

Machine Learning

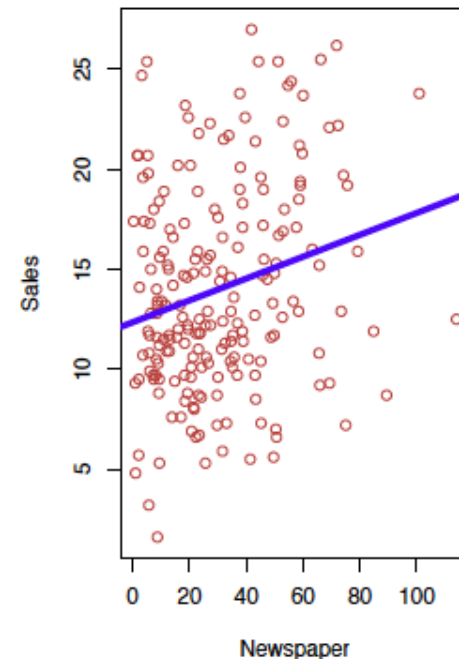
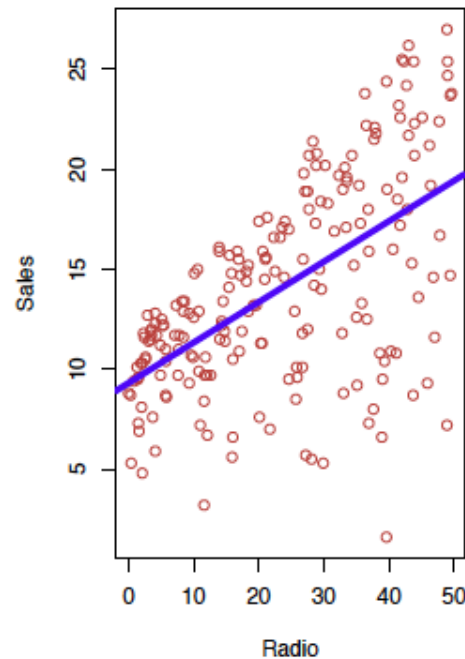
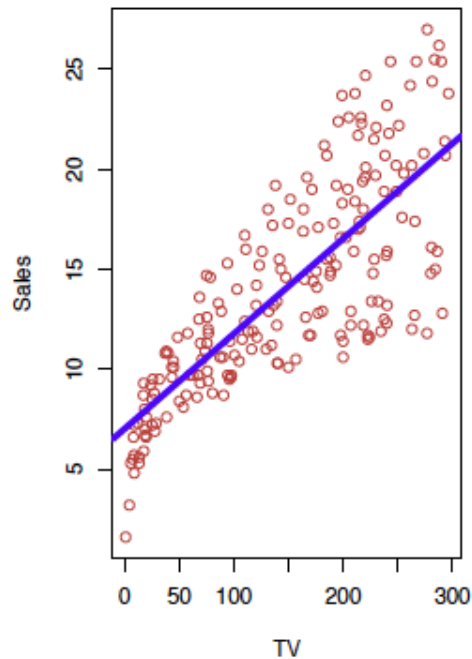
Supervised Machine Learning – Regression

Part 1: Parametric and Non-Parametric Approaches

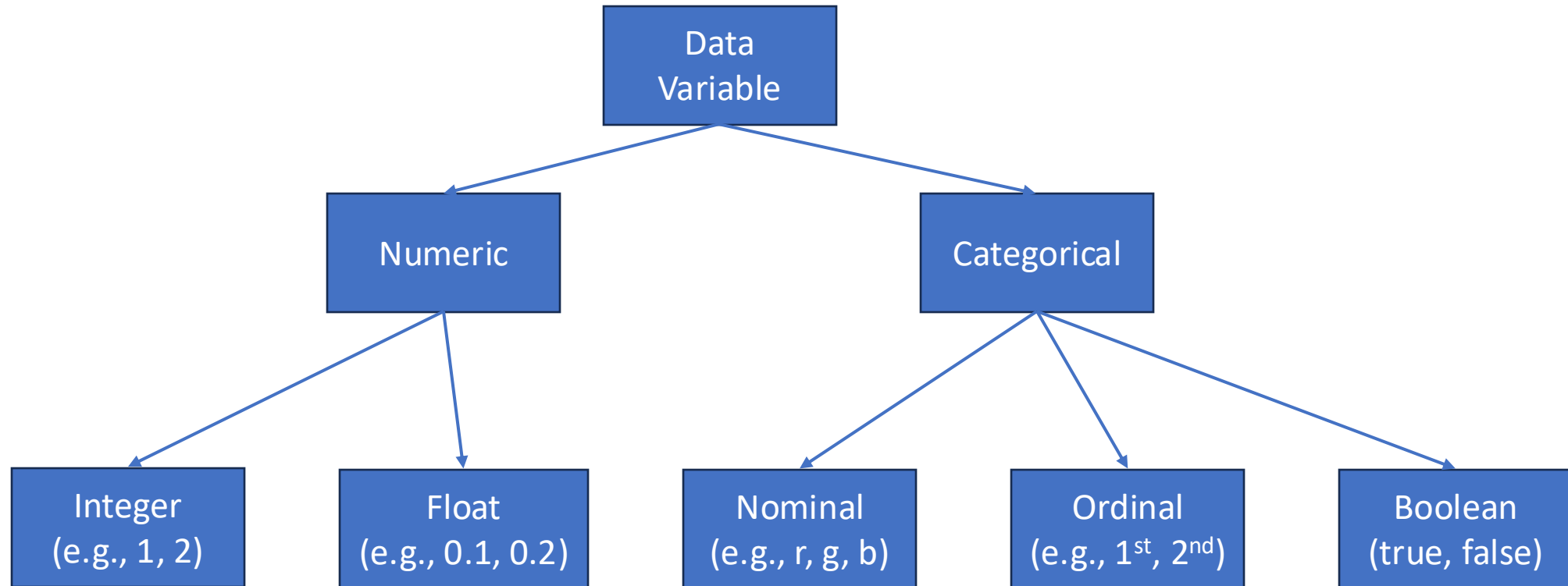
Regression, Data, Error, Performance

Prediction

- ❑ Shown are Sales vs amount spent on TV, Radio and Newspaper ads, with a blue linear-regression line fit separately to each.
- ❑ Can we predict Sales using these three?
- ❑ Perhaps we can do better using a model that combines all three: $\text{Sales} \approx f(\text{TV}, \text{Radio}, \text{Newspaper})$



Numerical Data



- In regression, inputs and outputs are all numerical.
- In classification, output can be categorical (but in practice will be mapped to numerical)

Regression - Data Model

- In this example, Sales is a response/target/output; we denote it by y . TV is a feature/input/predictor; we denote it by x_1 . Likewise, we denote Radio by x_2 , and so on.
- In general, we can refer to the input vector as $\mathbf{x} \in I$, containing d input features ($I \subseteq \mathbb{R}^d$ is the “input space”).
- We assume that there exists function f so that:

$$y(\mathbf{x}) = f(\mathbf{x}) + \epsilon \in \mathbb{R}_+$$

where ϵ is a zero-mean random variable (independent of \mathbf{x}) and captures measurement errors and noise.

- $f(\mathbf{x})$ is the unknown true model that describes input-output for our problem.

Regression - Goal

What is our goal here?

Our ideal goal is to design model $\hat{f}(x)$ that estimates the unknown $f(x)$ accurately as possible.

To do that, we will use multiple input/output examples.

The prediction of $y(x) = f(x) + \epsilon$ by $\hat{f}(x)$ is called **regression**.

