

Welcome !!



BITS Pilani
Pilani Campus

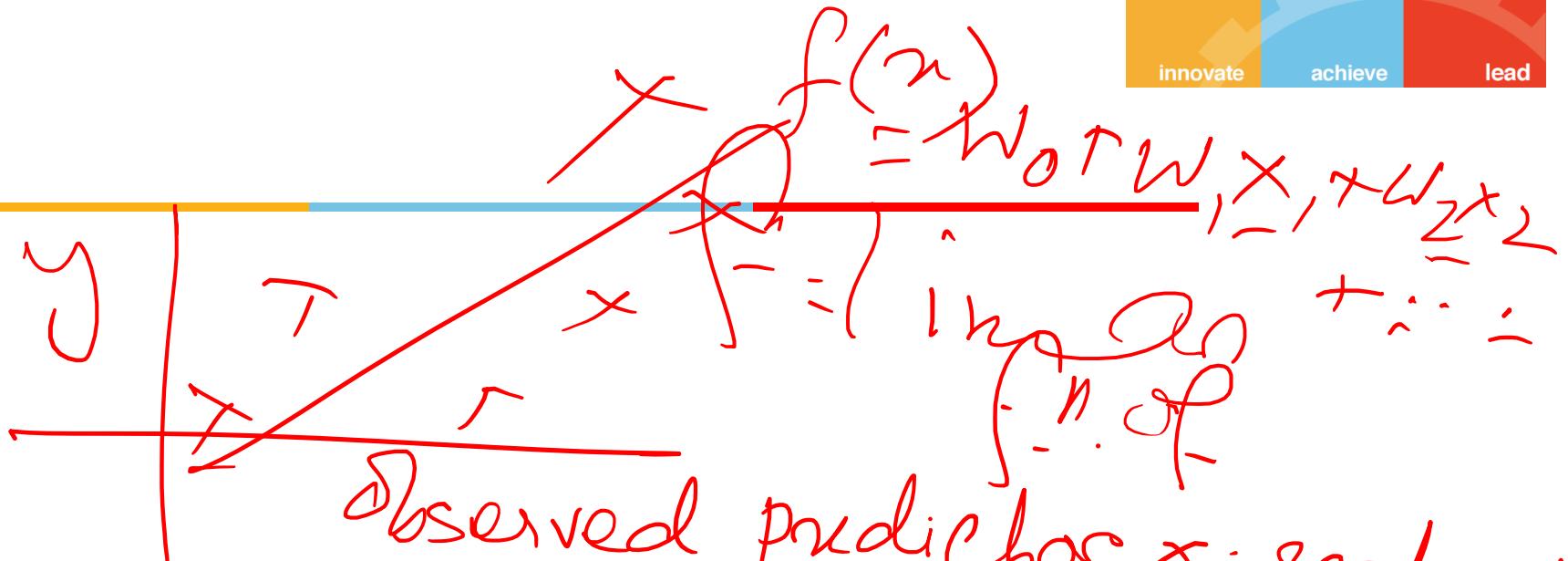
Machine Learning

AIML CLZG565

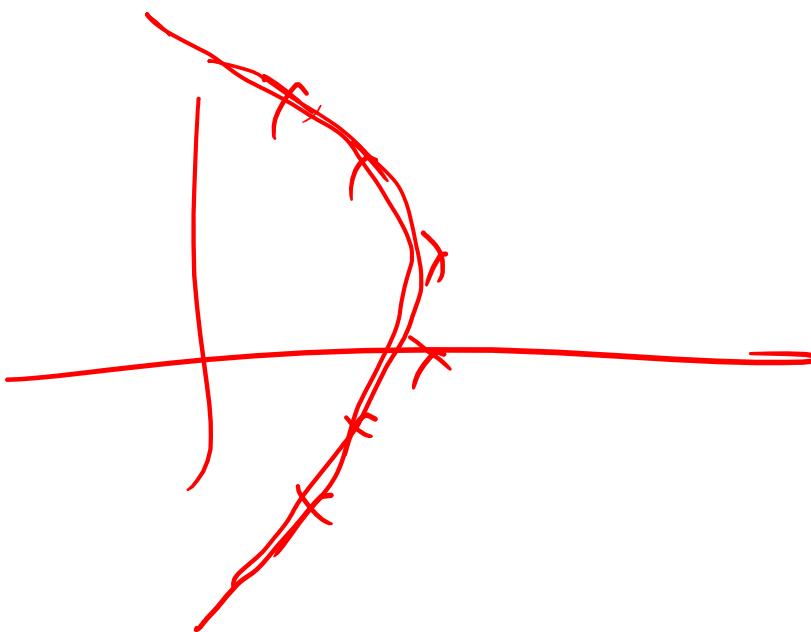
M3 : Linear Models for Regression

Dr. Sugata Ghosal
sugata.ghosal@pilani.bits-pilani.ac.in

Linear Regression



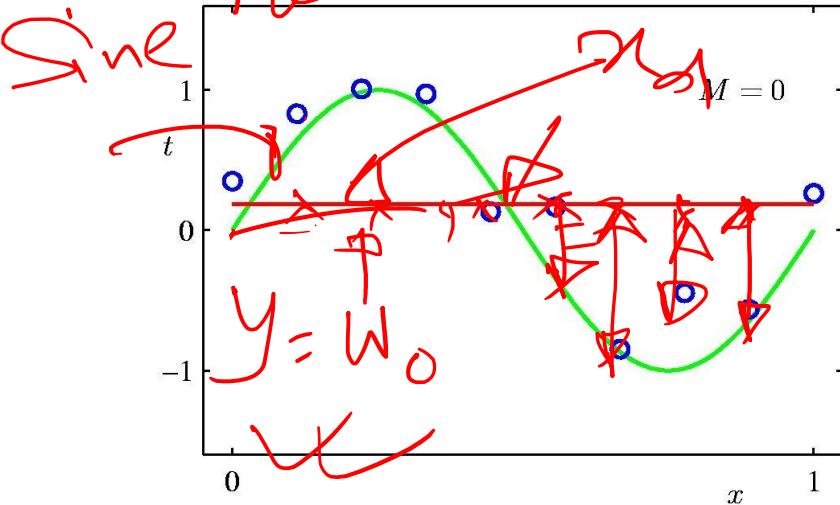
Linear Basis Models



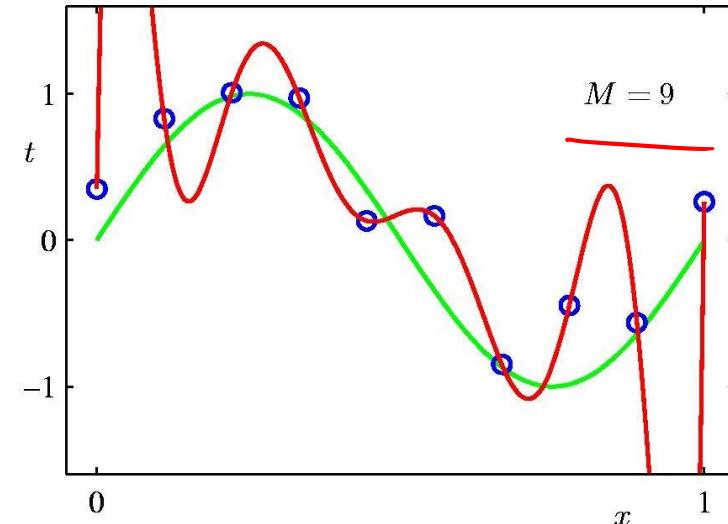
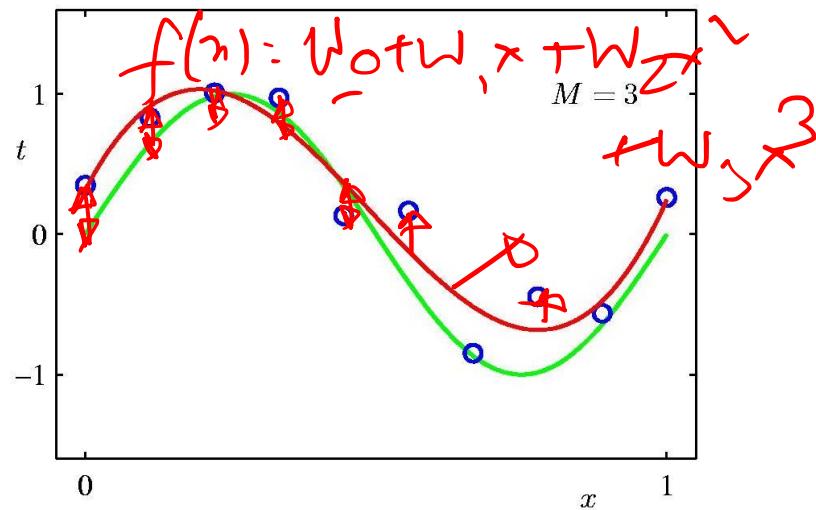
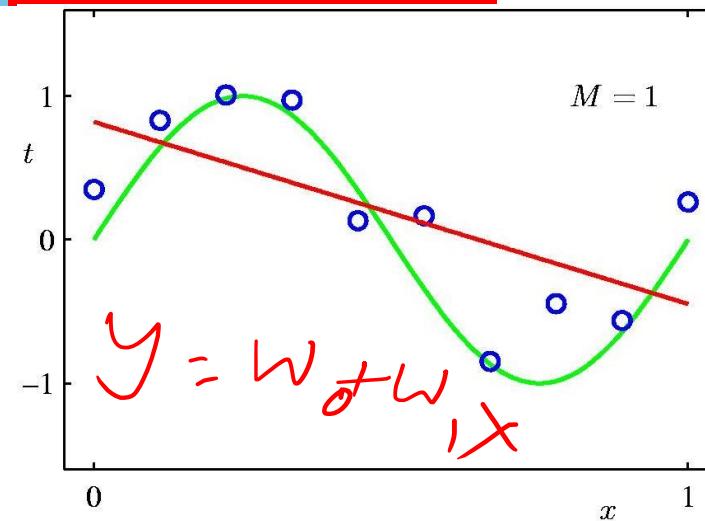
Polynomial Regression



Sine Relationship



$$y = \sin(\pi x)$$



Linear Basis Function Models

- The inputs **X** for linear regression can be:
 - Original quantitative inputs
 - Transformation of quantitative inputs
 - e.g. log, exp, square root, square, etc.
 - Polynomial transformation
 - example: $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot x^3$
 - Basis expansions
 - Dummy coding of categorical inputs
 - Interactions between variables
 - example: $x_3 = x_1 \cdot x_2$

This allows use of linear regression techniques to fit non-linear datasets.

X No.of.Years of Experience (in Years)	X^2	Y Salary Of the Employee (in Lakhs)
1	1	2
2	4	3
3	9	4
4	16	5
5	25	6

X1 = Graduate	X2 = PostGraduate	X3 = Others	Y Salary Of the Employee
0	0	1	2
1	0	0	3
0	0	1	4
0	1	0	5
1	0	0	6

$f(x) = w_0 + w_1 x$ *is linear in w*

Linear Basis Function Models

$$f(x) = w_0 + w_1 x + w_2 x^2 \quad \text{is a quadratic model}$$

Example: an M-th order polynomial function of one dimensional feature x :

$$y(x, w) = w_0 + \sum_{j=1}^M w_j x^j$$

$$\phi_0(u) = 1$$

x No.of.Years of Experience (in Years)	x^2	y Salary Of the Employee (in Lakhs)
1	1	2
2	4	3
3	9	4
4	16	5
5	25	6

where $x^j = j\text{-th power of } x$

$$y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

$\phi_j(x)$ are known as *basis functions*. Typically, $\phi_0(x) = 1$, so that w_0 acts as a bias.

