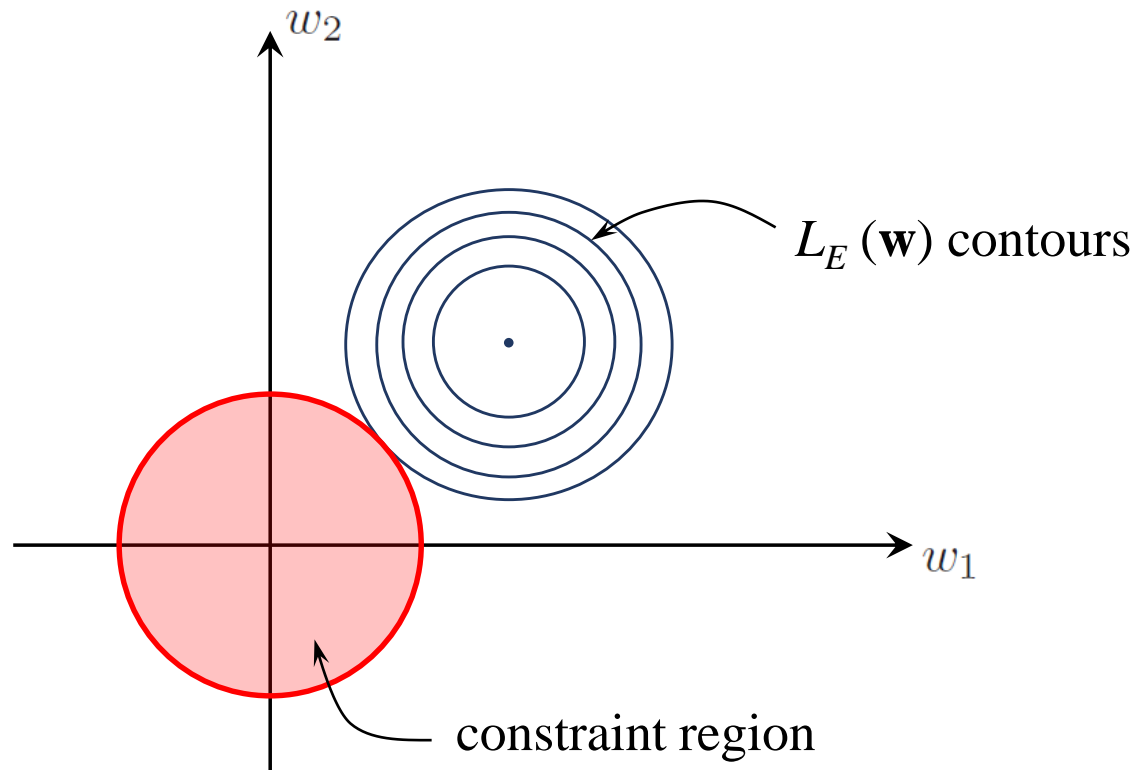
# L$_2$ regularization

- This is known as **L$_2$ Regularization** and also as **Ridge Regression**.

- Goal is to minimize the loss function $L(\mathbf{w})$. Note, since $L(\mathbf{w})$ is convex it has a unique solution.

- Taking derivative of $L(\mathbf{w})$ with respect to $\mathbf{w}$ and equating it to zero $\left(\text{i.e.} \dfrac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0\right)$ yields:

$$\mathbf{w} = \left(\Phi^{\mathrm{T}}\Phi + \lambda \mathbf{I}\right)^{-1}\Phi^{\mathrm{T}}\mathbf{y}$$

- If $\lambda = 0$, we get the least squares solutions.

- If $\lambda \to \infty$, we get $\mathbf{w} \to 0$.

- So $\lambda > 0$ will give weights of lower magnitudes than that obtained using least squares.

# L₁ regularization

- Use $L_1$ norm of the weight vector.

$$\text{minimize} \quad \left(\mathbf{y} - \Phi\mathbf{w}\right)^{\mathrm{T}}\left(\mathbf{y} - \Phi\mathbf{w}\right) \quad \text{such that} \quad \sum_{j=1}^{M} |w_i| \leq p$$

- Known as the **LASSO** (least absolute shrinkage and selection operator) algorithm (*Tibshirani*, 1996).