Regression - Prediction Error


How good is my prediction?

True output: $y(\mathbf{x}) = f(\mathbf{x}) + \epsilon$

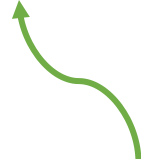
Prediction: $\hat{f}(\mathbf{x}) \approx f(\mathbf{x})$

Prediction error: $e = y(\mathbf{x}) - \hat{f}(\mathbf{x}) = \boxed{f(\mathbf{x}) - \hat{f}(\mathbf{x})} + \boxed{\epsilon}$

Model estimation error
(reducible)



Noise error (irreducible)



Regression - Prediction Error (cont'd)

Prediction error: $e = y(\mathbf{x}) - \hat{f}(\mathbf{x}) = (f(\mathbf{x}) - \hat{f}(\mathbf{x})) + \epsilon$

Model estimation error: $f(\mathbf{x}) - \hat{f}(\mathbf{x})$

- This is reducible.
- Our mission in ML is to make this as low as possible. There are many ways that we'll discuss.

Noise error: ϵ

- Irreducible. Nothing we can do about it.
- *The performance of any model is limited by the quality of the data it operates on.*

Regression - Model Performance

Prediction error $e = y(\mathbf{x}) - \hat{f}(\mathbf{x}) = (f(\mathbf{x}) - \hat{f}(\mathbf{x})) + \epsilon$ is a random variable.

- Input \mathbf{x} is random.
- Error ϵ is random
- We assume that f is a deterministic function (same input gives same output). The same for \hat{f} .

Even for given input \mathbf{x} , e is random due to ϵ . We need to calculate the statistics of e .

Regression - Model Performance

For any given input, the error is a random variable with some mean and variance.

Zero mean and zero variance would be ideal (deterministic error zero). But what is feasible?

Mean: $E[e|\mathbf{x}] = (f(\mathbf{x}) - \hat{f}(\mathbf{x}))$

Variance: $Var[e|\mathbf{x}] = E[(e - E[e|\mathbf{x}])^2|\mathbf{x}] = E[\epsilon^2|\mathbf{x}] = Var[\epsilon] =: \sigma_\epsilon^2$

Mean Squared Error: $E[e^2|\mathbf{x}] = E\left[\left(y - \hat{f}(\mathbf{x})\right)^2|\mathbf{x}\right]$

$$MSE = E\left[\left((f(\mathbf{x}) - \hat{f}(\mathbf{x})) - \epsilon\right)^2|\mathbf{x}\right] = E\left[\left(f(\mathbf{x}) - \hat{f}(\mathbf{x})\right)^2|\mathbf{x}\right] + E[\epsilon^2|\mathbf{x}] - E\left[\epsilon(f(\mathbf{x}) - \hat{f}(\mathbf{x}))|\mathbf{x}\right]$$

$$= \left(f(\mathbf{x}) - \hat{f}(\mathbf{x})\right)^2 + \sigma_\epsilon^2$$

If you find these derivations challenging, you should review the provided notes on probability and random variables.

Reducible

Irreducible (given data quality)

Parametric Modeling

Parametric Modeling

Goal: Train model \hat{f} such that

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}) \Leftrightarrow \left(f(\mathbf{x}) - \hat{f}(\mathbf{x})\right)^2 \approx 0 \quad \forall \mathbf{x} \in I$$

Basic assumption: True function f has a **parametric** form $f(\mathbf{x}) = m(\mathbf{x}; \mathbf{a})$ for some finite-size parameter vector \mathbf{a} in parameter space H .

Example:

$$f(\mathbf{x}) = a_5 x_1^2 + a_4 x_2^2 + a_3 x_1 x_2 + a_2 x_1 + a_1 x_2 + a_0$$

- The parametric model m is the bivariate 2nd-degree polynomial.
- The parameters are in parameter vector $\mathbf{a} = [a_0, a_1, a_2, a_3, a_4, a_5]^T \in H = \mathbb{R}^6$.

Training the Parametric Model

Assumption: $f(\mathbf{x}) = m(\mathbf{x}; \mathbf{a})$

Training parametric model $\hat{f} \approx f$ breaks in two steps:

- 1) Identify m
- 2) Identify \mathbf{a}

Training the Parametric Model

Assume the simplified case that true m is known.

To train \hat{f} is to train the parameters of the model $\mathbf{a} \in H$.

To do that, we rely on a set of training input/output data

$$S = \{(y_n, \mathbf{x}_n)\}_{n=1}^N$$

We train $\hat{f}(\mathbf{x}) = m(\mathbf{x}; \hat{\mathbf{a}})$, by solving

$$\hat{\mathbf{a}} = \operatorname{argmin}_{\mathbf{b} \in H} L(\mathbf{b}; S)$$

L is the loss function: It measures how well $m(\mathbf{x}; \hat{\mathbf{a}})$ represents the available training data.

MSE Training

Recall performance metric:

$$\text{MSE} = E \left[\left(y(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 | \mathbf{x} \right] = E[(y - m(\mathbf{x}; \hat{\mathbf{a}}))^2 | \mathbf{x}]$$

Mean can be estimated by sample-average.

In this case, we can use the training data: $\widehat{\text{MSE}} = \frac{1}{N} \sum_{n=1}^N (y_n - m(\mathbf{x}_n; \hat{\mathbf{a}}))^2$

MSE training:

$$\hat{\mathbf{a}} = \operatorname{argmin}_{\mathbf{b} \in H} L(\mathbf{b}; S) = \operatorname{argmin}_{\mathbf{b} \in H} \frac{1}{N} \sum_{n=1}^N (y_n - m(\mathbf{x}_n; \mathbf{b}))^2$$

Parameter Training Considerations

The quality of my parameter training is determined by 3 variables.

- Number of parameters.
- Number of available training data N .
- Noise variance in training data σ_{ϵ}^2 .

Good training for more model parameters in m needs more training data of better quality.

If m is known, better parameter estimates \hat{a} means better model \hat{f} .