In the simplest case, we use linear basis functions : $\phi_d(x) = x_d$.

They are called linear models because this function is linear in w .

$$y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x)$$

$$\phi_1 = x$$

$$\phi_2 = x^2$$

$$\phi_3 = x^3$$

$$x_d = \text{Feature}$$

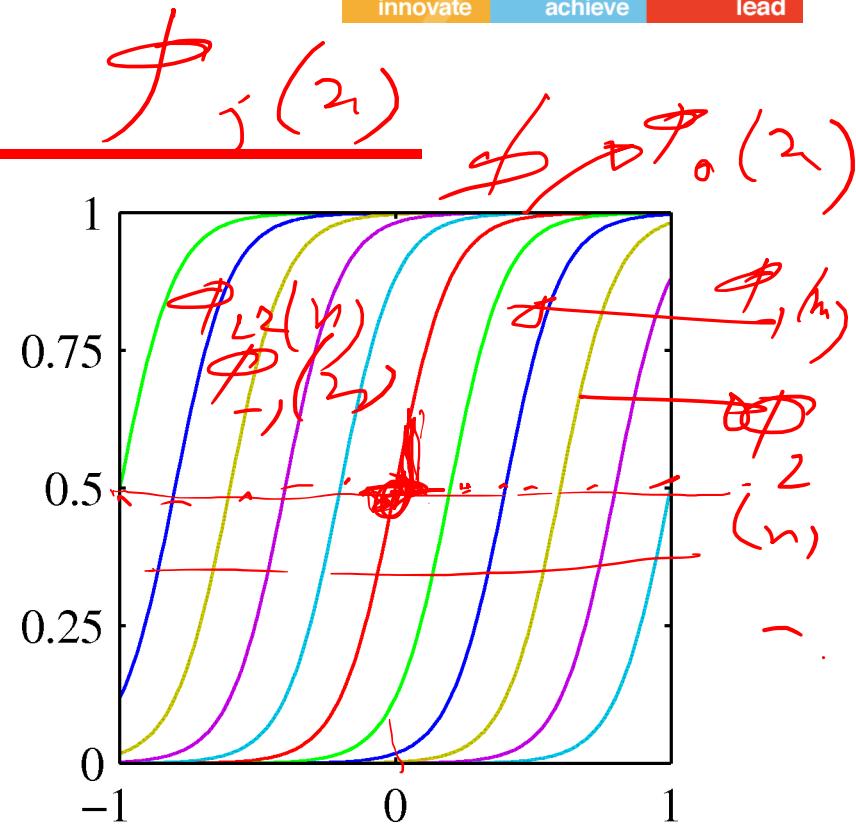
Linear Basis Function Models - Examples

Sigmoidal basis functions:

where $\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Also these are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (slope).



