



Lecture 3: Linear Regression (Part 2)

CIS 4190/5190

Spring 2023

Administrivia

- Class currently full at 225 students (211 enrolled, 16 permits out) with about 180 students in the waitlist, many from 3 months ago. Limited movement expected. Unused permits will be revoked as announced earlier, unless you have informed us.
 - If you're on the waitlist, submitting HW1 early will increase priority, but no guarantees.
- HW1 due **Wednesday 8 p.m.**, and HW2 will be posted that evening, on linear regression.
- In most cases, you should use EdSTEM to contact the course team, where you are much more likely to receive a fast response. But if your message must be kept private from TAs:
 - **Always** email both instructors together.
 - Start subject line with “[**CIS 4190 / 5190 Spring 2023**]”.
- Canvas link (for recordings) is on the class webpage > files:
<https://canvas.upenn.edu/courses/1704503> (video recordings posted on day of lecture/next day)
- If you're on the waitlist and submitting HW1, submit directly to gradescope.
- If switching from CIS 4190 to CIS 5190, contact cis-undergrad-advising@seas.upenn.edu
 - More work, and possibly higher grade cutoffs for CIS 5190.
- TA office hours start today. No help for HW1, but you can ask for help with the course material.
- Recitations on Thursday at 5 p.m. on Python, Numpy, Pandas, Scikit-Learn. See EdSTEM post.

Recap: Loss Minimization View of ML

- **To design an ML algorithm:**

- Choose model family $F = \{f_{\beta}\}_{\beta}$ (e.g., linear functions)
- Choose loss function $L(\beta; Z)$ (e.g., MSE loss)

- **Resulting algorithm:**

$$\hat{\beta}(Z) = \arg \min_{\beta} L(\beta; Z)$$

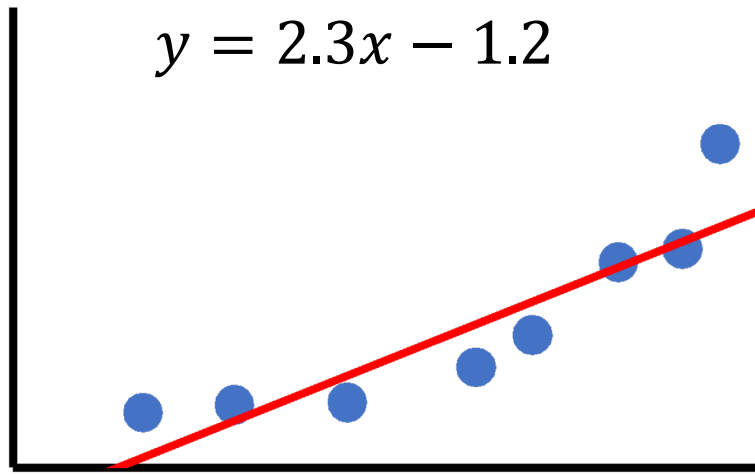
Recap: Linear Regression

- **Input:** Dataset $Z = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- Compute

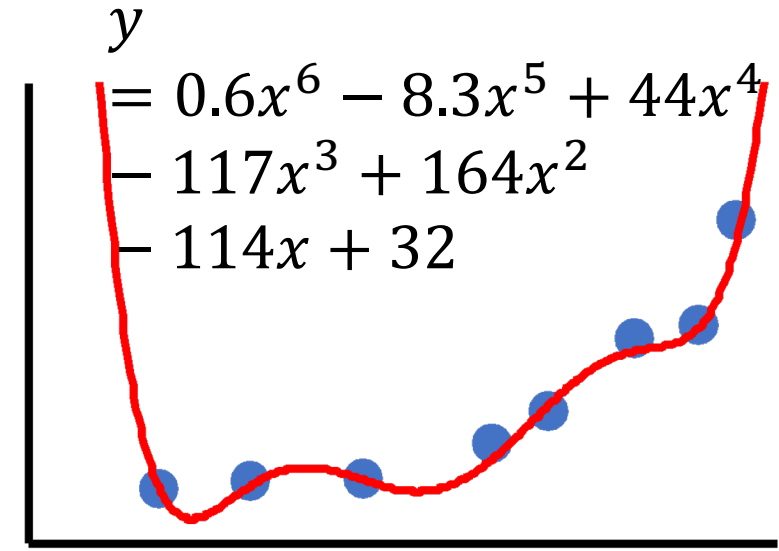
$$\hat{\beta}(Z) = \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \beta^\top x_i)^2$$

- **Output:** $f_{\hat{\beta}(Z)}(x) = \hat{\beta}(Z)^\top x$
- Discuss algorithm for computing the minimal β later (next class)

Recall: Overfitting



Underfitting



Overfitting

Overfitting occurs when:

- The learned hypothesis fits the training set very well e.g. $\mathcal{L}(\boldsymbol{\theta}) \approx 0$, but fails to generalize to new examples

Today's Class

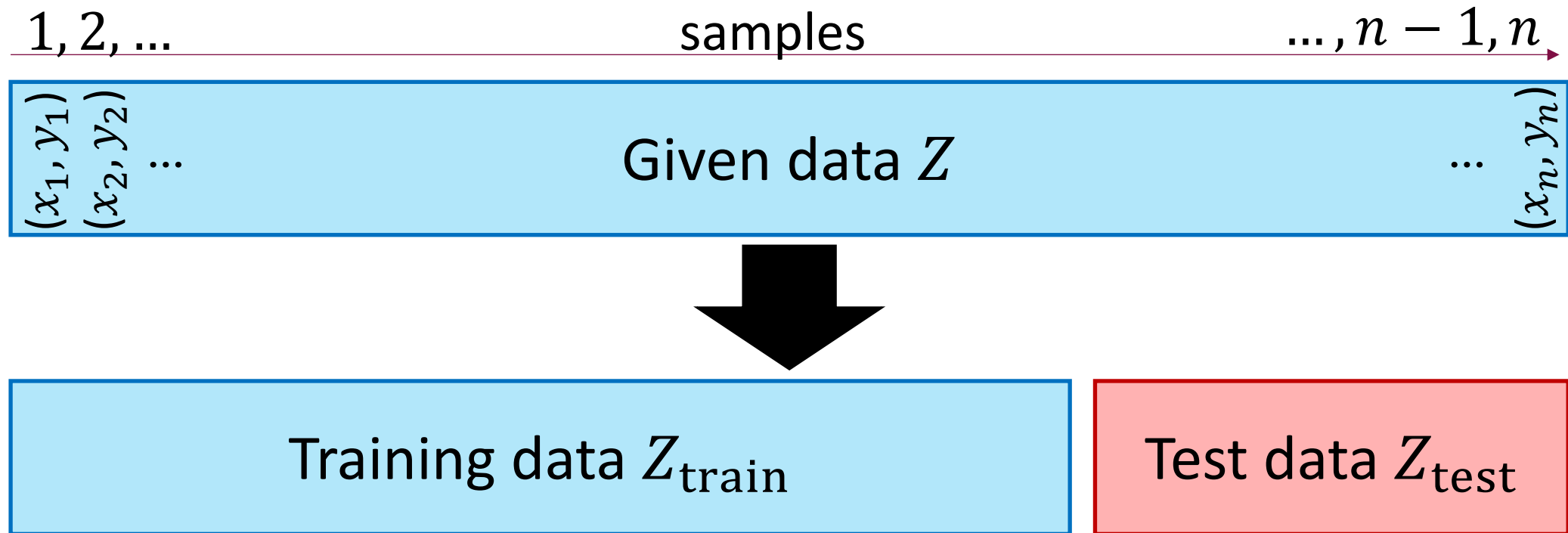
- Understanding, diagnosing, and combating overfitting:
 - Bias and Variance of hypothesis classes
 - Regularized linear regression
 - Cross-Validation
- Feature Selection and Preprocessing
 - Sparse linear regression



Assessing Overfitting

Training/Test Split

- **Issue:** How to detect overfitting vs. underfitting?
- **Solution:** Use **held-out test data** to estimate loss on new data
 - Typically, randomly shuffle data first



Training/Test Split Protocol for Machine Learning

- **Step 1:** Split Z into Z_{train} and Z_{test}

Training data Z_{train}

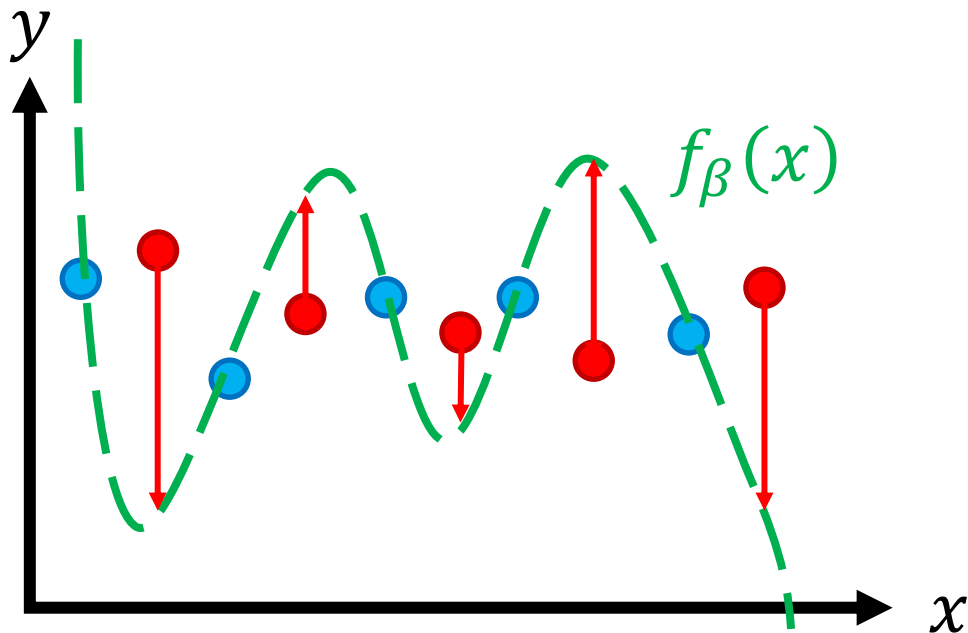
Test data Z_{test}

- **Step 2:** Run linear regression with Z_{train} to obtain $\hat{\beta}(Z_{\text{train}})$
- **Step 3:** Evaluate
 - **Training loss:** $L_{\text{train}} = L(\hat{\beta}(Z_{\text{train}}); Z_{\text{train}})$
 - **Test (or generalization) loss:** $L_{\text{test}} = L(\hat{\beta}(Z_{\text{train}}); Z_{\text{test}})$, (plus other performance metrics besides the loss function)