Example: Univariate Polynomial Model

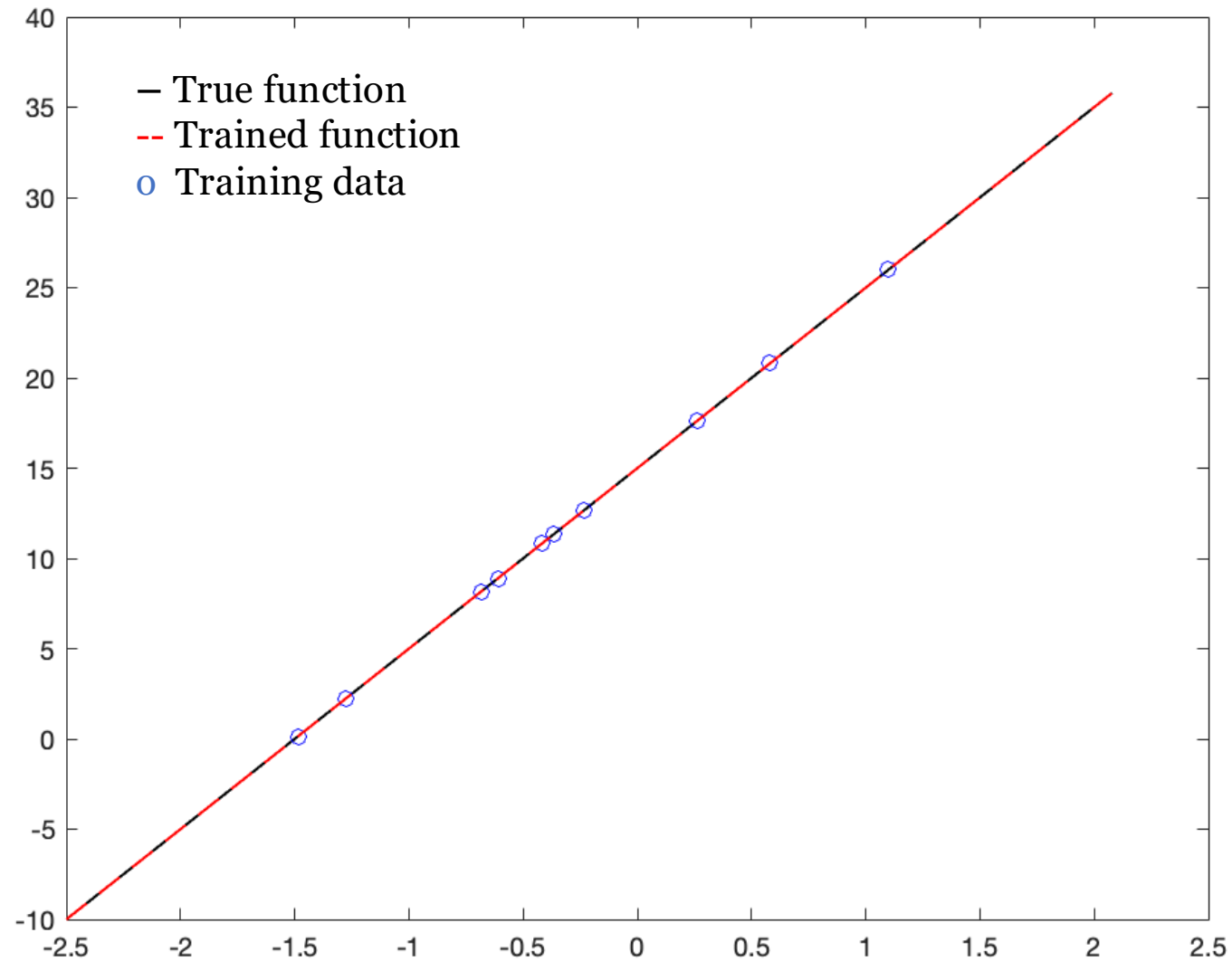
For simplicity, assume 1-D input and polynomial f .

- Constant (degree 0): $f(x) = a_0$ (no interest)
- Linear (degree 1): $f(x) = a_1x + a_0$
- Quadratic (degree 2): $f(x) = a_2x^2 + a_1x + a_0$
- Cubic (degree 3): $f(x) = a_3x^3 + a_2x^2 + a_1x + a_0$

For degree- k polynomial, we need to estimate $k + 1$ parameters $\mathbf{a} = [a_0, a_1, \dots, a_k]^T$.