$$\sigma^{(n)} = \frac{1}{1 + e^{-x}} = \frac{1}{1 + e^{-x/s}}$$

$$\phi_0(x) = \sigma\left(\frac{x}{s}\right)$$

Linear Basis Function Models - Examples

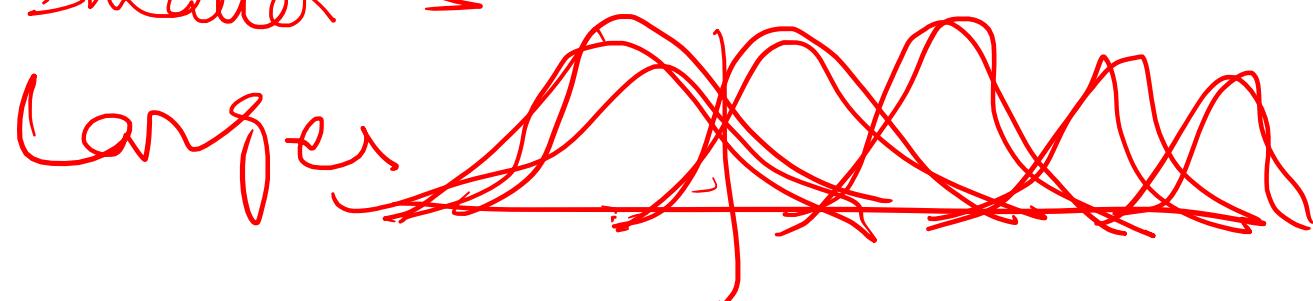
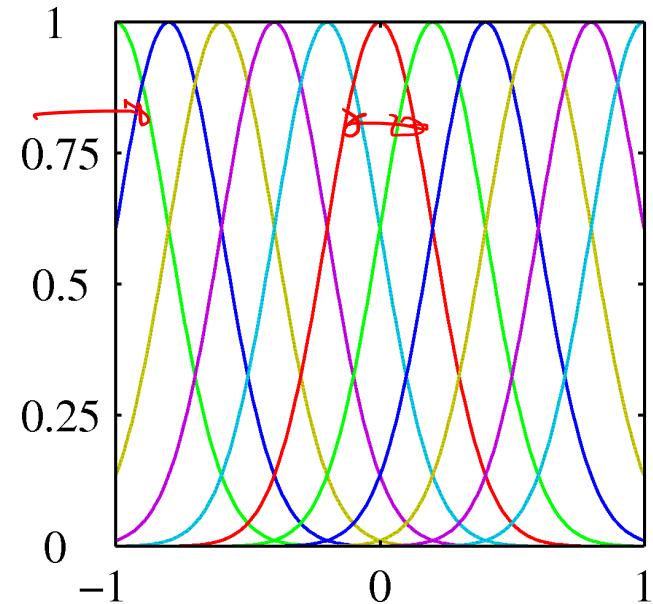
s is constant

Gaussian basis functions:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

These are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (width).

Smaller s
larger s

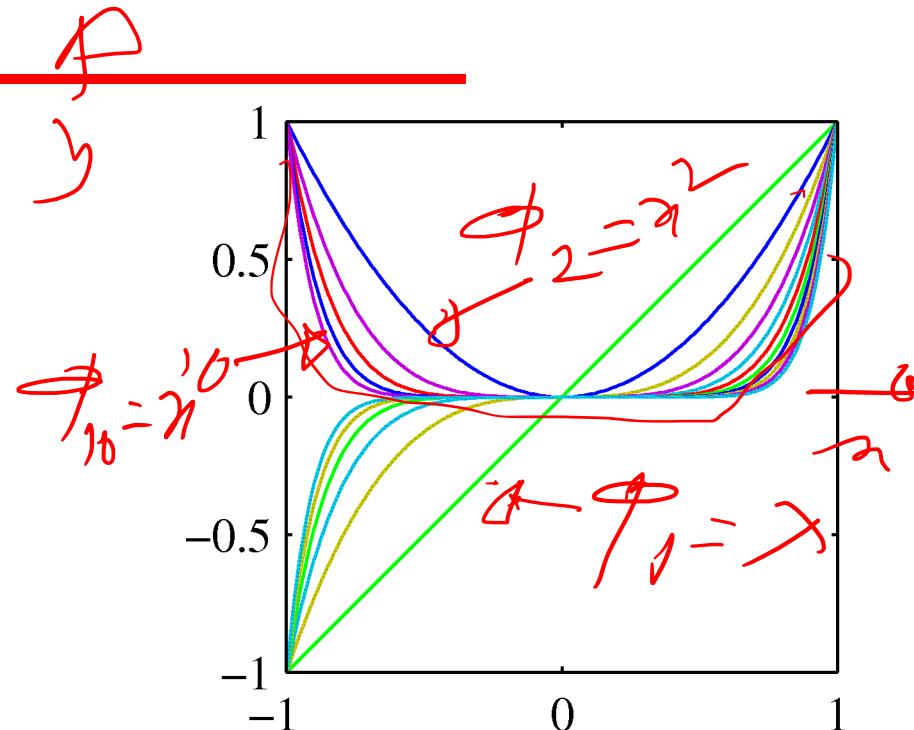


Linear Basis Function Models - Examples

Polynomial basis functions:

$$\phi_j(x) = x^j.$$

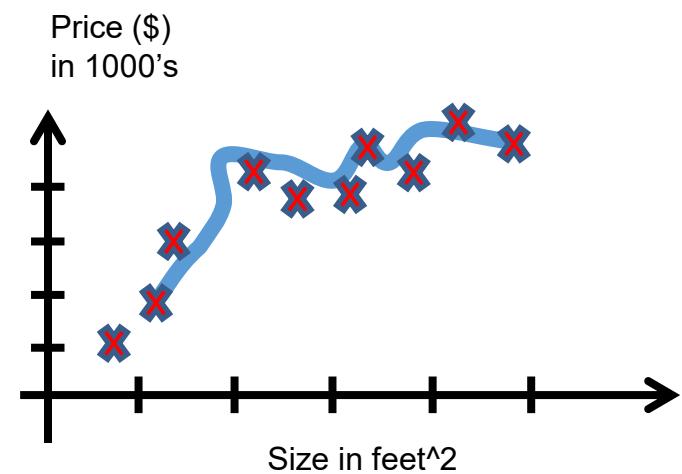
These are global; a small change in x affect all basis functions.



Notion of Bias - Variance

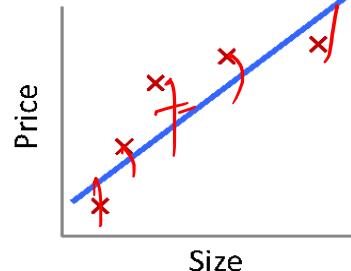
Addressing overfitting

- x_1 = size of house
- x_2 = no. of bedrooms
- x_3 = no. of floors
- x_4 = age of house
- x_5 = average income in neighborhood
- x_6 = kitchen size
- :
- x_{100}



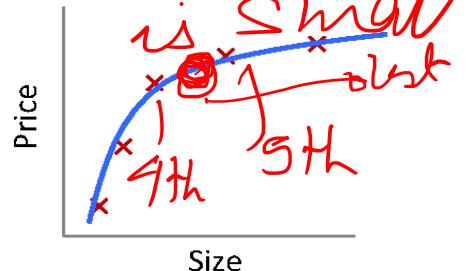
Quality of Fit

bias large



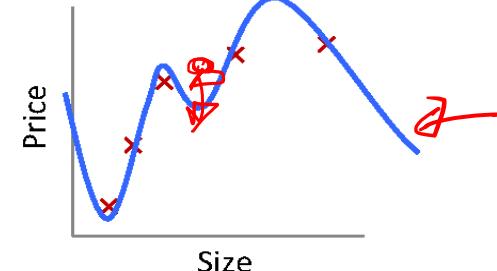
$\theta_0 + \theta_1 x$
Underfitting
(high bias)

bias



$\theta_0 + \theta_1 x + \theta_2 x^2$
Correct fit

bias = 0



$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
Overfitting
(high variance)