Training/Test Split Protocol for Machine Learning

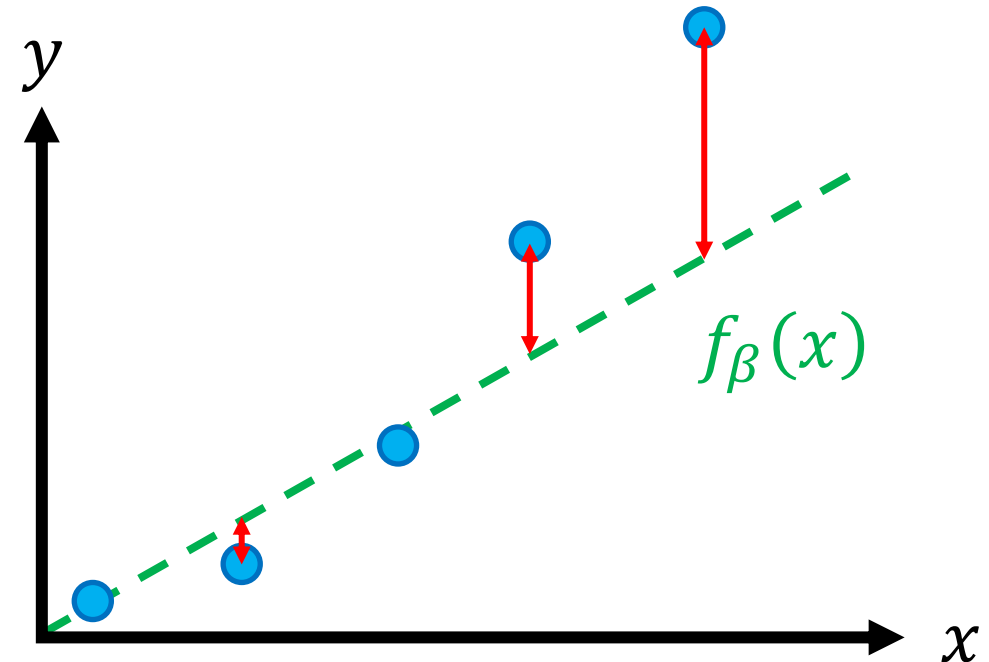
- **Overfitting**

- Fit the **training data** Z well
- Fit new **test data** (x, y) poorly



- **Underfitting**

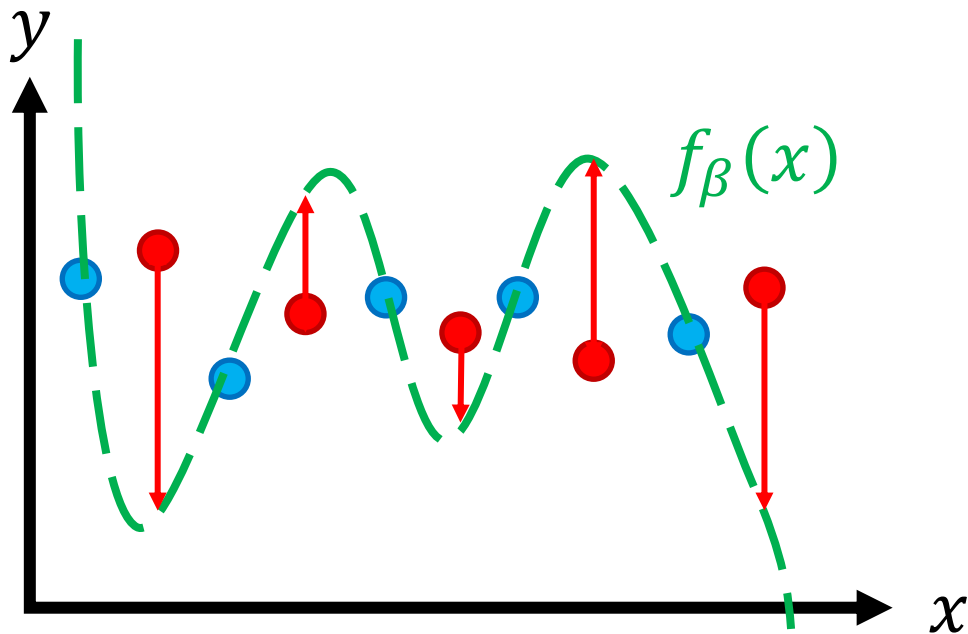
- Fit the **training data** Z poorly
- (Necessarily fit new **test data** (x, y) poorly)



Training/Test Split Protocol for Machine Learning

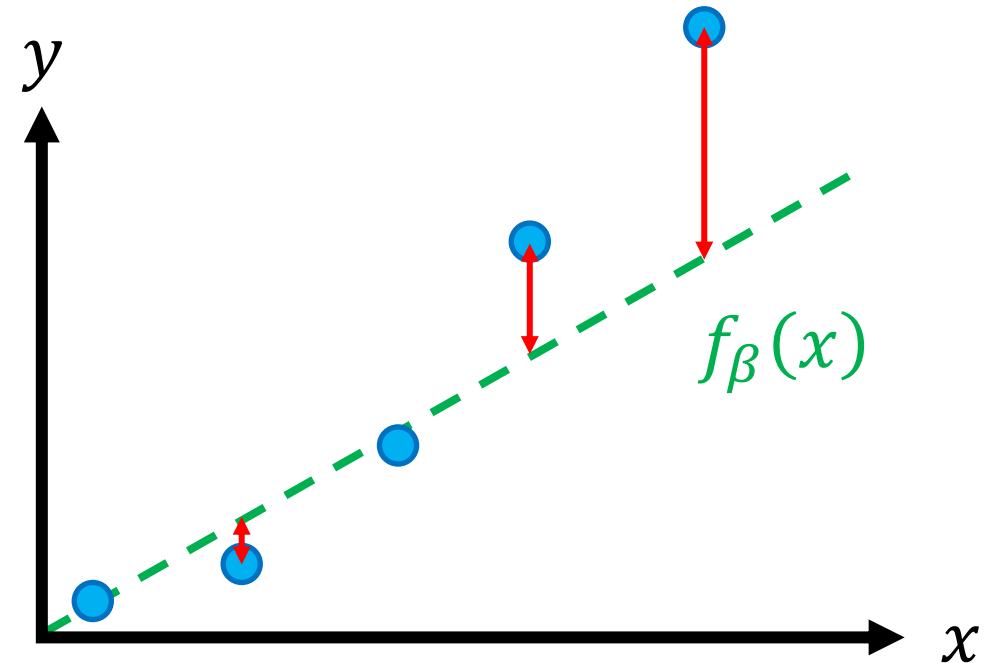
- **Overfitting**

- L_{train} is small
- L_{test} is large



- **Underfitting**

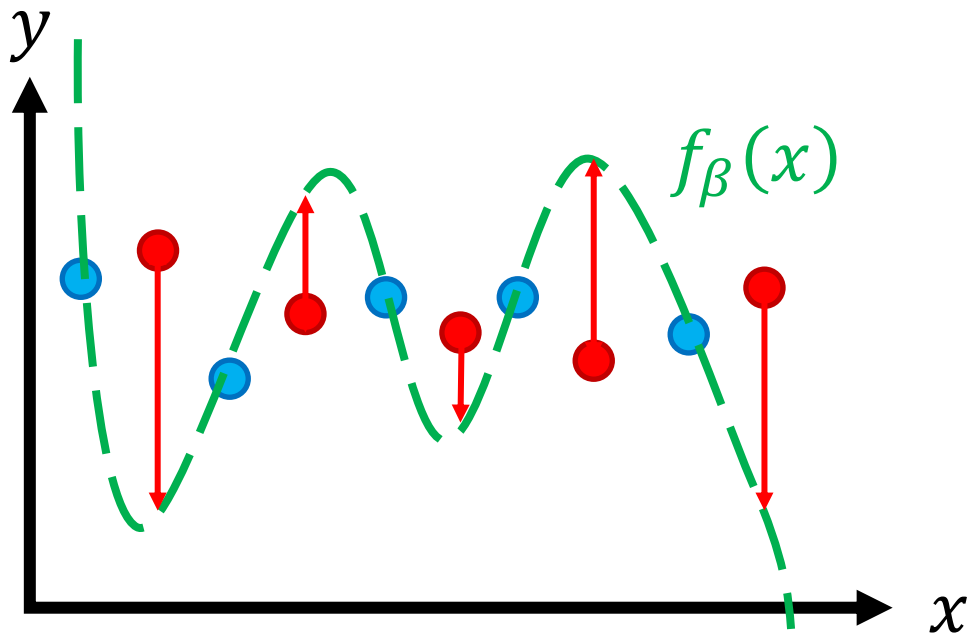
- Fit the **training data** Z poorly
- (Necessarily fit new **test data** (x, y) poorly)



Training/Test Split Protocol for Machine Learning

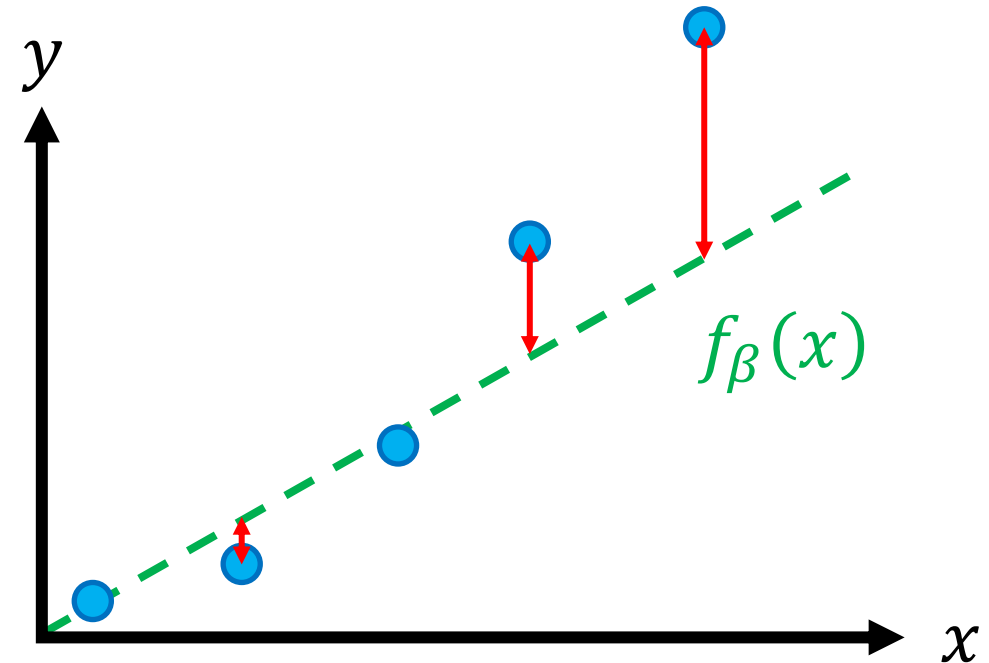
- **Overfitting**

- L_{train} is small
- L_{test} is large



- **Underfitting**

- L_{train} is large
- L_{test} is large



“Independent and Identically Distributed”

- **The “IID” assumption**

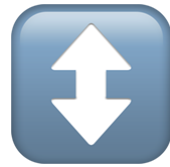
- “Test” data Z_{test} are drawn IID from same data distribution $P(x, y)$ as Z_{train}
- IID = independent and identically distributed
- This is a strong (but common) assumption!