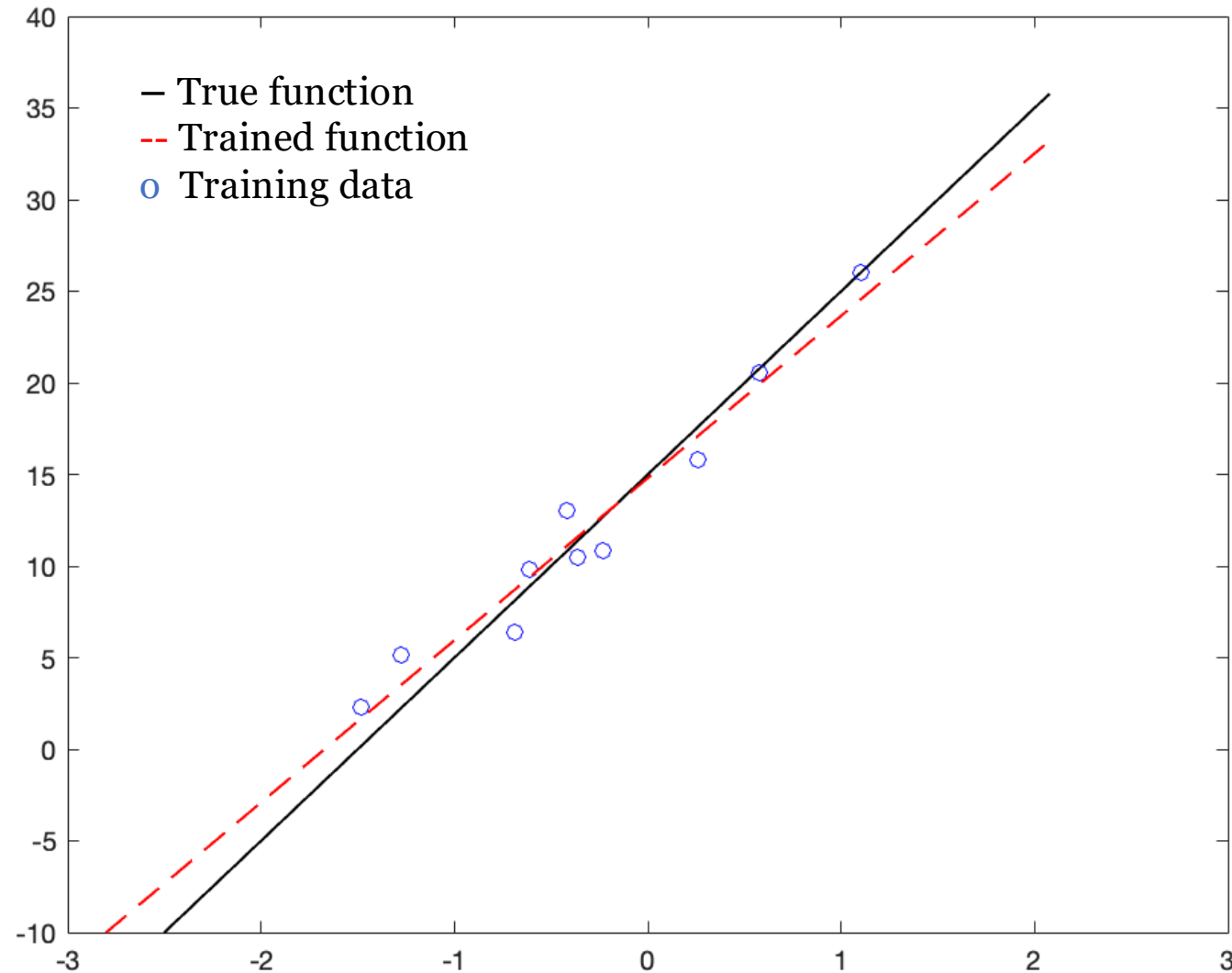
Example: Varying k, N, σ_ϵ^2

- $k = 1$
- $N = 10$
- $\sigma_\epsilon^2 = 0$



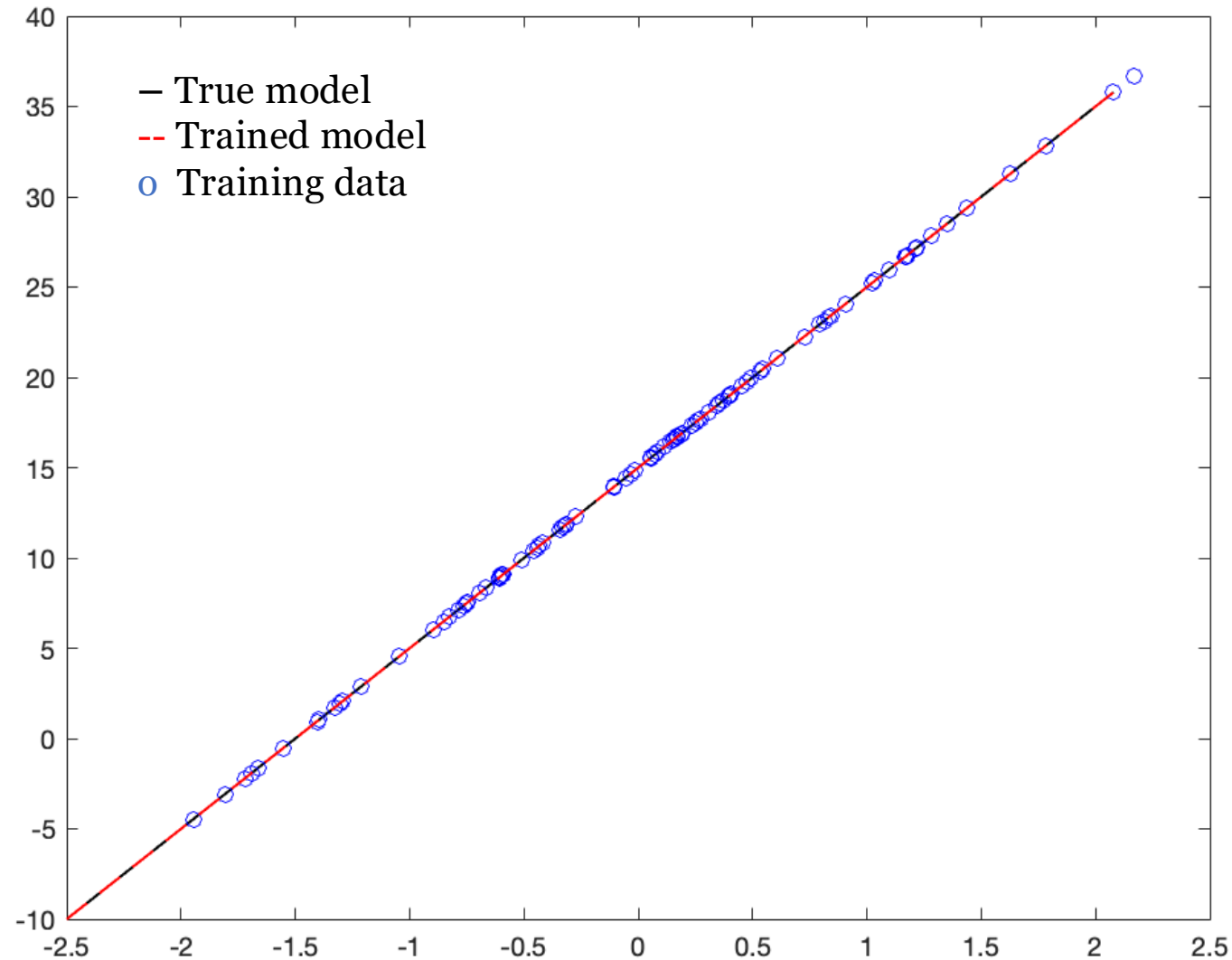
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 1$
- $N = 10$
- $\sigma_\epsilon^2 = 4$



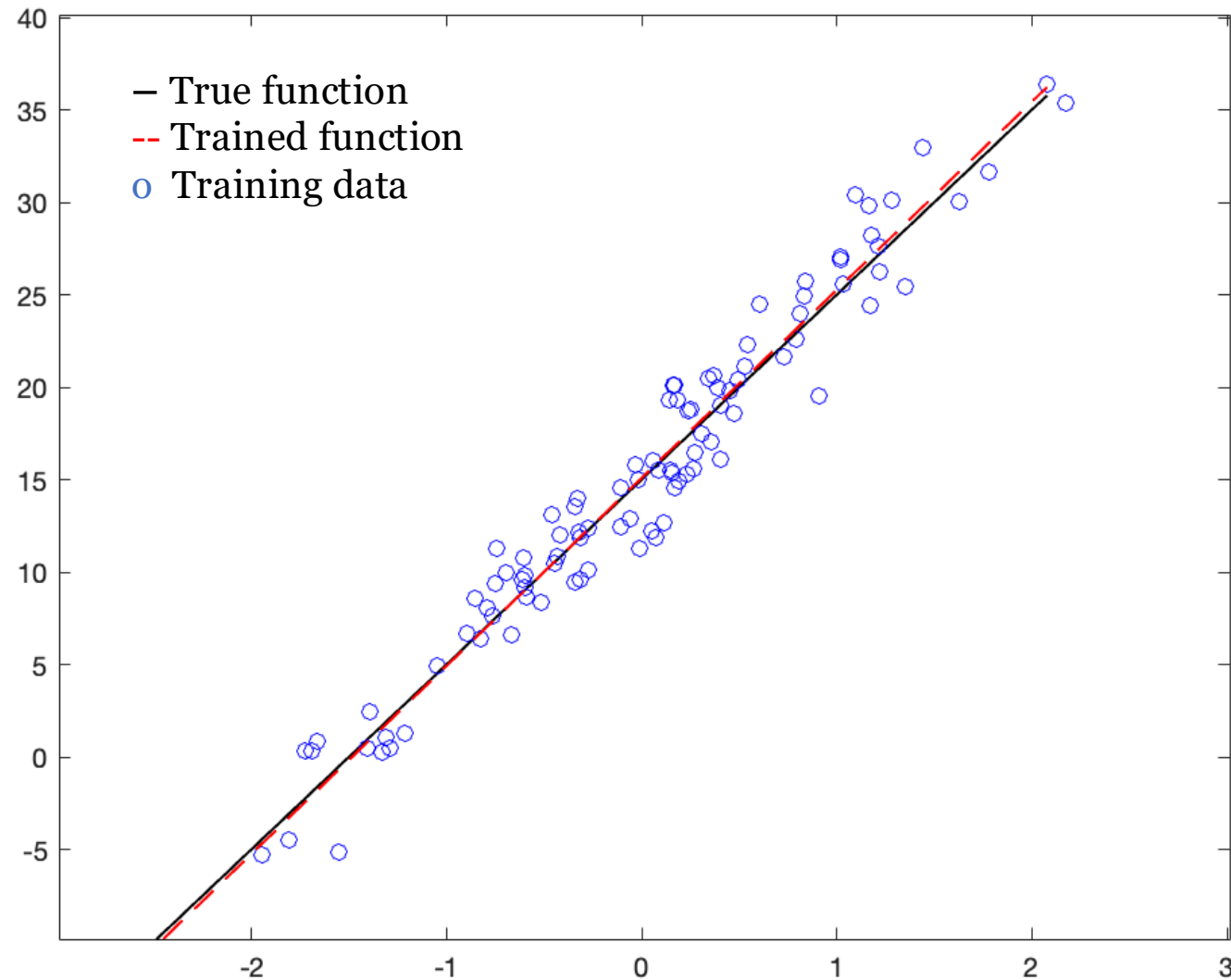
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 1$
- $N = 100$
- $\sigma_\epsilon^2 = 0$