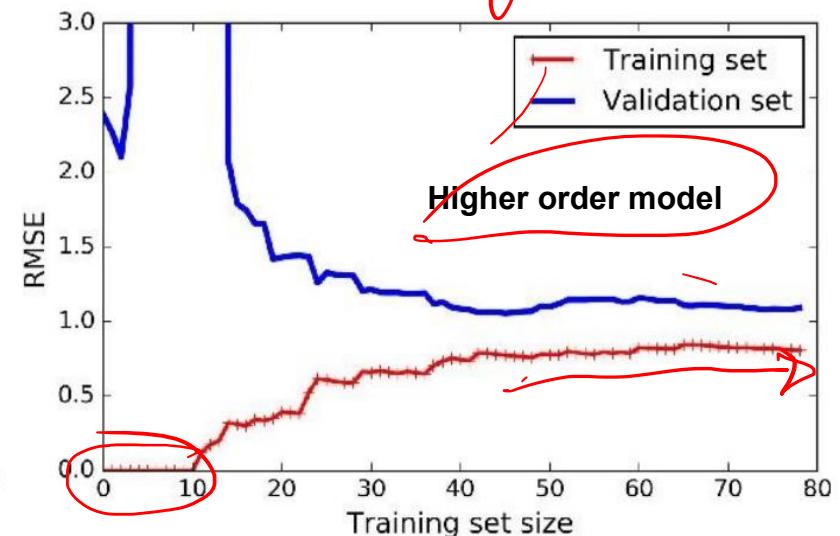
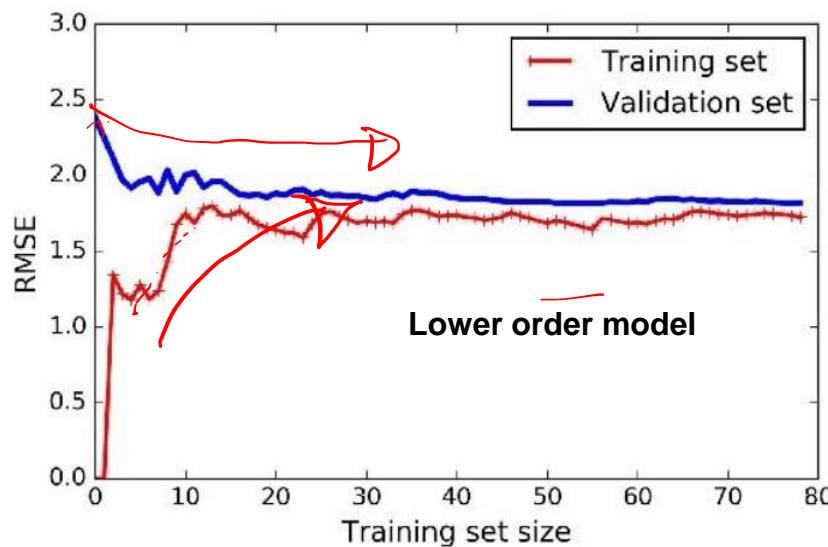
Overfitting:

- The learned hypothesis may fit the training set very well ($J(\theta) \approx 0$)
- ...but fails to generalize to new examples

Effect of Training Size on Over fitting

- Size of training dataset needs to be large to prevent when higher order model is used.

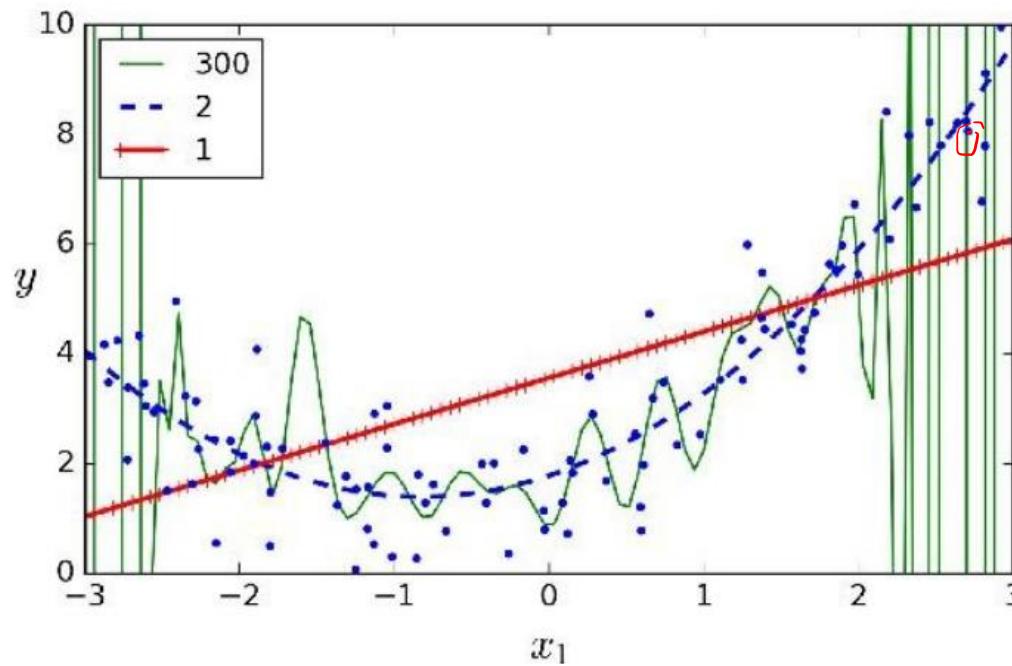
$\frac{1}{n} \sum_i (pred - act)^2$
 $i \neq \text{of training}$