- **Time series data**

- Particularly important failure case since data distribution may shift over time
- **Solution:** Split along time (e.g., data before 9/1/20 vs. data after 9/1/20)



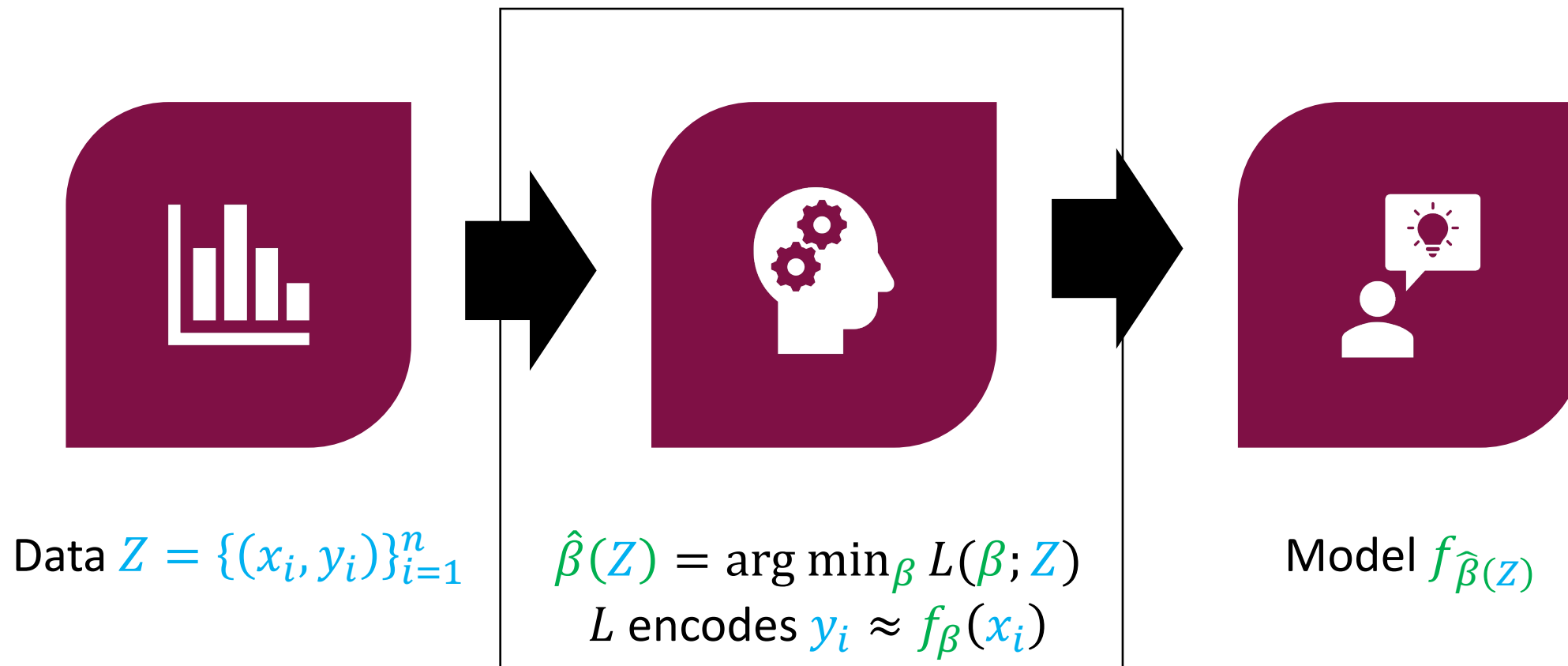
Underfitting & Overfitting



Bias & Variance

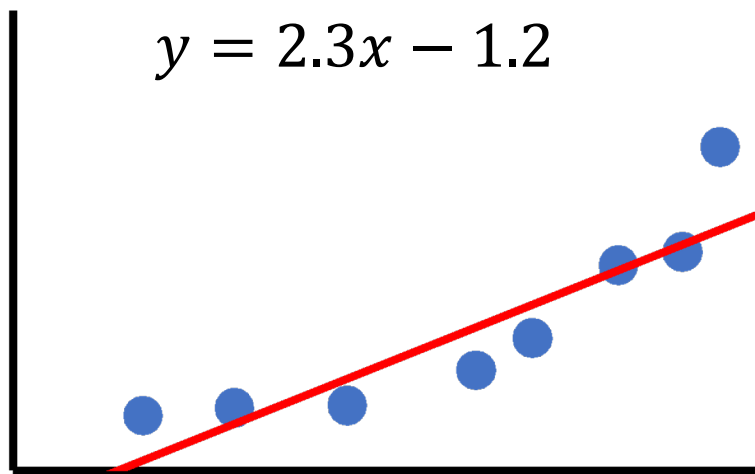
Recall

Reminder: we don't yet know how to find the argmin. For the moment, we will continue to assume it will be found.

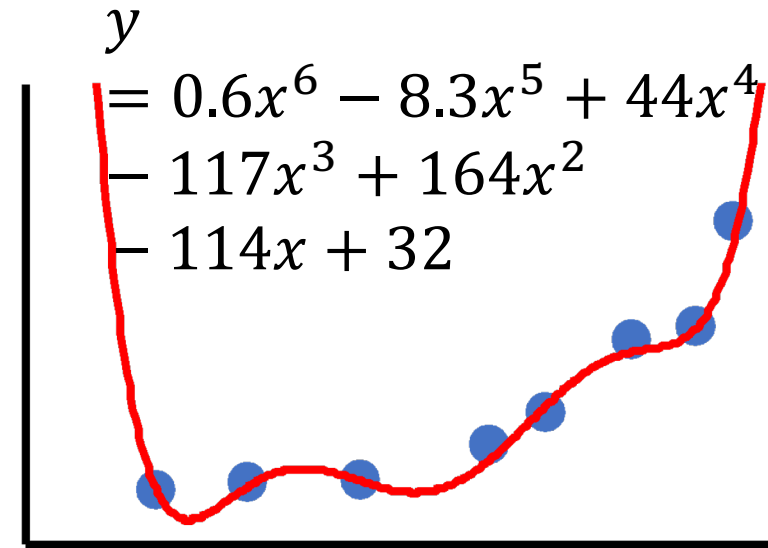


We are still thinking about two main design choices that influence this box:
Model class (linear regression, feature map etc.), and loss function (MSE)

Recall: Underfitting and Overfitting



Underfitting



Overfitting

We will understand these phenomena now through two properties of a model family, “bias”, and “variance”.

Language for thinking about the ways in which model families can be bad.

How to Fix Underfitting/Overfitting?

Three main options:

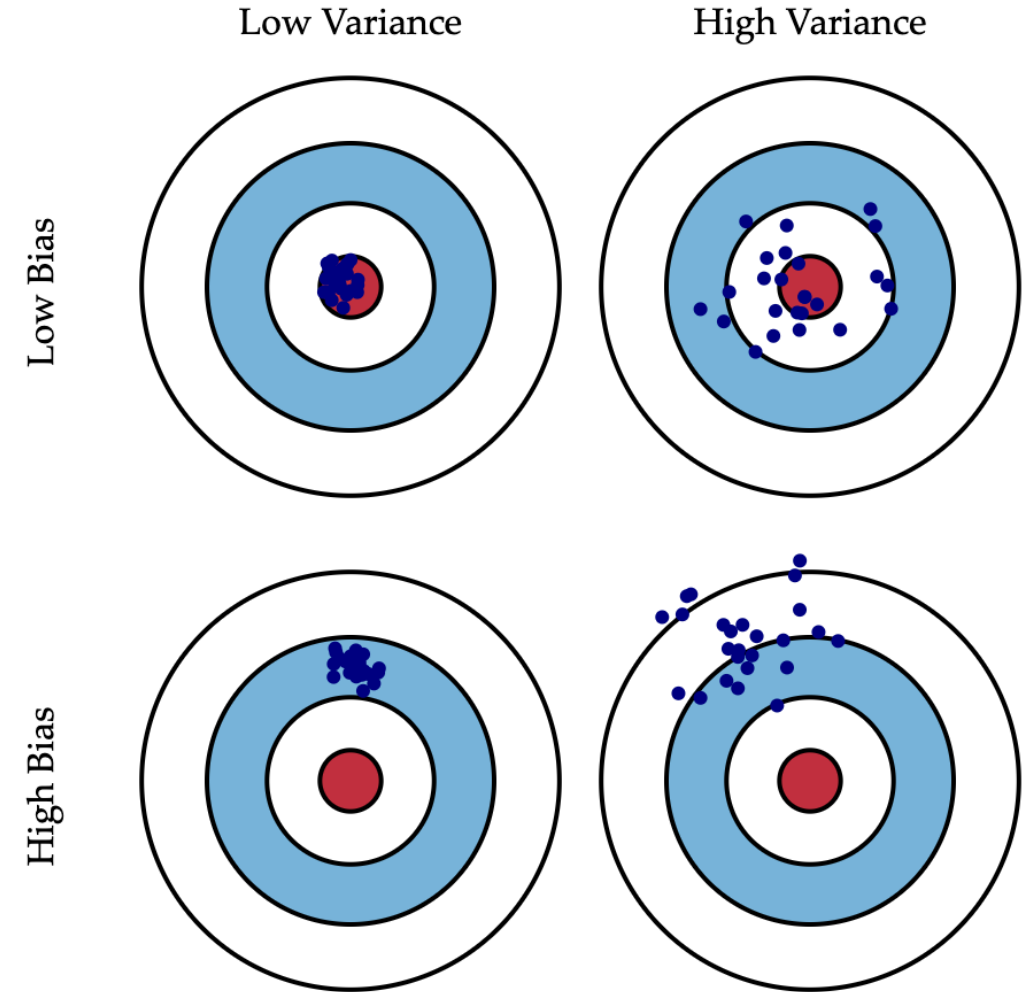
- Improve the training dataset
- Choose the right model family
- Choose the right loss function

We will explore these in some detail over the next few slides.

Definitions: “Bias” and “Variance”

Imagine you draw k training datasets from the same probability distribution over data, and each time fit your model $\{f_{\beta}\}_{1:k}$ to it.

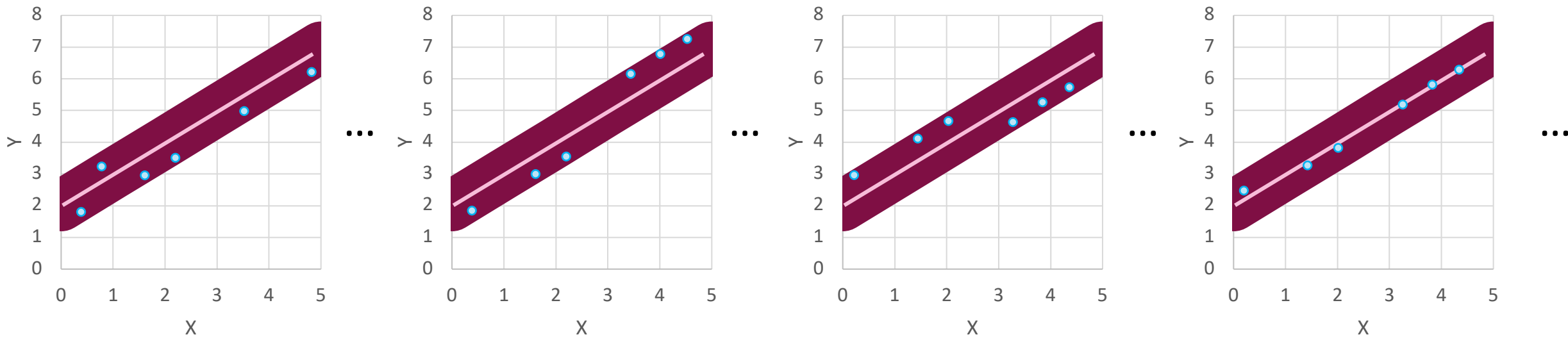
- “**Variance**”: how much do those fitted functions $\{f_{\beta}\}_{1:k}$ differ amongst each other, on average over the data distribution?
- “**Bias**”: how much does the average of all those fitted functions $\{f_{\beta}\}_{1:k}$ deviate from the “true” function over the data distribution?