- **LASSO** has no closed form solution unlike ridge regression.

- Can be solved using quadratic programming techniques.

- Often want some of the weights $w_j$'s to be 0.

- **LASSO** looks for a sparse solution and so likely to yield some of the weights to be 0. But why?

# L₁ regularization

- Consider a problem with two features $x_1$ and $x_2$.

- In this case we are trying to solve a optimization problem with respect to weights $w_1$ and $w_2$:

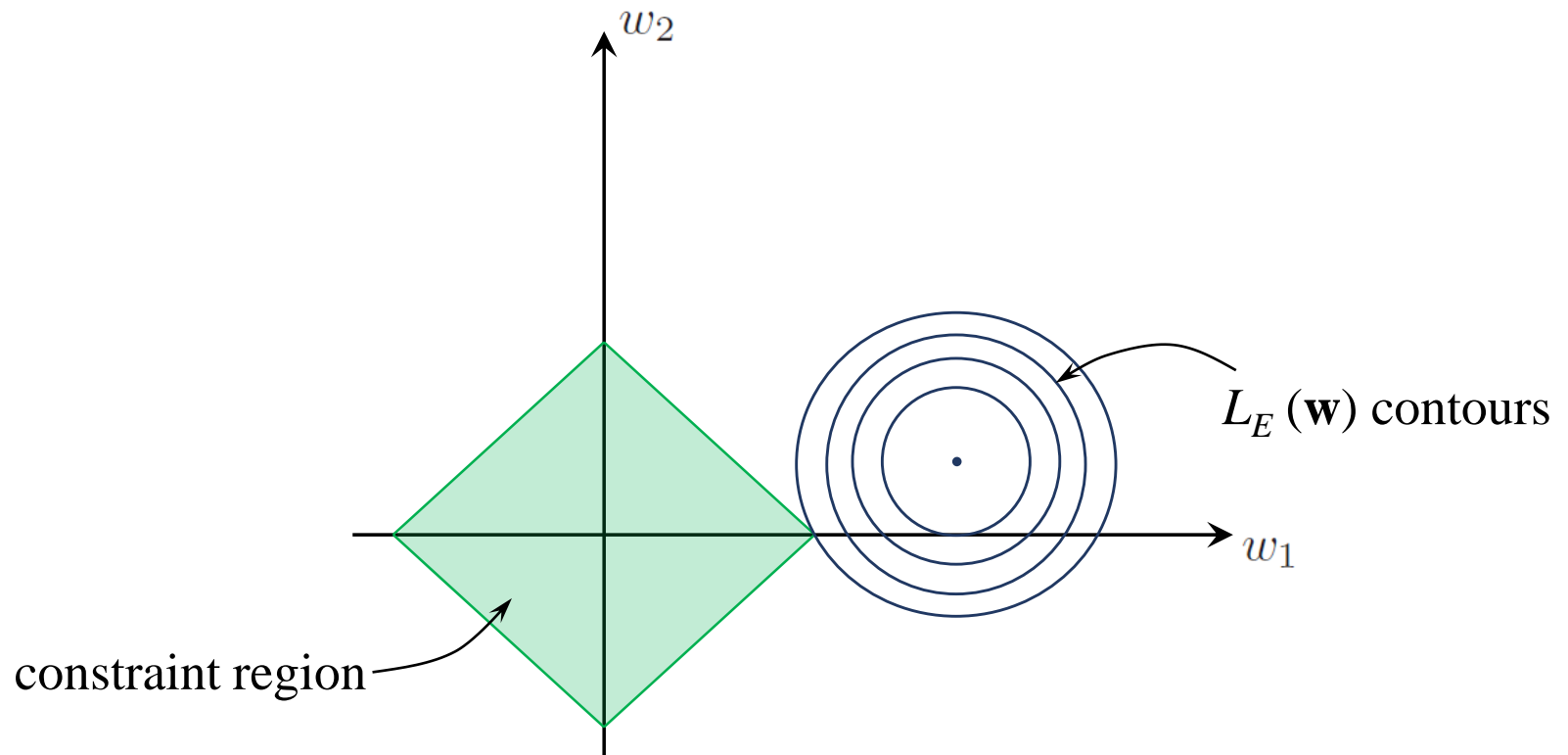$$\text{minimize} \quad \sum_{i=1}^{N} \left( y^{(i)} - w_1 x_1^{(i)} - w_2 x_2^{(i)} \right) \right)^2$$

