

Lecture 9: Generalization Gap & Bias-variance decomposition

- We introduced E_{new} \leftarrow error on fresh unseen data
- Designing a method with small E_{new} is the central goal in supervised learning
- Cross-validation helps in estimating E_{new}

(Unless one uses it for tuning hyperparameters)

- E_{new} can be further analyzed even more mathematically

error
function
 $E(y, \hat{y})$

Recall that $E_{\text{train}} = \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}(x_i; \tau))$

$E_{\text{new}} = \int E(y_*, \hat{y}(x_*; \tau)) p(x_*, y_*) dx_* dy_*$
multi-dimensional integral

- Note the values of E_{train} and E_{new} were calculated while keeping the training set τ fixed

- We now emphasize the fact that E_{train} and E_{new} are functions of the training set
 - As the training set \mathcal{T} changes, the values of E_{train} and E_{new} changes too!

function
of training
data

$$E_{\text{train}}(\mathcal{T}) = \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}(x_i; \mathcal{T}))$$
$$E_{\text{new}}(\mathcal{T}) = \int E(y_*, \hat{y}(x_*; \mathcal{T})) p(x_*, y_*) dx_* dy_*$$

- We now emphasize the fact that E_{train} and E_{new} are functions of the training set
 - As the training set \mathcal{T} changes, the values of E_{train} and E_{new} changes too!

function of training data

$$E_{\text{train}}(\mathcal{T}) = \frac{1}{N} \sum_{i=1}^N E(y^{(i)}, \hat{y}(x^{(i)}; \mathcal{T}))$$

$$E_{\text{new}}(\mathcal{T}) = \int E(y^*, \hat{y}(x^*; \mathcal{T})) p(x^*, y^*) dx^* dy^*$$

- So we introduce another level of abstraction
 - **Training-data averaged** versions of E_{new} and E_{train}

$$\bar{E}_{\text{new}} = \mathbb{E}_{\mathcal{T}} [E_{\text{new}}(\mathcal{T})]$$

$$\bar{E}_{\text{train}} = \mathbb{E}_{\mathcal{T}} [E_{\text{train}}(\mathcal{T})]$$

- $\mathbb{E}_{\mathcal{T}}$ denotes the expected value w.r.t. the training set $\mathcal{T} = \left\{ x^{(i)}, y^{(i)} \right\}_{i=1}^N$
- We assume \mathcal{T} consists of independent samples from $p(x, y)$

\bar{E}_{new} is the average E_{new} if we were to train the model multiple times on
on different training datasets of size N . (Same goes for E_{train})

$$\bar{E}_{\text{new}} = \mathbb{E}_{\mathcal{T}} [E_{\text{new}}(\mathcal{T})]$$

$$\bar{E}_{\text{train}} = \mathbb{E}_{\mathcal{T}} [E_{\text{train}}(\mathcal{T})]$$

- $\mathbb{E}_{\mathcal{T}}$ denotes the expected value w.r.t. the training set $\mathcal{T} = \left\{ \underline{x}^{(i)}, y^{(i)} \right\}_{i=1}^N$
- We assume \mathcal{T} consists of independent samples from $p(\underline{x}, y)$

\bar{E}_{new} is the average E_{new} if we were to train the model multiple times on different training datasets of size N .

— Why introduce these average quantities \bar{E}_{new} and \bar{E}_{train} ?

- It is easier to reason about the average behavior \bar{E}_{new} and \bar{E}_{train} than about the errors E_{new} and E_{train} obtained when the model is trained on one specific training dataset \mathcal{T}
- Insights from \bar{E}_{new} are useful, even though we most often care about E_{new} in the end (as the training data is usually fixed)

— We have already seen that $E_{\text{new}} \neq E_{\text{train}}$

• Usually, $E_{\text{train}} < E_{\text{new}}$ and $\bar{E}_{\text{train}} < \bar{E}_{\text{new}}$



On an average, a method usually performs worse on new, unseen data than on training data

— A method generalizes well if it performs well on unseen data after training

— We call the difference between \bar{E}_{new} and \bar{E}_{train} as the generalization gap

— Generalization gap = $\bar{E}_{\text{new}} - \bar{E}_{\text{train}}$ (more correctly as the expected generalization gap)

↙
is the difference between the expected performance on new unseen data and the expected performance on training data

— So \bar{E}_{new} can be decomposed as:

$$\bar{E}_{\text{new}} = \bar{E}_{\text{train}} + \text{generalization gap}$$

What affects the Generalization Gap?