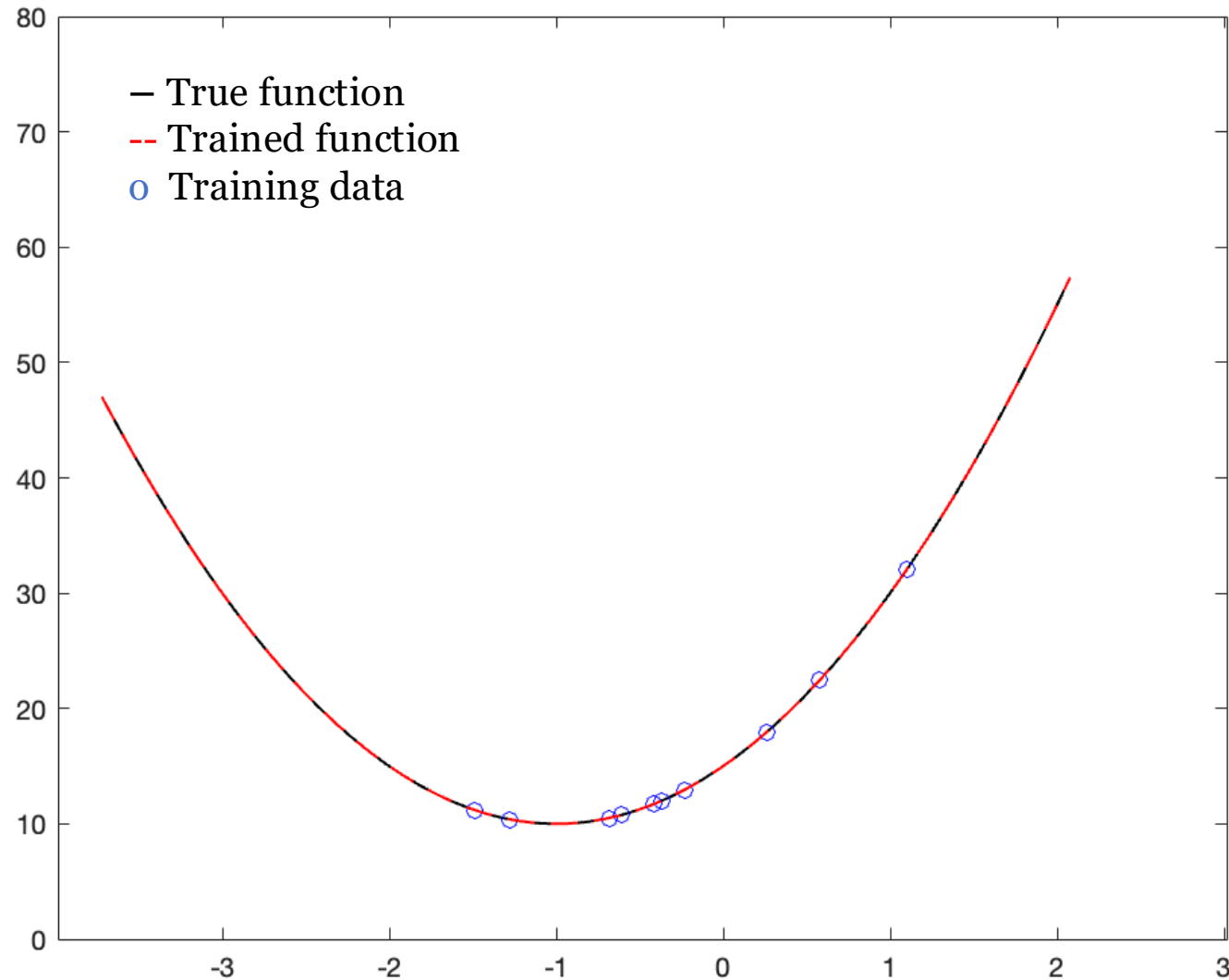
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 1$
- $N = 100$
- $\sigma_\epsilon^2 = 4$



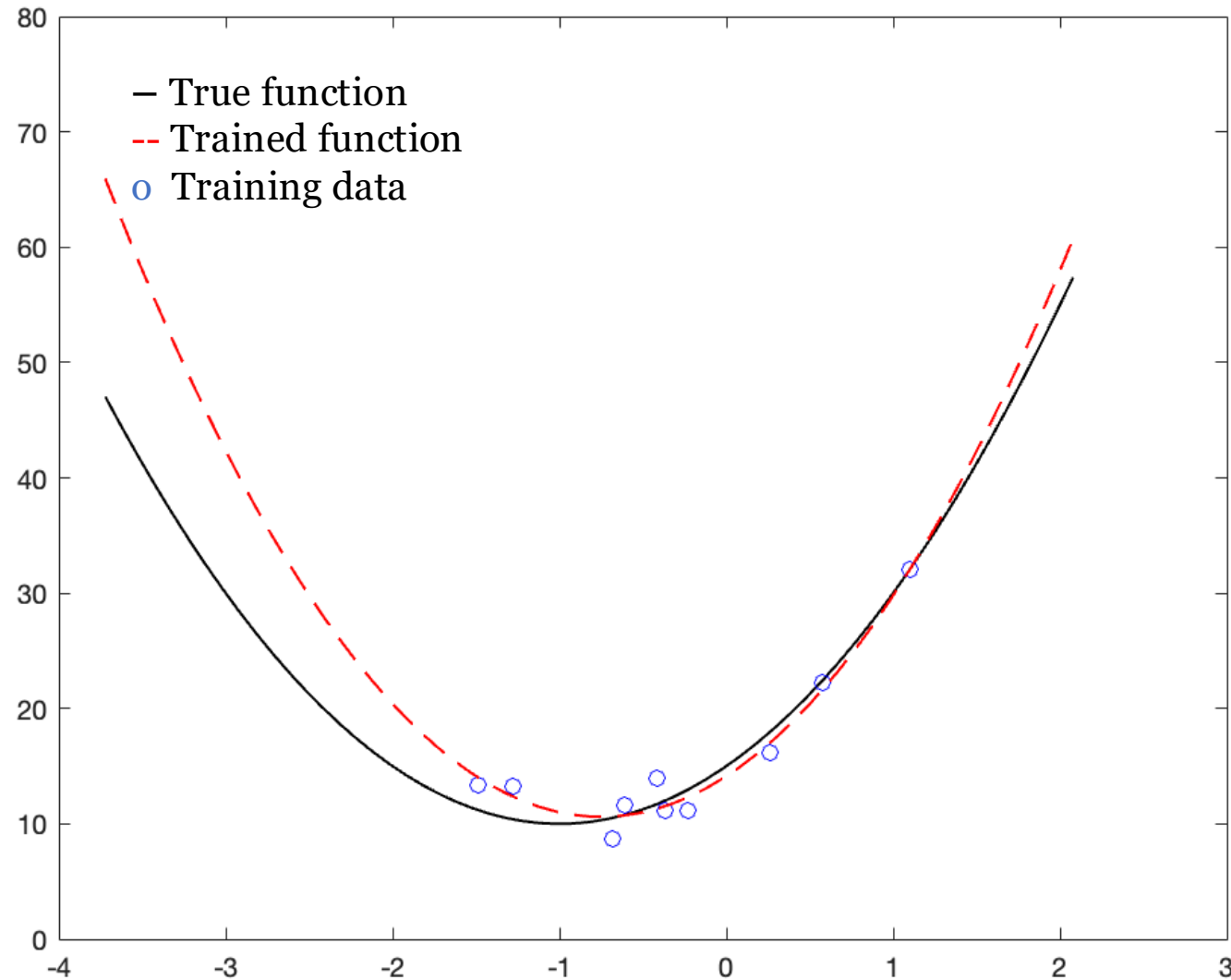
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 2$
- $N = 10$
- $\sigma_\epsilon^2 = 0$