such that

$$w_1 + w_2 \leq p$$

$$-w_1 + w_2 \leq p$$

$$w_1 - w_2 \leq p$$

$$-w_1 - w_2 \leq p$$

- When $\lambda$ is large, then among the contours satisfying the constraints, the contour with the least value of the objective function is likely to intersect the constraints' boundary at a corner.
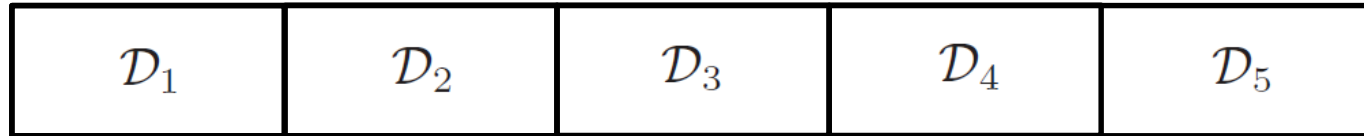
# *K*-fold cross validation

- Training data is subdivided into $K$ separate subsets – $\mathcal{D}_1, \mathcal{D}_2, ...., \mathcal{D}_K$ of equal size (say $n_K$).

- For $k = 1, 2, .., K$

  - Leave out the $k$th fold data $\mathcal{D}_k$ and train the model on the remaining $k - 1$ folds.

  - Use the trained model to make prediction on the $k$th fold data $\mathcal{D}_k$ and compute the (cross validation) error for this fold

$$E_k^{(\lambda)} = \frac{1}{n_K} \sum_{i=1}^{n_K} \left( y_{k,i} - f_{-k}^{(\lambda)}(\mathbf{x}_i) \right)^2$$

  where $f_{-k}^{(\lambda)}$ is the model trained excluding the $k$th fold data with a specific value of $\lambda$.

# K-fold cross validation

- Training data is subdivided into $K$ separate subsets – $\mathcal{D}_1, \mathcal{D}_2, ...., \mathcal{D}_K$ of equal size (say $n_K$). Let's take $K = 5$.

| $\mathcal{D}_1$ | $\mathcal{D}_2$ | $\mathcal{D}_3$ | $\mathcal{D}_4$ | $\mathcal{D}_5$ |
|---|---|---|---|---|

- Can generate $K$ training-test datasets using the $K$ subsets

Validation      Left out: $\mathcal{D}_1$

Validation      Left out: $\mathcal{D}_2$

Validation      Left out: $\mathcal{D}_3$

Validation      Left out: $\mathcal{D}_4$

Validation      Left out: $\mathcal{D}_5$

Regularization

# *K*-fold cross validation

- For $k = 1, 2, .., K$

    - Leave out the $k$th fold data $\mathcal{D}_k$ and train the model on the remaining $k - 1$ folds.

    

    - Use the trained model to make prediction on the $k$th fold data $\mathcal{D}_k$ and compute the (cross validation) error for this fold