- One can in some sense say, the more a method adapts to training data, the larger the generalization gap
- Vapnik-Chervonenkis (VC) dimension is a theoretical framework which assesses how much a method adapts to training data
 - The framework provides probabilistic bounds on generalization gap
 - However, the bounds are quite conservative
 - Therefore, the method is not very useful in practice
- Instead, we will only use vague terms "Model Complexity" (or Model flexibility)
 - ability of a method to adapt to patterns in training data
 - High complexity models — e.g. very deep decision trees, KNN with small 'k'
 - Low complexity models — e.g. linear regression, logistic regression

— Instead, we will only use vague terms “Model Complexity” (or Model flexibility)

ability of a method
to adapt to patterns
in training data

- High complexity models — e.g. very deep decision trees, KNN with small ‘K’
- Low complexity models — e.g. linear regression, logistic regression
basic parametric methods
- For parametric models, the model complexity is related to the number of parameters, but is also affected by regularization techniques

$$\text{Generalization gap} = \bar{E}_{\text{new}} - \bar{E}_{\text{train}}$$

Model complexity ↑

\bar{E}_{train} ↓

Generalization gap ↑

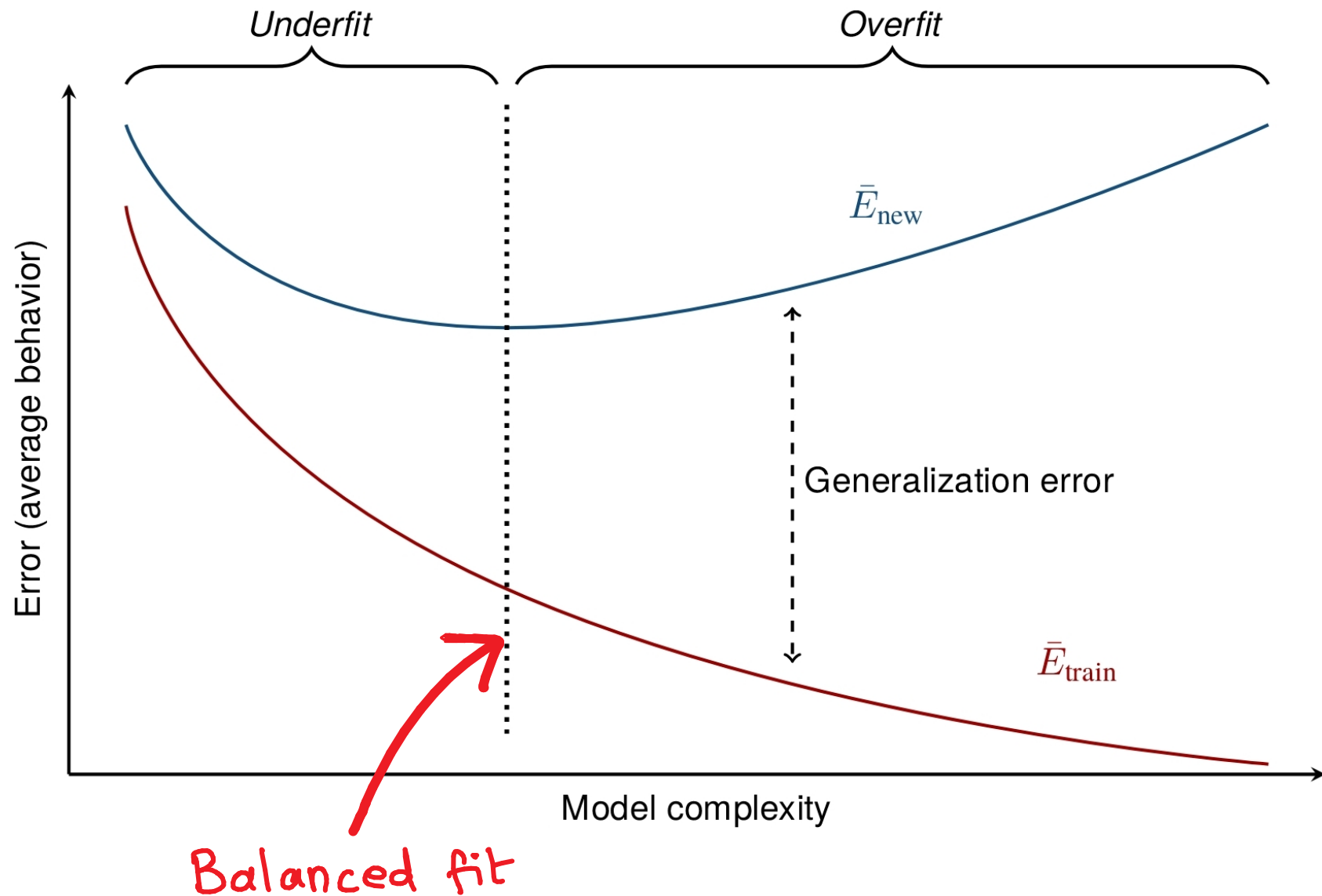
Model complexity ↓

\bar{E}_{train} ↑

Generalization gap ↓

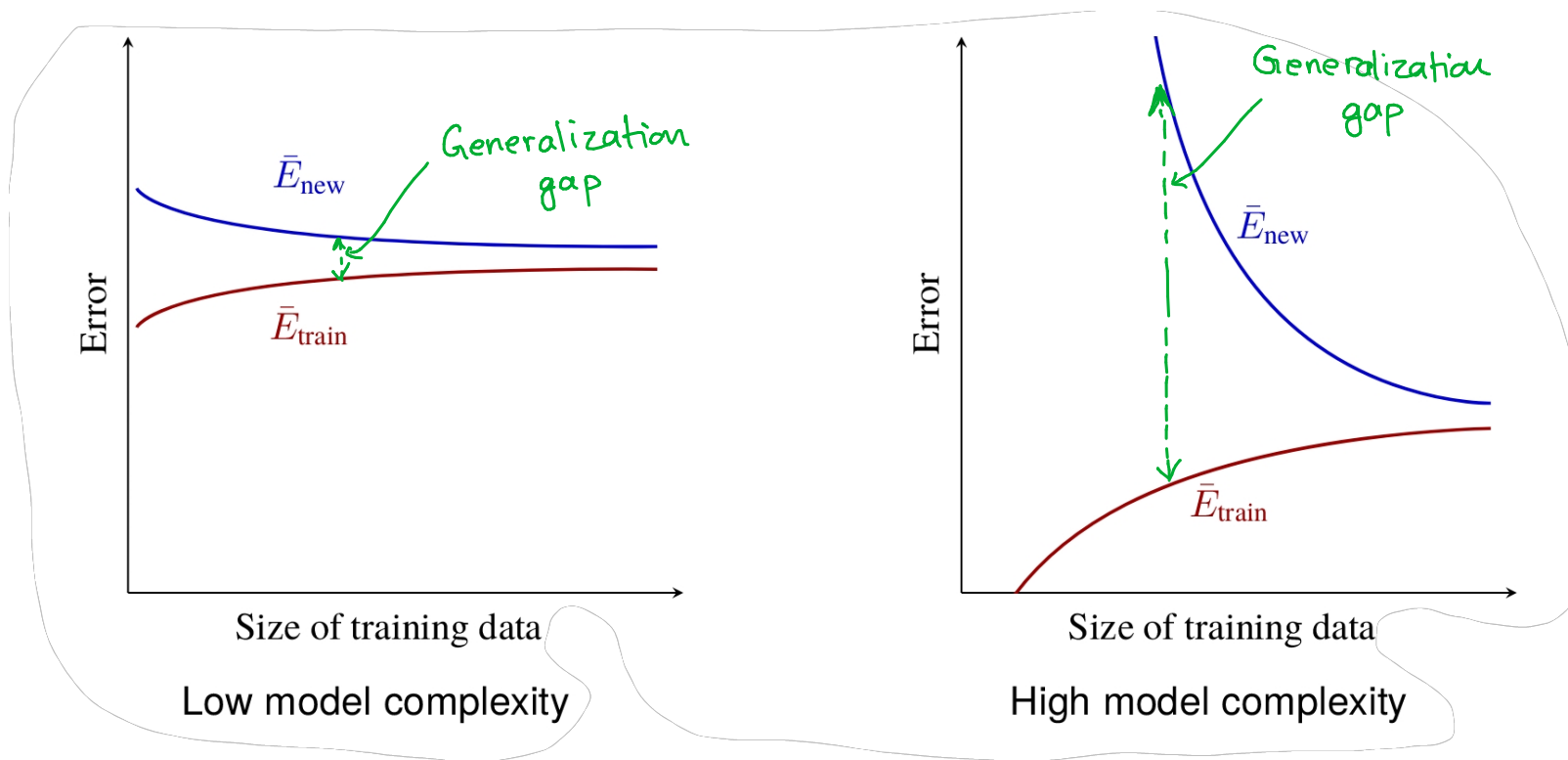
E_{new} usually attains a minimum at some intermediate complexity

Behavior of \bar{E}_{train} and \bar{E}_{new} with model complexity



When we use cross-validation to select hyperparameters, we search for balanced fit