Drawing Multiple Training Datasets

Consider a linear “true function” $f^*(x) = x + 2$ that generates labels y_i for training data with uniform measurement noise in $[-1, +1]$.

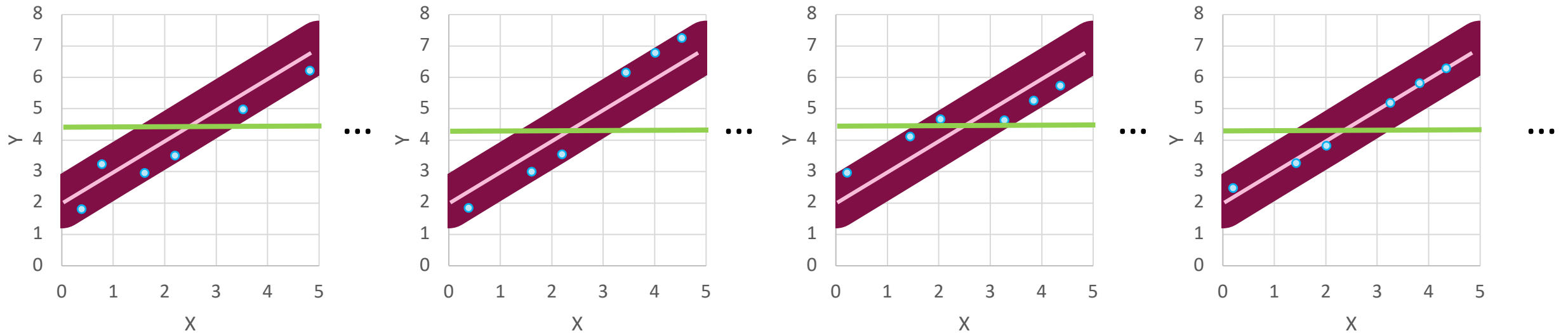
Let us draw $k \rightarrow \infty$ training sets of $n = 6$ samples each, drawn from $P(X, Y)$.



Different *Constant* Fits

What if the hypothesis class was the constant function class

$$f_{\theta}(x) = \theta_0$$

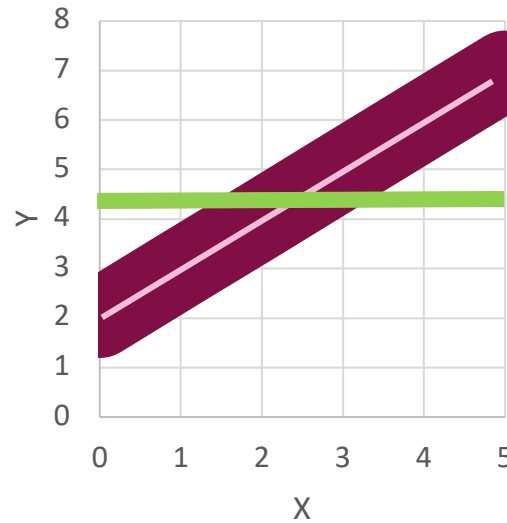


Different *Constant* Fits

What if the hypothesis class was the constant function class

$$f_{\theta}(x) = \theta_0$$

Almost identical fits
“low variance”



Average fit far from the true
function
“high bias”

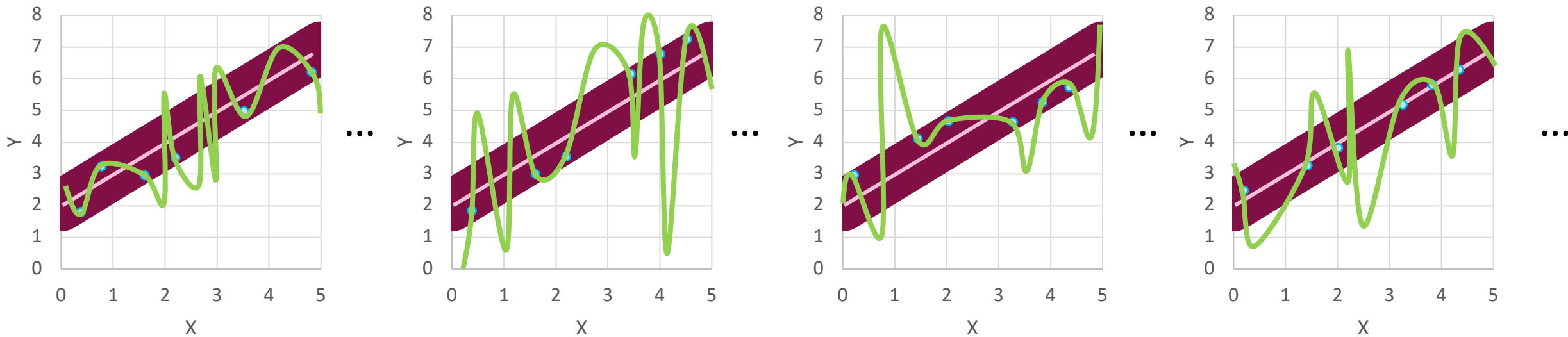


Theoretical result: Generalization MSE \approx “Bias” + “Variance”

Different 10th Degree Curve Fits

What if the hypothesis class was instead a 10th degree monomial

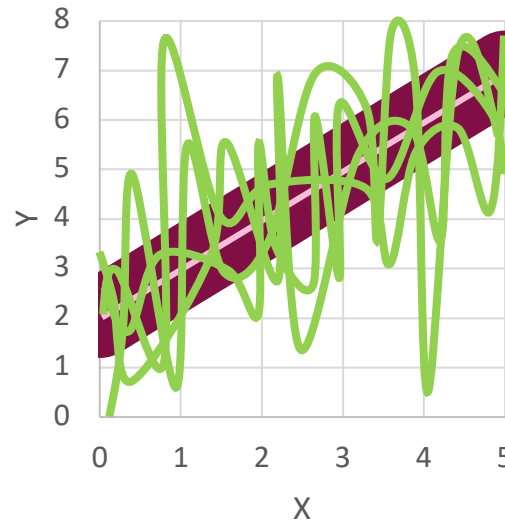
$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \dots \theta_{10} x^{10}$$



Different 10^{th} Degree Fits

What if the hypothesis class was instead a 10^{th} degree monomial
$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 + \dots \theta_{10} x^{10}$$

Very different fits
“high variance”



Average fit close to the true function
“low bias”



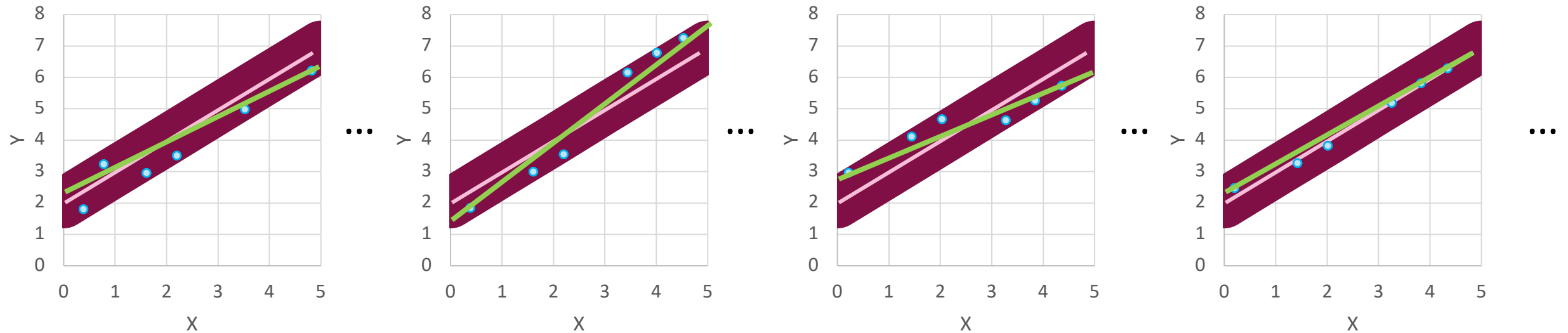
Theoretical result: Generalization MSE \approx “Bias” + “Variance”

Different *Linear* Fits

Say, our hypothesis class is a line:

$$f_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Fit by minimizing MSE with any optimizer. What would the resulting line look like?



Slightly different fits

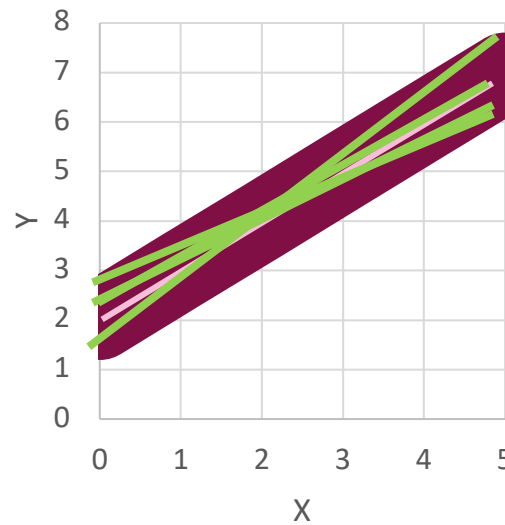
Different *Linear* Fits

Say, our hypothesis class is a line:

$$f_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Fit by minimizing MSE with any optimizer. What would the resulting line look like?

Quite similar fits
“low variance”



Average fit close to the true
function

“low bias”



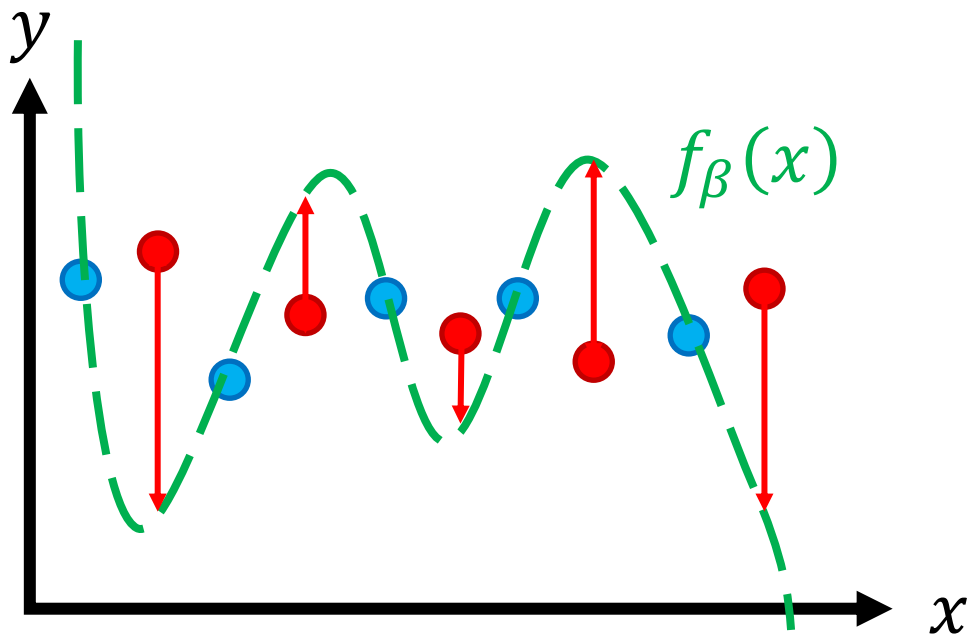
This is the type of model complexity you want to have