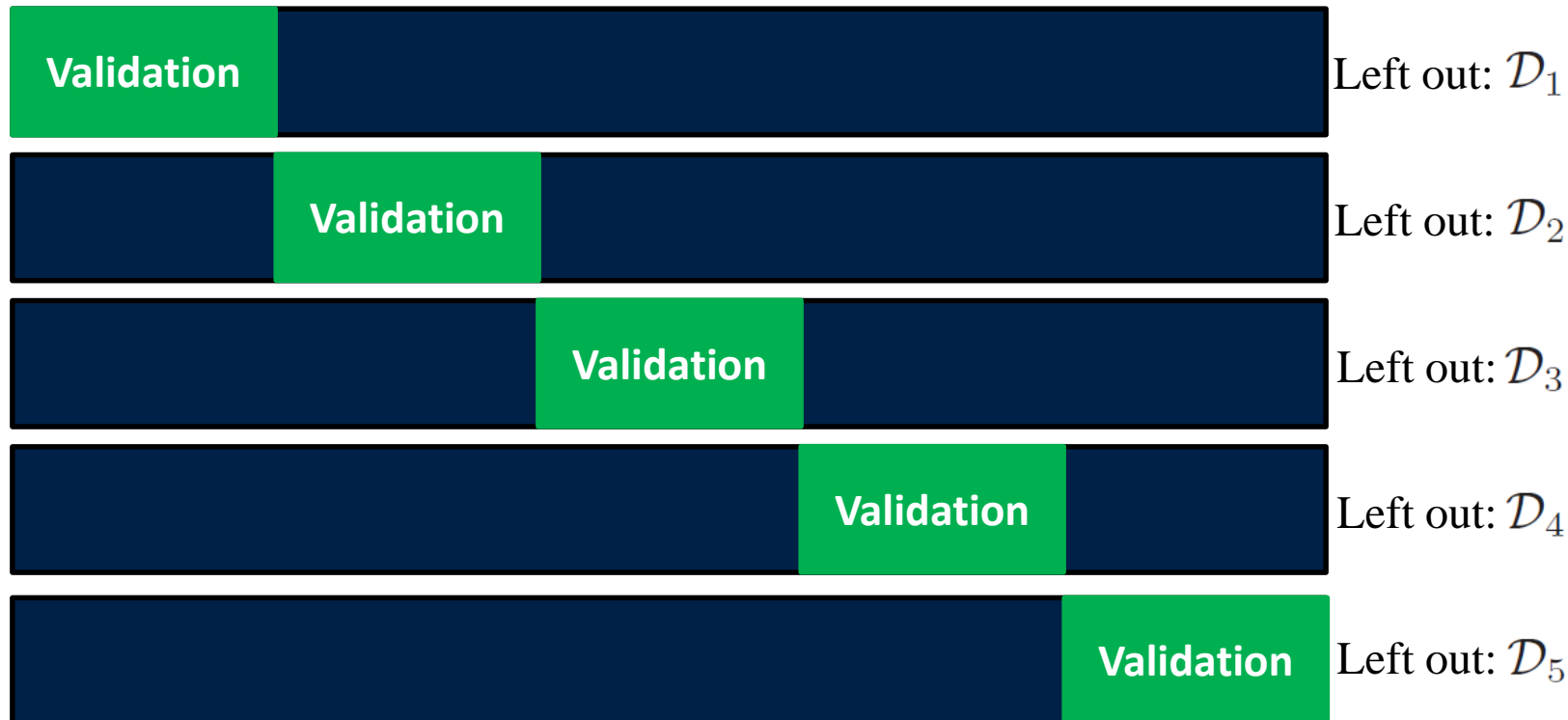    $$E_k^{(\lambda)} = \frac{1}{n_K} \sum_{i=1}^{n_K} \left( y_{k,i} - f_{-k}^{(\lambda)}(\mathbf{x}_i) \right)^2$$

    where $f_{-k}^{(\lambda)}$ is the model trained excluding the $k$th fold data with a specific value of $\lambda$.

# *K*-fold cross validation

- Estimated generalization error:

$$\mathbf{E}^{(\lambda)} = \frac{1}{K} \sum_{k=1}^{K} E_k^{(\lambda)}$$

- The optimal value of $\lambda$ (say $\lambda^*$) is the one yielding the least value of $\mathbf{E}^{(\lambda)}$.

- Using $\lambda^*$ train the model on the entire training dataset.

- When $K = N$ (size of the training dataset), the approach is known as leave-one-out cross-validation.