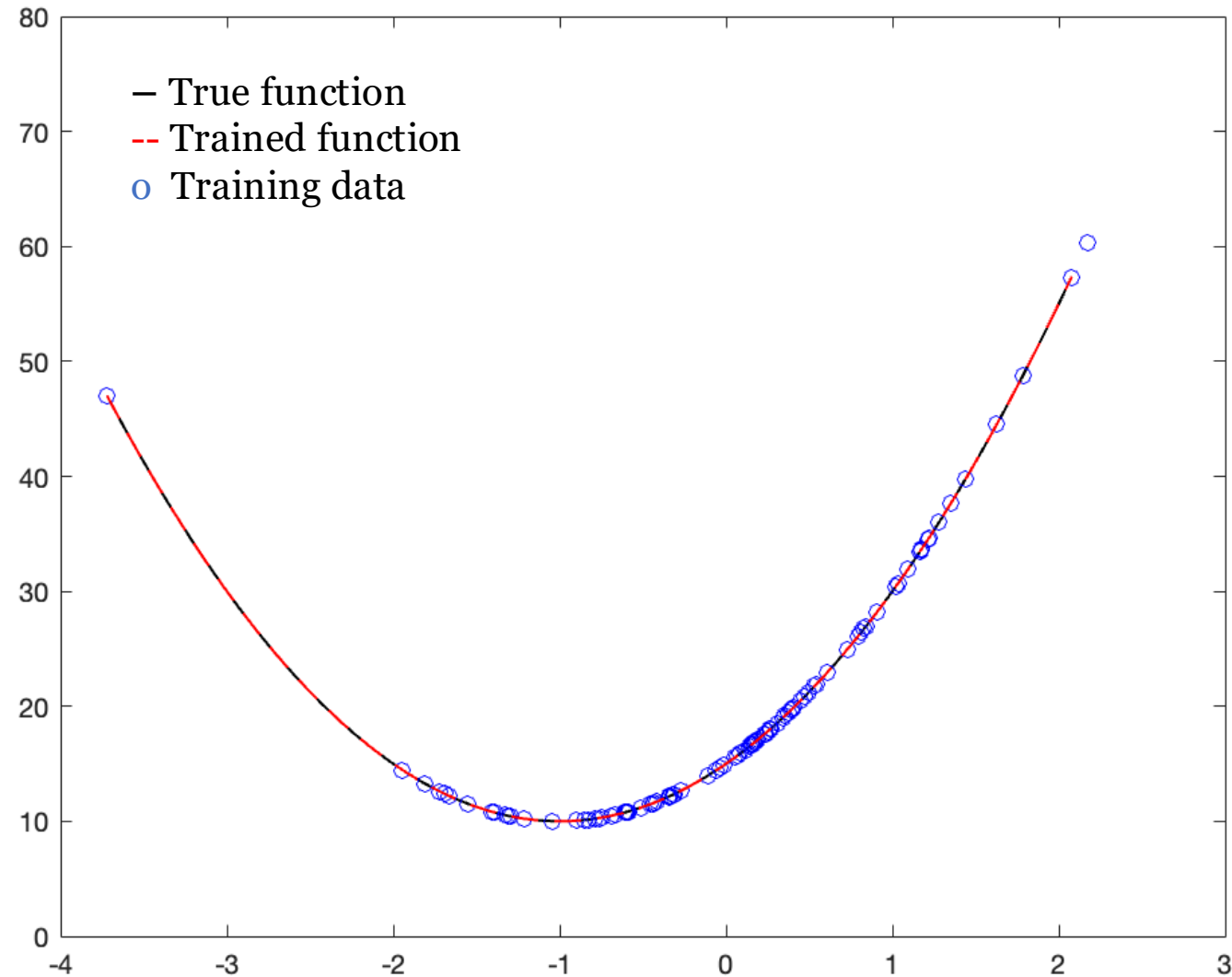
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 2$
- $N = 10$
- $\sigma_\epsilon^2 = 4$



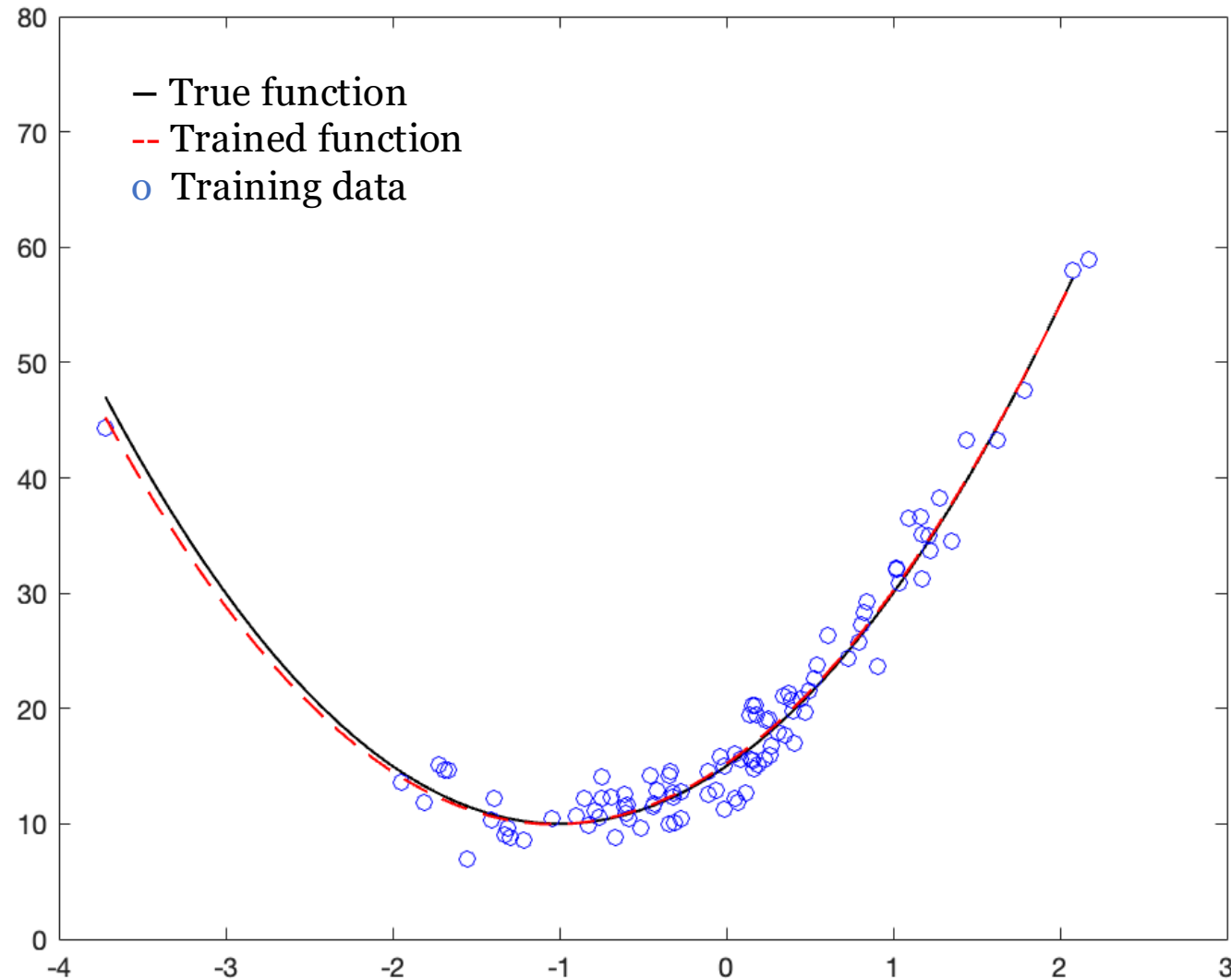
Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 2$
- $N = 100$
- $\sigma_\epsilon^2 = 0$



Example: Varying k, N, σ_ϵ^2 (cont'd)

- $k = 2$
- $N = 100$
- $\sigma_\epsilon^2 = 4$



Example: Varying k, N, σ_ϵ^2 (cont'd)

- In all cases, $\sigma_\epsilon^2 = 0$ resulted $\hat{f} = f$, for all tested values of k and N .
- For $\sigma_\epsilon^2 > 0$, for all k , higher N resulted in better training.
- For $\sigma_\epsilon^2 > 0$, for any given value of N , we attained better training for lower k .