Theoretical result: Generalization MSE \approx “Bias” + “Variance”

Bias-Variance Tradeoff

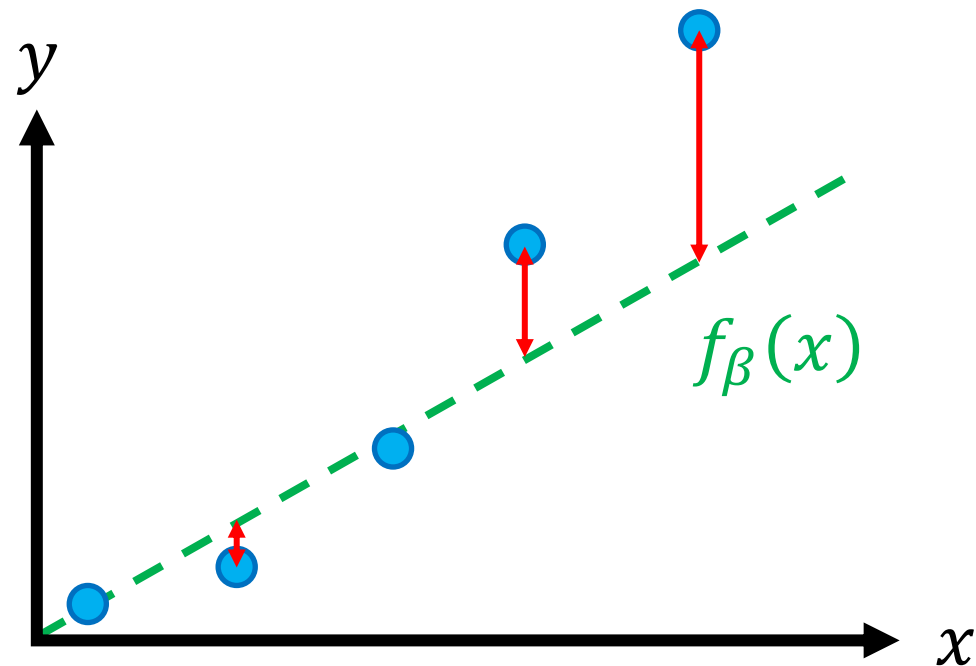
- **Overfitting (high **variance**)**

- High capacity model capable of fitting complex data
- Insufficient data to constrain it

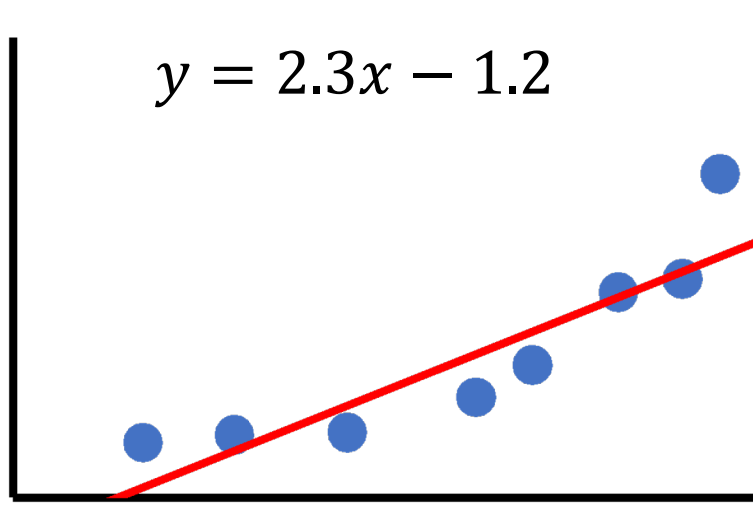


- **Underfitting (high **bias**)**

- Low capacity model that can only fit simple data
- Sufficient data but poor fit

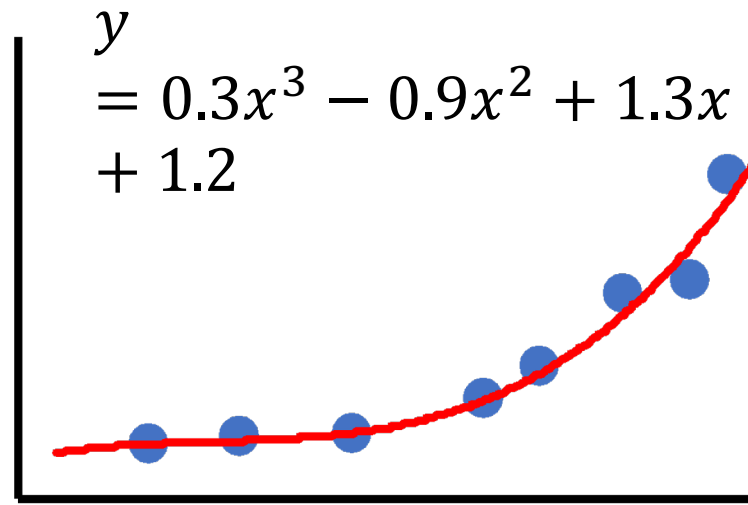


Underfitting & Overfitting, Bias & Variance



Underfitting

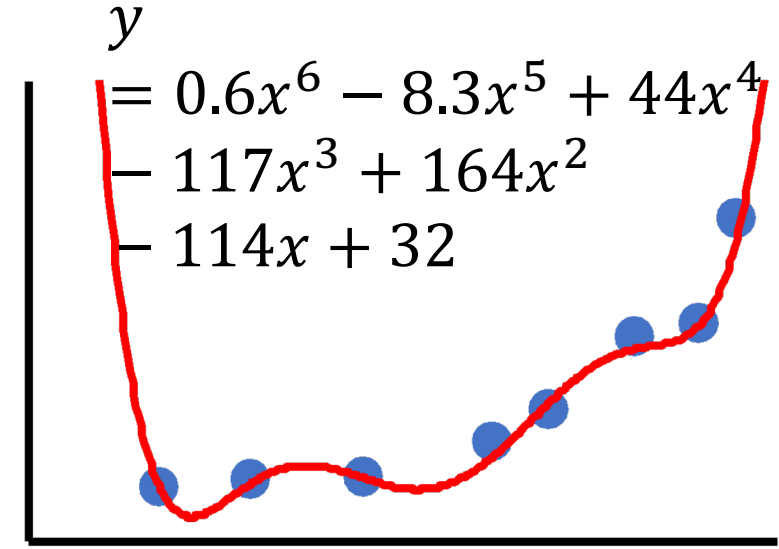
This hypothesis space has high bias



Correct fit

The data decides which hypothesis class choices are good/bad.

(Train/Test Split Protocol)

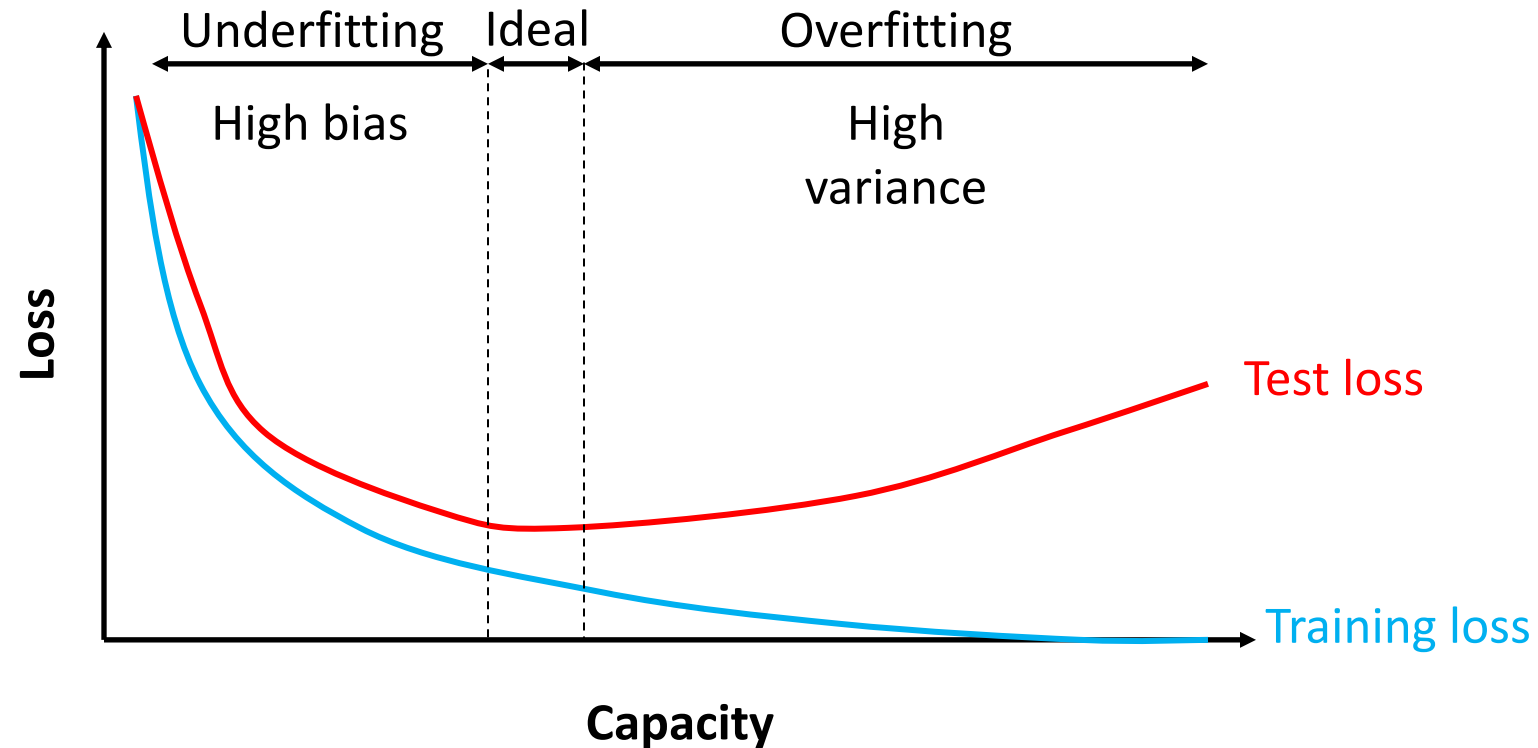


Overfitting

This hypothesis space has high variance

Under/Over -Fitting & Model Capacity

Expanding the hypothesis class usually leads to higher variance, lower bias.
(e.g. when adding new dimensions to the feature map)



Bias-Variance Tradeoff For Linear Regression

- For linear regression with feature maps, increasing feature dimension d' ...
 - Tends to **increase capacity**
 - Tends to **decrease bias** but **increase variance**
- Need to construct ϕ to balance tradeoff between bias and variance
 - **Rule of thumb:** You will need $n \approx d' \log d'$ samples, if your ϕ has dimension d'
- A large fraction of data science work is data cleaning + feature engineering. We will see some common rules of thumb for feature engineering soon.

The Effect of Dataset Size

Recall, we said:

Let us draw $k \rightarrow \infty$ training sets of $n = 6$ samples each, drawn from $P(X, Y)$.

Q: What if we had drawn larger training sets? Would it impact:

- Bias?

Usually no. As $k \rightarrow \infty$, the average of fits $\{f_\beta\}_{1:k}$ to k training sets of finite size n (for any n) approaches the fit to an $n' \rightarrow \infty$ -sized set.