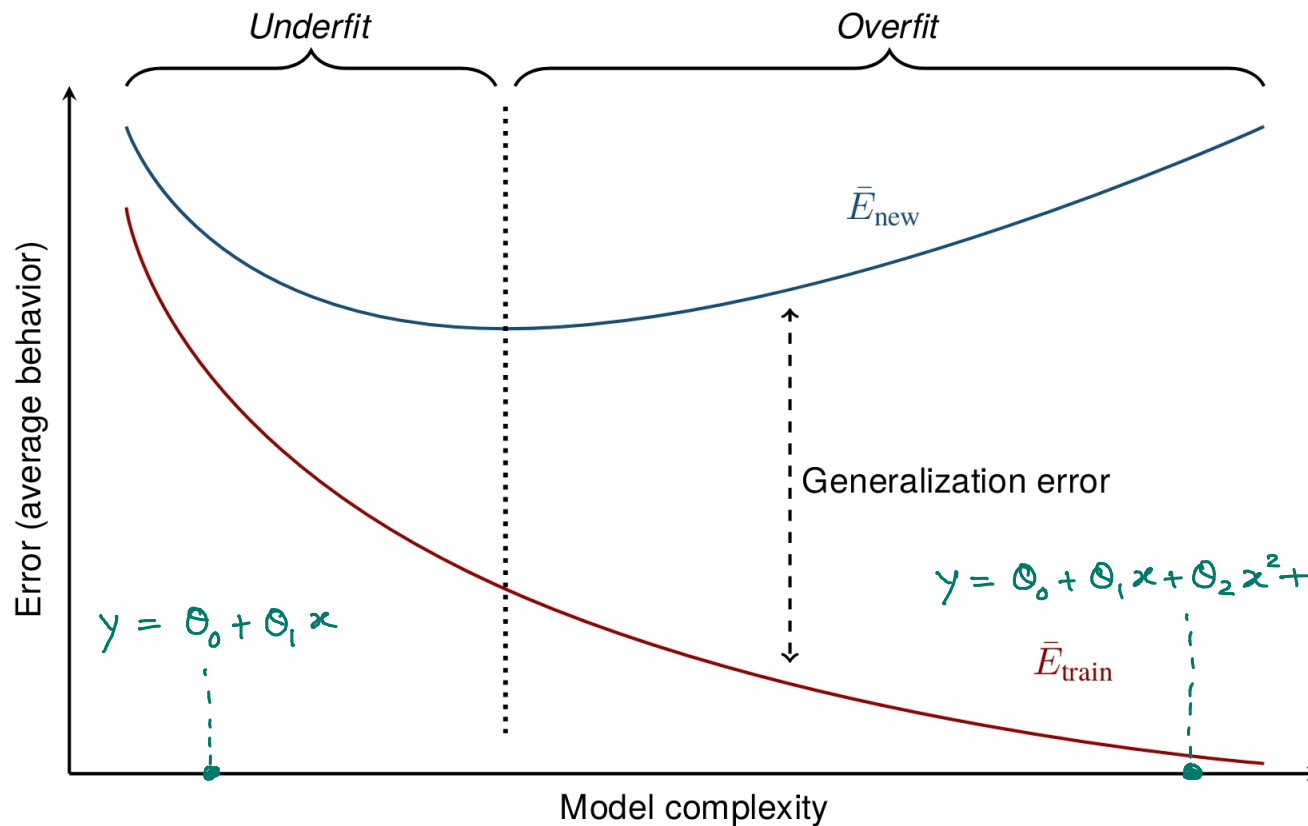
- Another important aspect is the **size of training dataset N**
- In general, the more the training data, the smaller the generalization gap
- \bar{E}_{train} typically increases as N increases
 - since most models are unable to fit all training data points well, more so if there are too many of them



- For large enough data, a more complex model will attain a smaller \bar{E}_{new} than a simpler model would

Shortcomings of using Model Complexity

- The figure below is relevant only when there is a single hyperparameter to choose



Say a hyperparameter to choose in polynomial regression is the order of polynomial

say kNN vs PolyReg

- However, when there are multiple hyperparameters (or multiple methods) to choose from, the above one-dimensional model complexity will not work well

- Let's take an example of jointly choosing the degree of polynomial regression and the regularization parameter

- higher degree of polynomial \Rightarrow more flexibility / complexity
- more regularization \Rightarrow less flexibility / complexity

- Example of a simulated problem: Sample data points from $p(x, y)$

Data
Generation

- $N = 10$ data points
- Input $x \sim \text{UniformDist}(-5, 10)$
- $y = \min(0.1x^2, 3) + \epsilon$
- $\epsilon \sim \text{NormalDist}(0, 1)$