Parametric Model Training – Open Questions

Basic assumption: $f(x) = m(x; a)$

Training parametric model \hat{f} breaks in two steps:

1) Identify m

If unknown (common), estimate \hat{m} .

How?

2) Identify a

Train \hat{a} .

What loss? How do we solve the fitting problem?

What are the trade-offs between estimation of \hat{m} and training of \hat{a} ?

We will return to this discussion...

Non-Parametric Modeling

Non-Parametric Regression

Recall, our goal is to train \hat{f} such that, for all \mathbf{x} in input space,

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}) \quad \forall \mathbf{x} \in I$$

Recall that $y(\mathbf{x}) = f(\mathbf{x}) + \epsilon$ with $E[\epsilon] = 0$. That is,

$$f(\mathbf{x}) = E[y(\mathbf{x})] \quad \forall \mathbf{x} \in I$$

Therefore, I want to train \hat{f} such as

$$\hat{f}(\mathbf{x}) \approx E[y(\mathbf{x})] \quad \forall \mathbf{x} \in I$$

- We do not have $E[y(\mathbf{x})]$ for any \mathbf{x} .
- We just have $y(\mathbf{x})$ for some values of \mathbf{x} (in S).

Nearest-Neighbor Regression

Idea: Estimate mean $E[y(\mathbf{x})]$ by averaging values of $y(\mathbf{x})$.

Coherent samples: For any $\mathbf{x} \in I$, to estimate $E[y(\mathbf{x})]$, I need to have many instances of $y(\mathbf{x}) = f(\mathbf{x}) + \epsilon$, or this same \mathbf{x} .

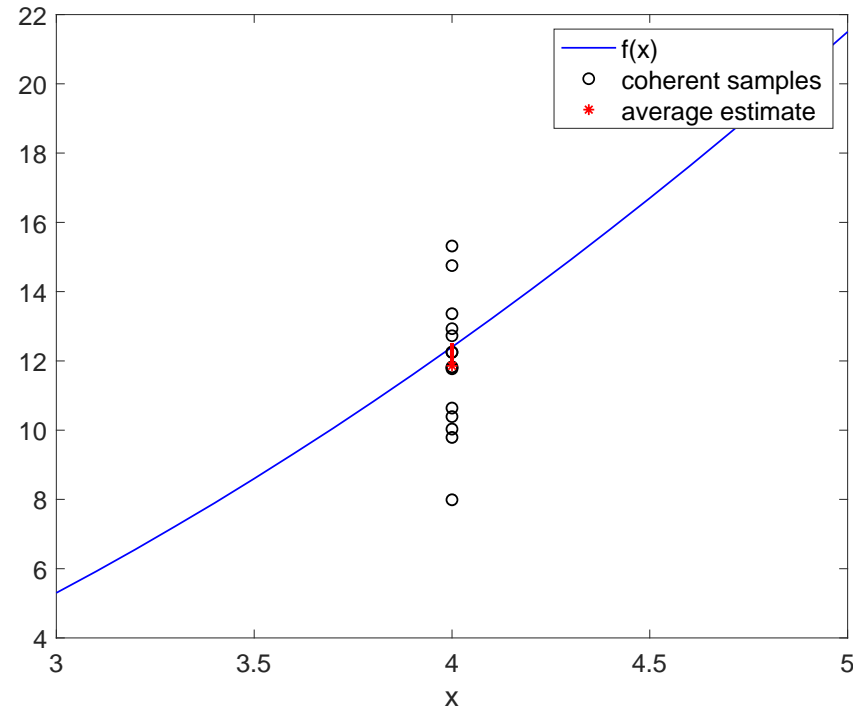
How many distinct instances $(y(\mathbf{x}), \mathbf{x})$ do we have in S for each $\mathbf{x} \in I$?

In S , we only have some/few/one pairs $(y(\mathbf{x}), \mathbf{x})$ for some values of \mathbf{x} .

Idea: I will relax coherence constraint. I will estimate $E[y(\mathbf{x})]$ by averaging $y(\mathbf{z})$ for all available \mathbf{z} in “the neighborhood” of \mathbf{x} .

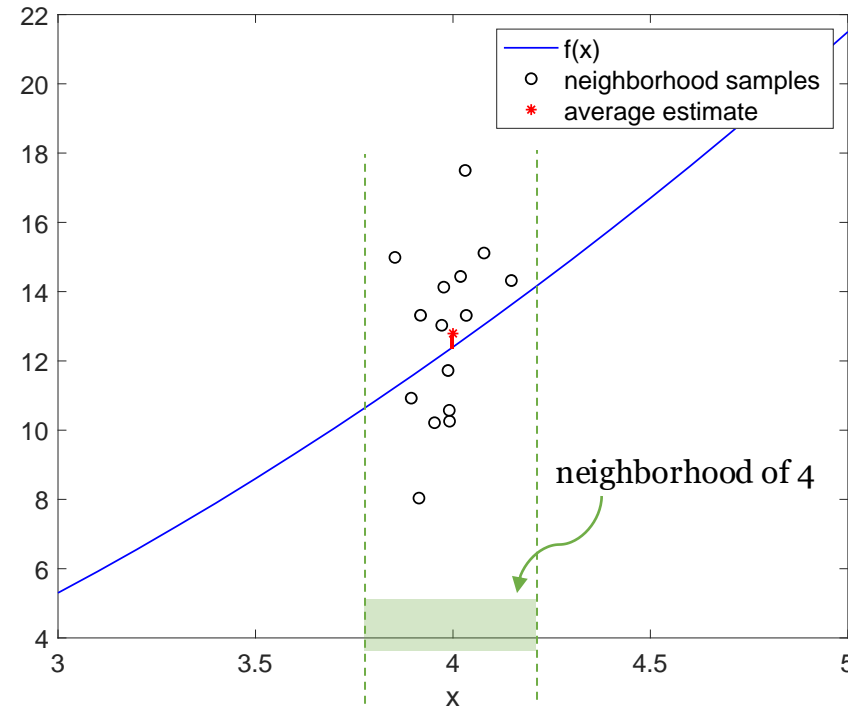
Nearest-Neighbor Regression (cont'd)

Slow varying function. Average estimation on coherent samples $y(4) = f(4) + \epsilon$, $\sigma_\epsilon = 2$.