Often convenient to think of bias as the lowest achievable error corresponding to the “best model” within the hypothesis space.

- Variance?

Yes, as the dataset size grows, different i.i.d. datasets would induce similar fits \Rightarrow lower variance.

Convenient to think of variance as average error w.r.t. the “best model”.

The Effect of Dataset Size

As dataset size grows:

- Generalization error (\approx ``Bias'' + ``Variance'') is dominated by bias.
- To reduce error, we select high capacity, low bias models.

Larger datasets have room for expanded hypothesis classes.



Regularization

How to Fix Underfitting/Overfitting?

Recall, three main options:

- Improve the training dataset (collect more data)
- Choose the right model family (not too complex, not too simple)
- Choose the right loss function

We will explore this third option now.

Regularization: Modifying the Loss function

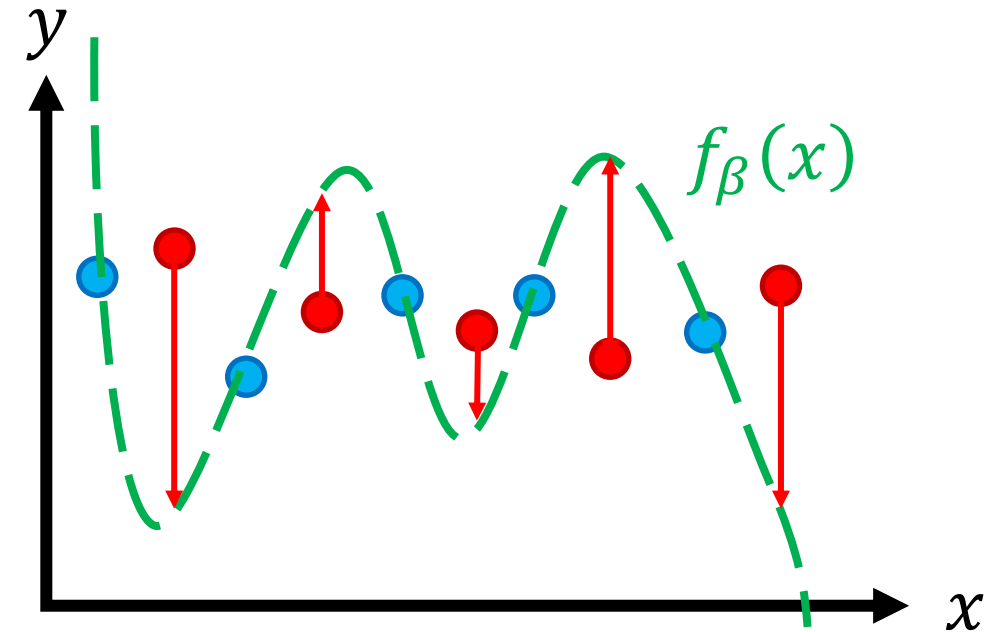
- **Intuition:** We *only* asked the ML algorithm to fit the training data as well as possible, so it produced overly complex fits → “Overfitting”

$$L(\beta; Z) = \text{Train MSE}$$

- **Solution:** we will ask the model to produce a “*simple fit*” to the training data.

$$L(\beta; Z) = \text{Train MSE} + \text{Fit complexity}$$

How to measure this?



Recall: Mean Squared Error Loss


- Mean squared error loss for linear regression:

$$L(\beta; Z) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta^\top x_i)^2$$

Linear Regression with L_2 Regularization

- **Original loss** + **regularization**:

One measure of fit complexity

$$\begin{aligned} L(\beta; Z) &= \frac{1}{n} \sum_{i=1}^n (y_i - \beta^\top x_i)^2 + \lambda \cdot \|\beta\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \beta^\top x_i)^2 + \lambda \sum_{j=1}^d \beta_j^2 \end{aligned}$$


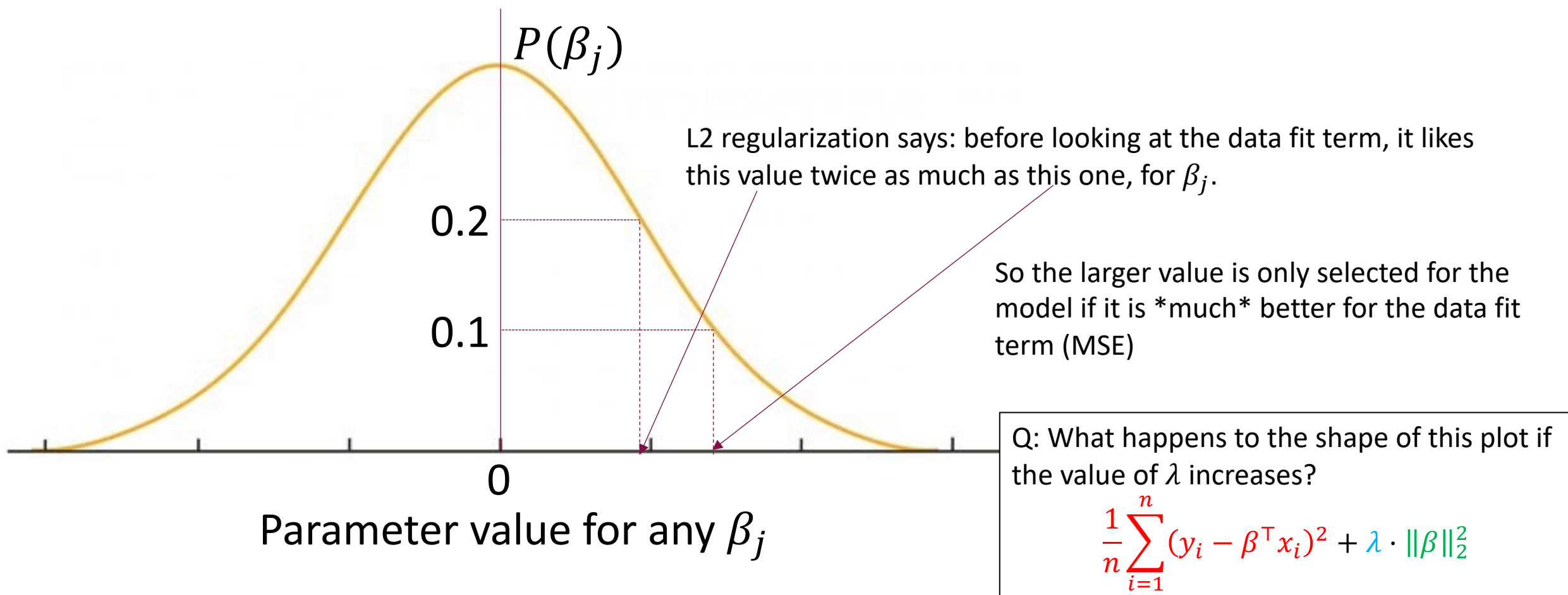
- λ is a **hyperparameter** that must be tuned (satisfies $\lambda \geq 0$)

Intuition on L_2 Regularization

- A thought experiment.
 - Consider a feature map with $d = 10$ features $[x_1, \dots, x_d]$
 - Suppose during linear regression that we forced $\beta_j = 0$ for all $j > 5$
 - This is exactly equal to choosing a smaller-capacity hypothesis class, induced by the smaller feature map with $d = 5$
- Thus, forcing β_j 's to be 0 induces a smaller-capacity hypothesis class.
- The soft version of this: encouraging β_j 's to have small magnitude also induces a smaller-capacity hypothesis class.
 - This is what L_2 regularization does: $\sum_{j=1}^d \beta_j^2 = \|\beta\|_2^2 = \|\beta - 0\|_2^2$
 - Pulls coefficients towards 0
 - As $\lambda \rightarrow \infty$, it forces $\beta = 0$

Intuition on L_2 Regularization: Gaussian Priors

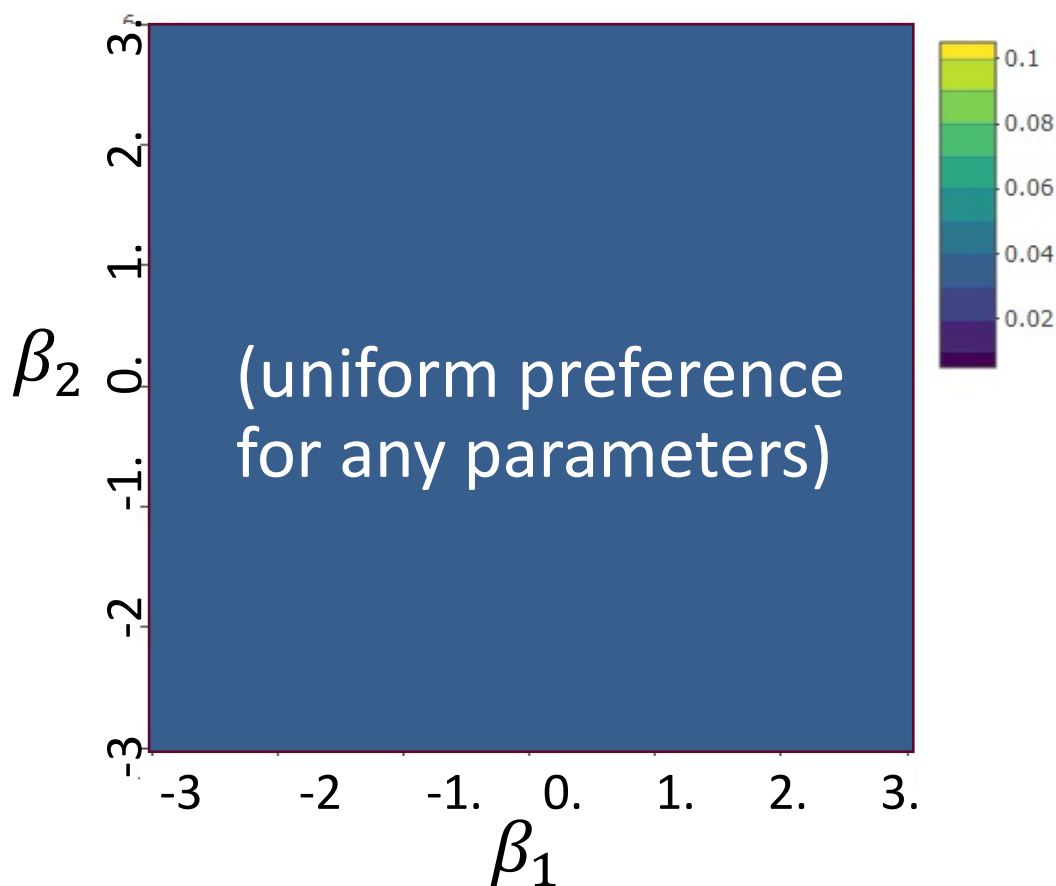
L_2 regularized linear regression amounts to preferring smaller weights according to a Gaussian pdf.