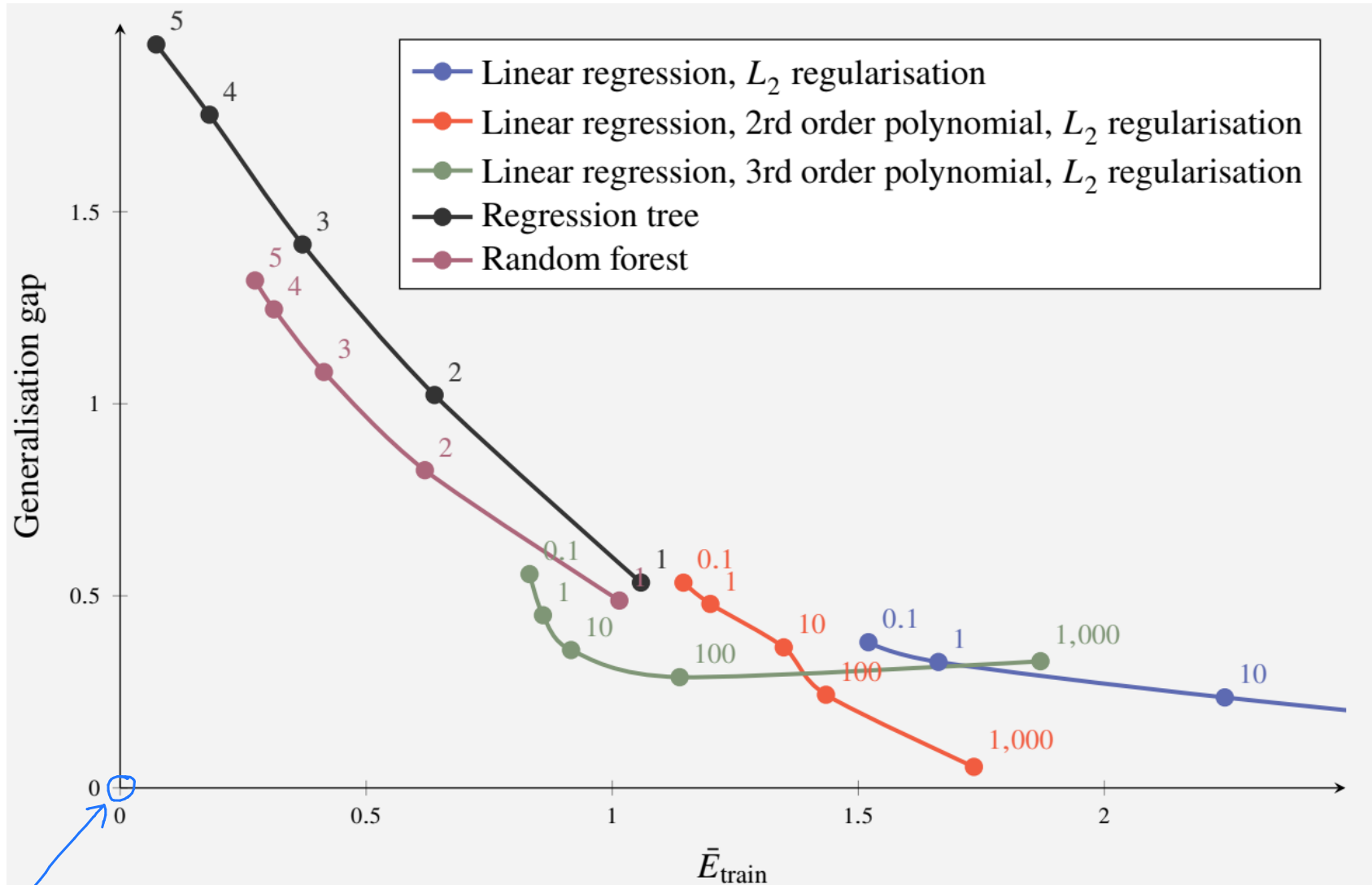
$$p(x, y) = \underbrace{p(y|x)} \underbrace{p(x)}$$

Then sample y given x \leftarrow First sample x

Now fit the input-output data $\{x^{(i)}, y^{(i)}\}_{i=1}^{10}$ using

- (a) Linear regression with L_2 -regularization
- (b) Linear regression with a quadratic polynomial and L_2 -regularization
- (c) Linear regression with a cubic polynomial and L_2 -regularization
- (d) Regression Tree, (e) A random forest with 10 regression trees

- For each of the methods, we try different values of hyperparameters (regularization parameter λ and tree depth) and compute \bar{E}_{train} and generalization gap



In practice,
we only
have limited
data and
we cannot
generate
these plots!