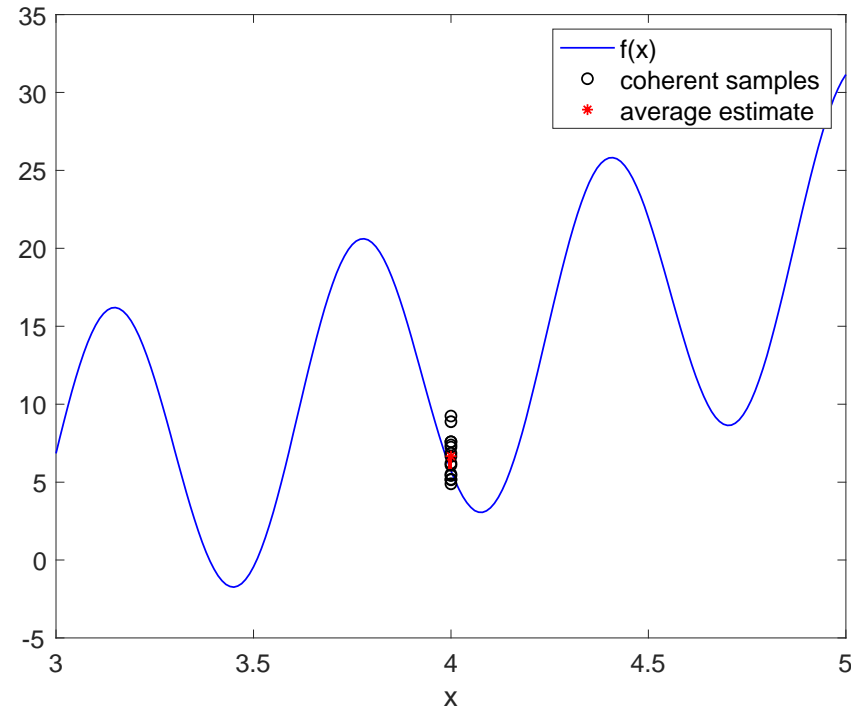
Nearest-Neighbor Regression (cont'd)

Slow varying function. Average estimation on neighborhood samples $y(4 + n) = f(4 + n) + \epsilon$, $\sigma_\epsilon = 2$.



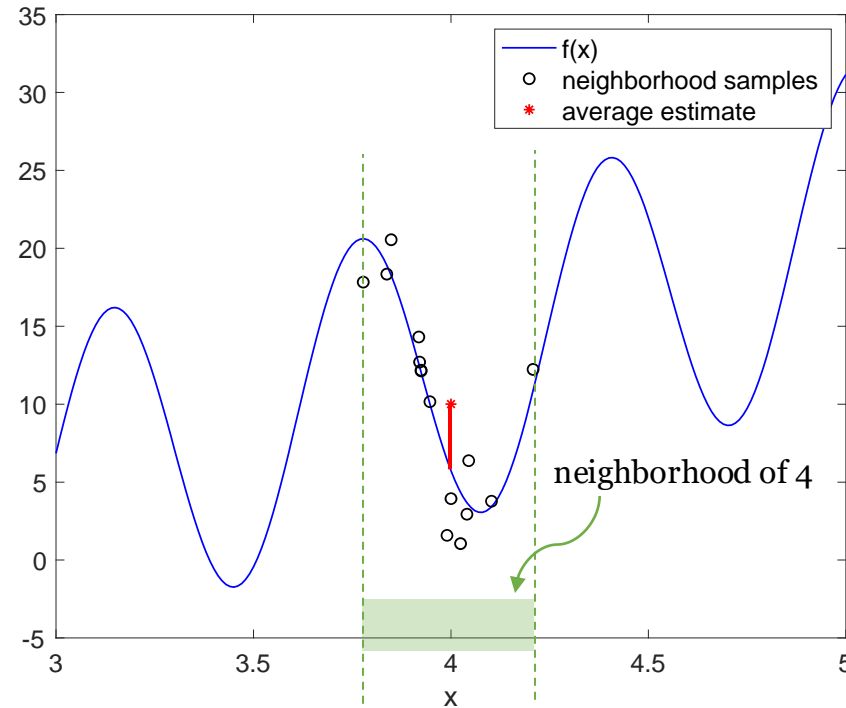
Nearest-Neighbor Regression (cont'd)

Fast varying function. Average estimation on coherent samples $y(4) = f(4) + \epsilon$, $\sigma_\epsilon = 2$.



Nearest-Neighbor Regression (cont'd)

Fast varying function. Average estimation on neighborhood samples $y(4 + n) = f(4 + n) + \epsilon$, $\sigma_\epsilon = 2$.



Nearest-Neighbor Regression (cont'd)

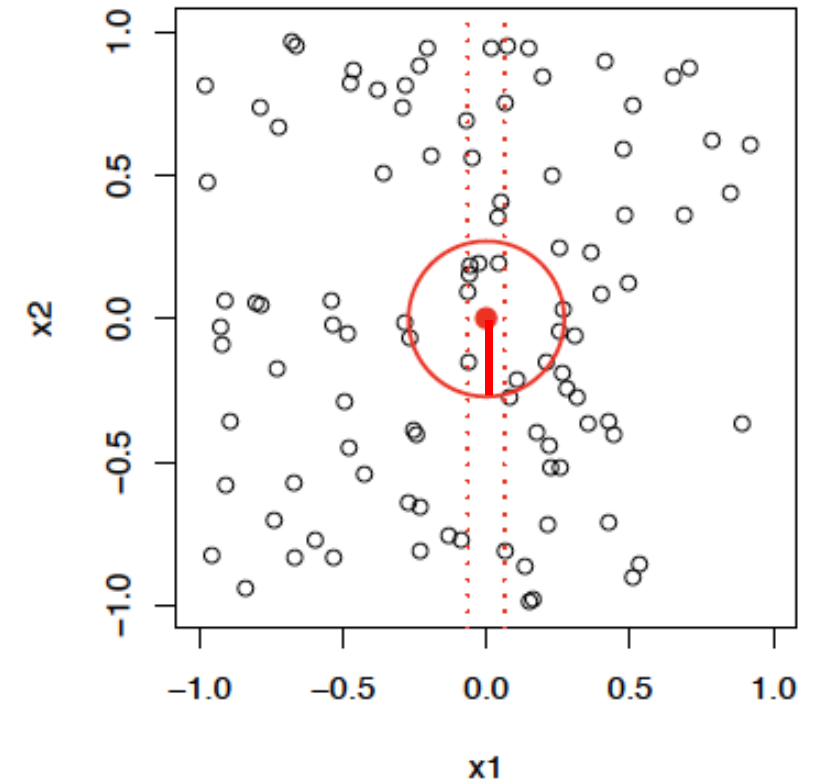
- For each $\mathbf{x} \in I$ let $V(\mathbf{x})$ be the set of points in S that are the nearest to \mathbf{x} .
 - nearest in terms of Euclidean distance, or any other “distance” metric
- We call these points the “nearest neighbors” of \mathbf{x}

- Option 1: Limit $K = |V(\mathbf{x})|$

For example, $K = 10\% |S| = 0.1 N$

- Option 2: $\max_{\mathbf{z} \in V(\mathbf{x})} \text{distance}(\mathbf{z}, \mathbf{x}) \leq C$

For example, $\|\mathbf{z} - \mathbf{x}\|_2 \leq C$



Nearest-Neighbor Regression (cont'd)

For every $\mathbf{x} \in I$, estimate $f(\mathbf{x}) = E[y(\mathbf{x})]$ by the nearest neighbor average:

- Identify input neighborhood $V(\mathbf{x}) = \{\mathbf{z}_k\}_{k=1}^K$ and corresponding outputs $\{y_k\}_{k=1}^K$.
- Estimate $\hat{f}(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^K y_k$.

The Curse of Dimensionality

- Nearest-neighbor averaging performance drops as d increases.
- **This is part of “the curse of dimensionality”**
- Nearest neighbors tend to be far away in high dimensions.
- We need to get a reasonable fraction of the N values of y_i to average to bring the variance down – e.g., 10%.
- A 10% neighborhood in high dimensions need no longer be local, so we lose the spirit of estimating $f(\mathbf{x})$ by local averaging.

The Curse of Dimensionality (cont'd)