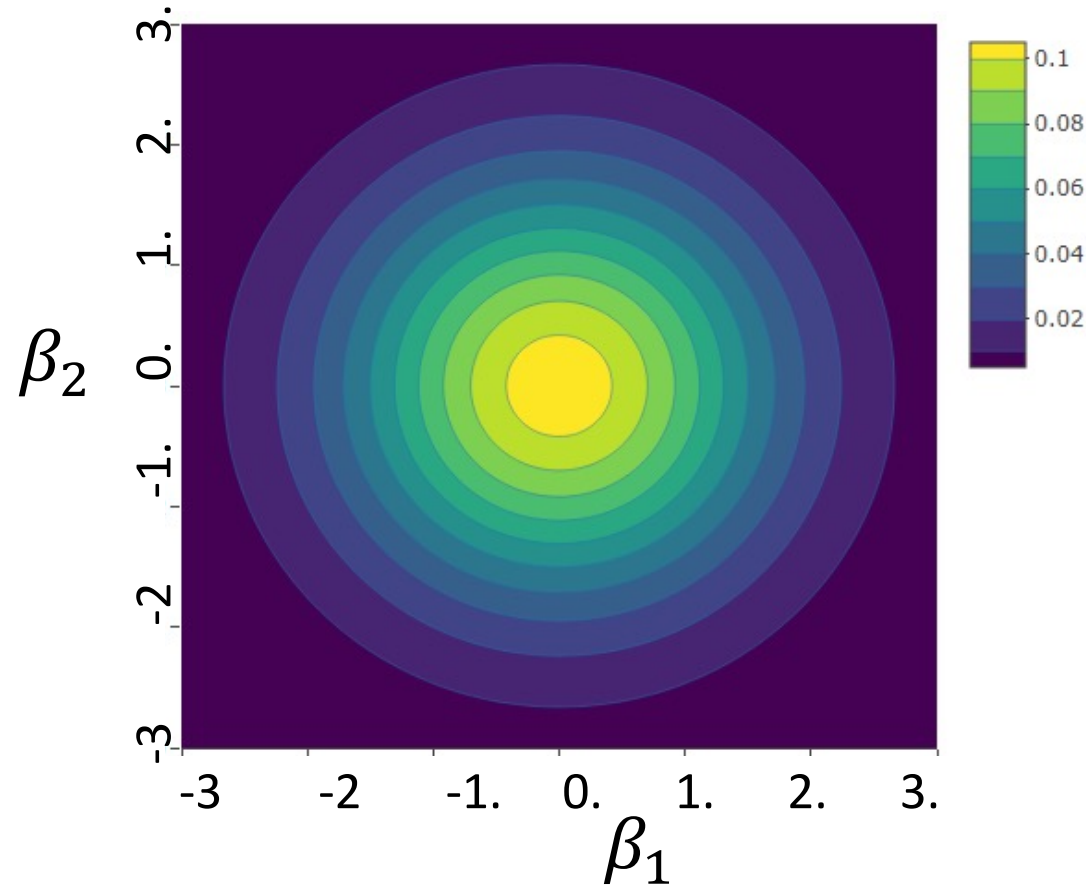
Intuition on L_2 Regularization: Gaussian Priors

- Extending the gaussian pdf over each individual weight β_j to the full parameter vector β , the hypothesis space is constrained.

Before regularization



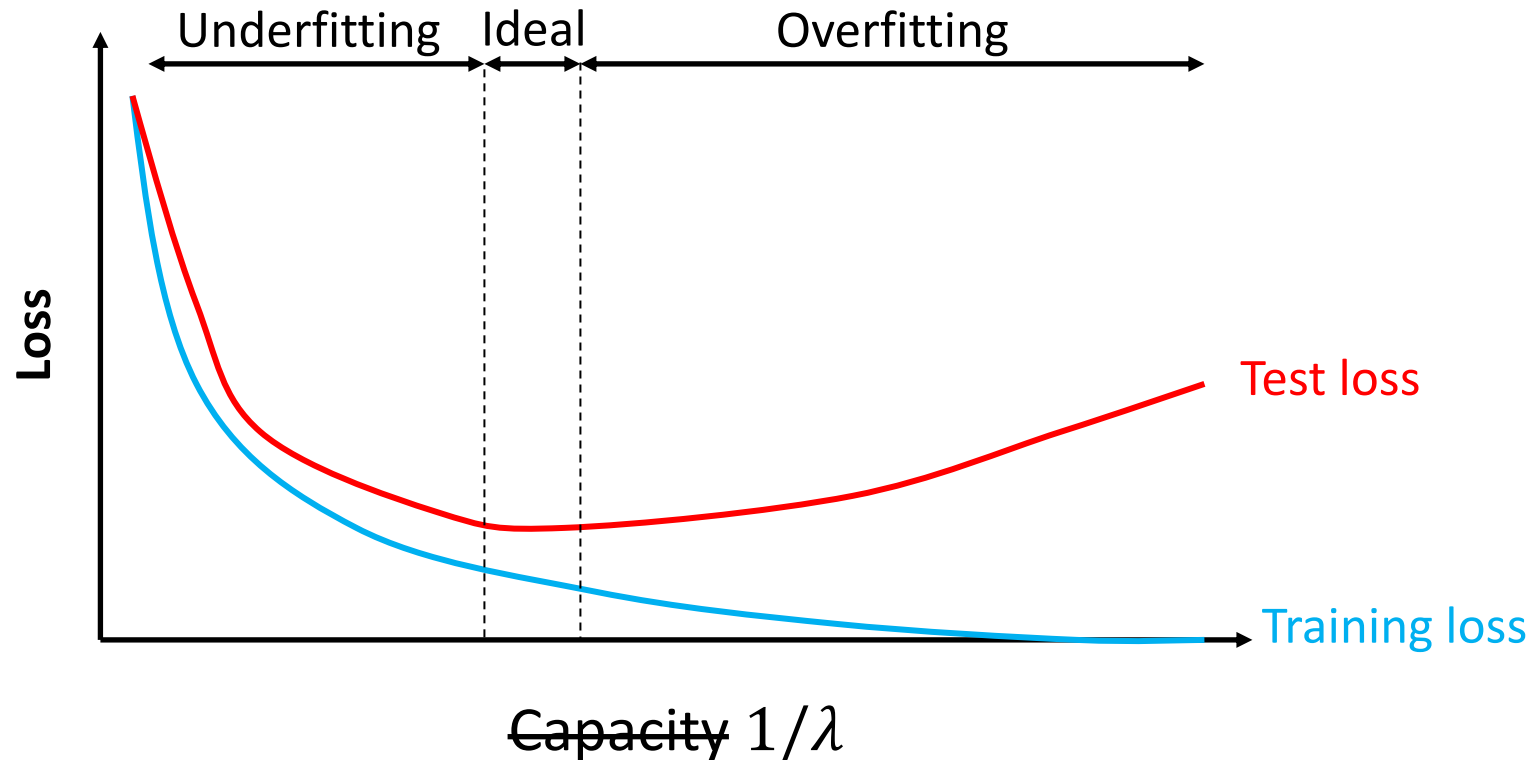
With L2 regularization



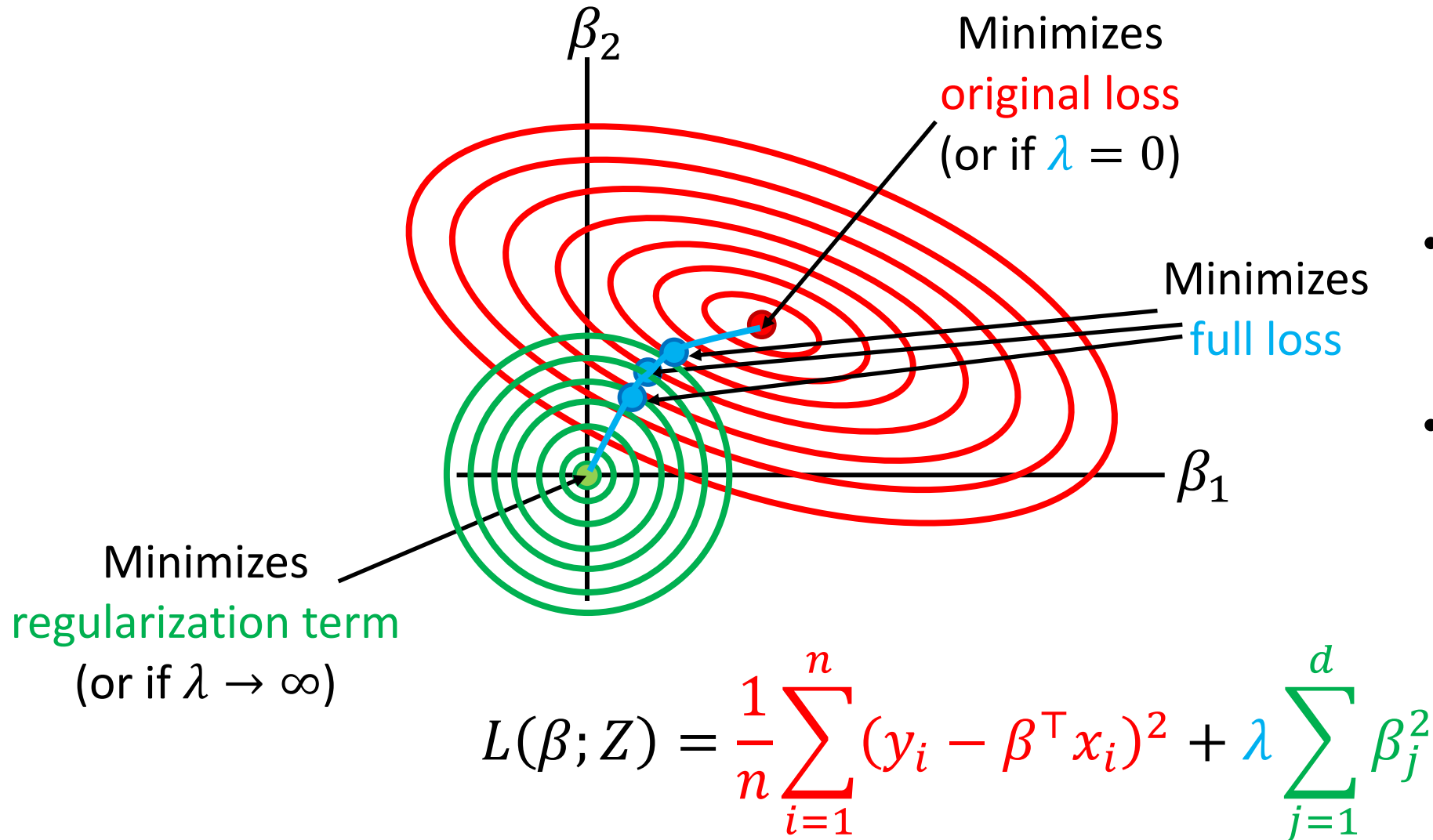
Intuition on L_2 Regularization

- **Why does it help?**
 - Encourages “simple” functions
 - E.g., Use λ to tune bias-variance tradeoff.

Q: How would you set λ to get higher bias / lower variance?



Intuition on L_2 Regularization



- At this point, the gradients are **equal** (with opposite sign)
- Tradeoff depends on choice of λ

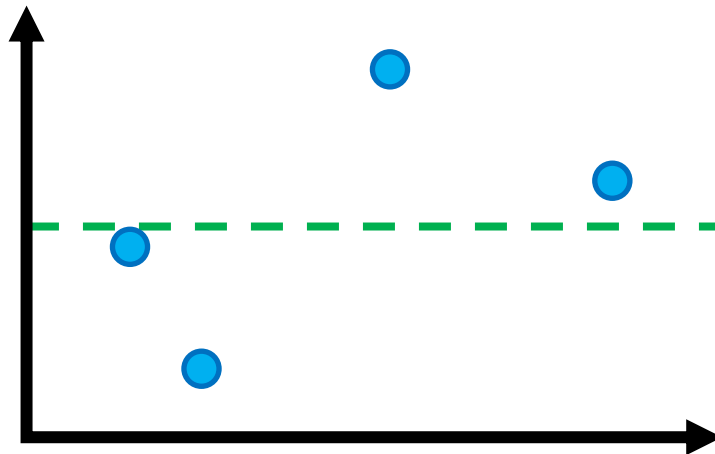
Regularization and Intercept Term / “Bias” Parameter

- **Common convention:** if using intercept term ($\phi(x) = [\textcolor{red}{1} \quad x_1 \quad \dots \quad x_d]^\top$), no penalty on the “bias parameter” β_1 :

$$L(\beta; Z) = \frac{1}{n} \sum_{i=1}^n (\textcolor{red}{y}_i - \beta^\top \textcolor{red}{x}_i)^2 + \textcolor{blue}{\lambda} \sum_{j=2}^d \textcolor{green}{\beta}_j^2$$

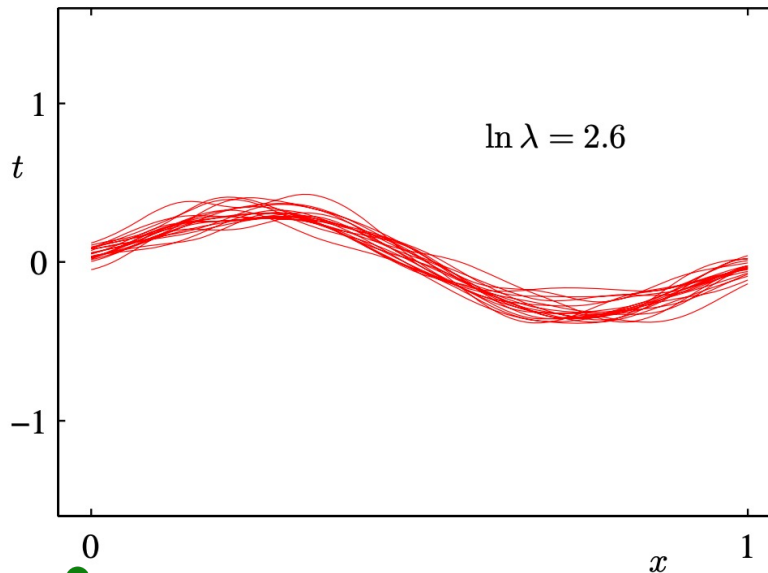
← Sum from $j = 2$

- As $\textcolor{blue}{\lambda} \rightarrow \infty$, we have $\textcolor{green}{\beta}_2 = \dots = \textcolor{green}{\beta}_d = 0$
 - I.e., only fit $\textcolor{green}{\beta}_1$ (which yields $\hat{\textcolor{green}{\beta}}_1(\textcolor{blue}{Z}) = \text{mean}(\{\textcolor{blue}{y}_i\}_{i=1}^n)$)



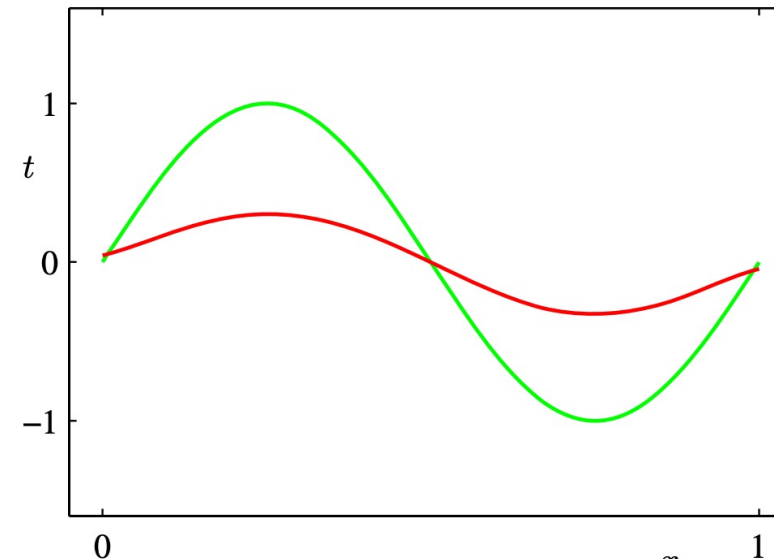
Effect of λ on Bias and Variance

Fits to various datasets sampled from an underlying sinusoid



Illustrates variance

Mean fit compared to the “true function”



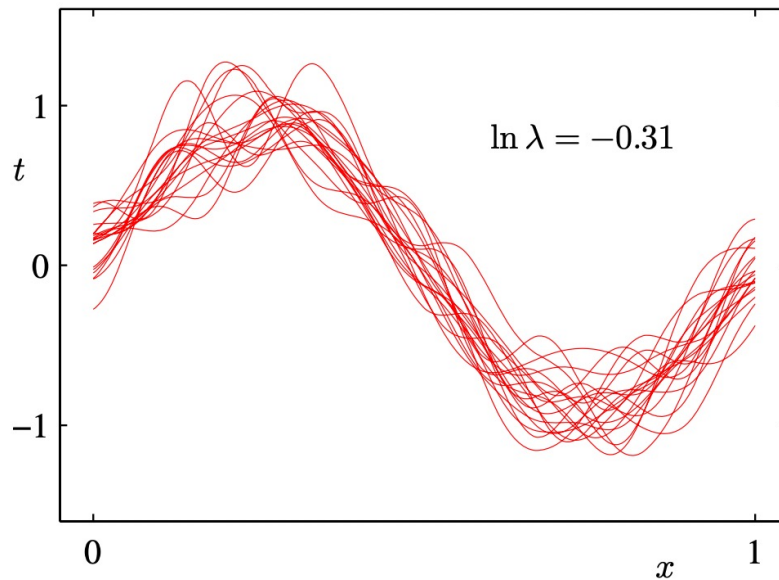
Illustrates bias

With λ too high, the algorithm finds overly simple solutions.

*dataset size n set to 25, and feature map dimensionality d' is 25 (gaussian feature map)

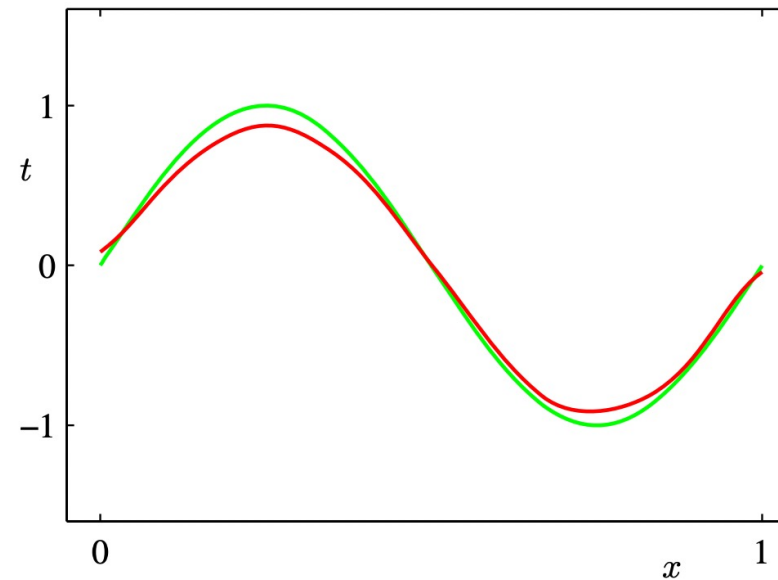
Effect of λ on Bias and Variance

Fits to various datasets sampled from an underlying sinusoid



Illustrates variance

Mean fit compared to the “true function”



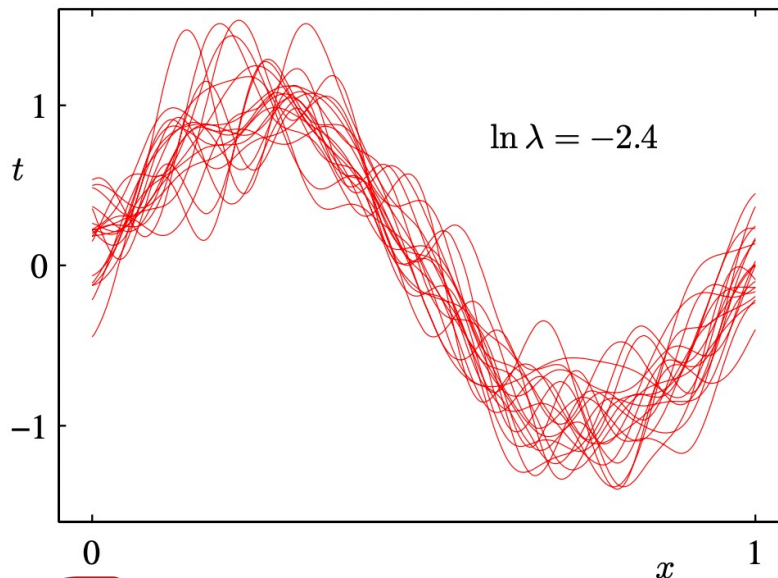
Illustrates bias



*dataset size n set to 25, and feature map dimensionality d' is 25 (gaussian feature map)

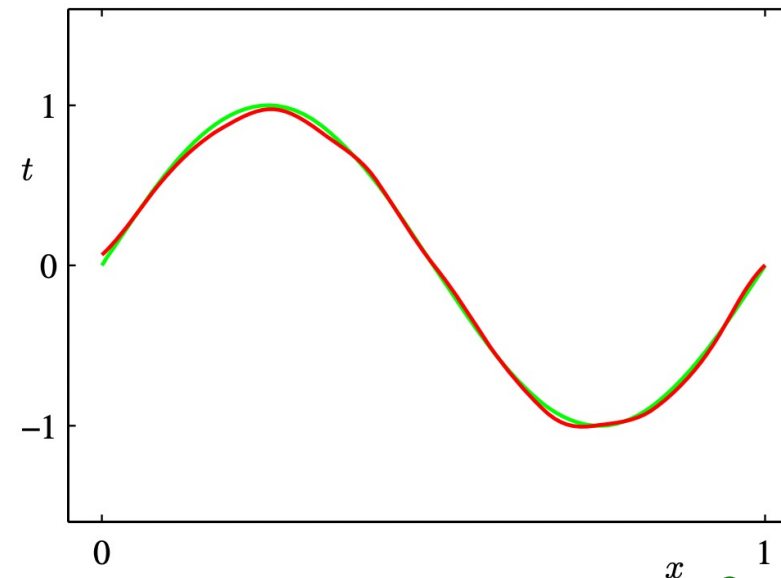
Effect of λ on Bias and Variance

Fits to various datasets sampled from an underlying sinusoid



Illustrates variance

Mean fit compared to the “true function”



Illustrates bias

With λ too low, the algorithm finds overly complex solutions.

*dataset size n set to 25, and feature map dimensionality d' is 25 (gaussian feature map)

General Regularization Strategy

- **Original loss** + **regularization**:

$$L_{\text{new}}(\beta; Z) = L(\beta; Z) + \lambda \cdot R(\beta)$$

- Offers a way to express a preference for “simpler” functions in family
- Typically, regularization is independent of data

Q: For the new parameters $\beta_{\text{new}}^* = \min_{\beta} L_{\text{new}}$, would their corresponding value of $L(\beta; Z)$ be smaller or larger than before regularization?