Polynomial Fitting can lead to Over fitting

- Underlying target function is quadratic
- Linear model results in under fitting with large bias
- Polynomial of order 300 results in a large variance

$\sum_{i=1}^m (act - pred)^2$
Validation



Regularization

Regularization

- A method for automatically controlling the complexity of the learned hypothesis
- **Idea:** penalize for large values of θ_j
 - Can incorporate into the cost function
 - Works well when we have a lot of features, each that contributes a bit to predicting the label
- Can also address overfitting by eliminating features (either manually or via model selection)

$$f_{h.al} = \sum_{i=1}^n w_i^2$$

$w_i \neq 0$

$$\min \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n w_i^2 \right)$$

Regularization

- Linear regression objective function

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

model fit to data regularization

- λ is the regularization parameter ($\lambda \geq 0$)
 - No regularization on θ_0 !

Ridge

for L_2 regularization $T = \dots$

Understanding Regularization

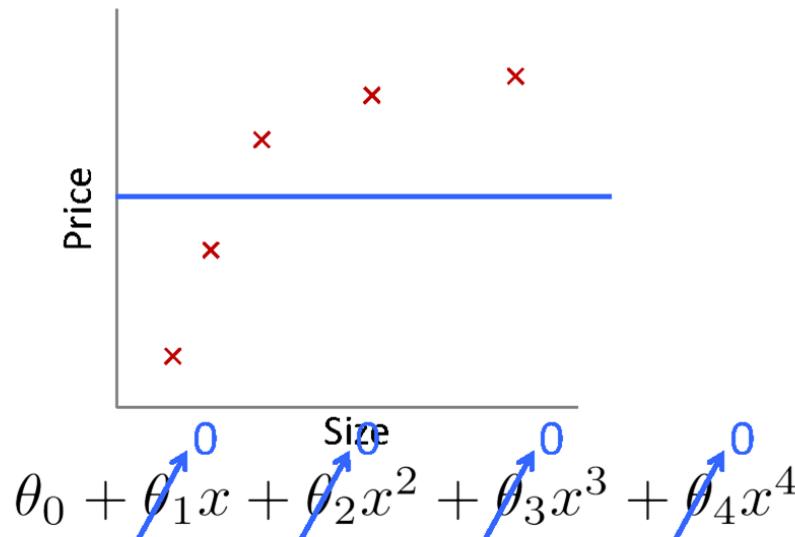
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- Note that $\sum_{j=1}^d \theta_j^2 = \|\boldsymbol{\theta}_{1:d}\|_2^2$
 - This is the magnitude of the feature coefficient vector!
- We can also think of this as:
$$\sum_{j=1}^d (\theta_j - 0)^2 = \|\boldsymbol{\theta}_{1:d} - \vec{0}\|_2^2$$
 - L₂ regularization pulls coefficients toward 0

Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- What happens if we set λ to be huge (e.g., 10^{10})?



Based on example by Andrew Ng

Regularized Linear Regression

- Cost Function

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

- Fit by solving $\min_{\theta} J(\theta)$

- Gradient update:

$$\begin{aligned} \frac{\partial}{\partial \theta_0} J(\theta) & \quad \theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) \\ \frac{\partial}{\partial \theta_j} J(\theta) & \quad \theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} \end{aligned}$$

$\lambda \theta_j$

- We can rewrite the gradient step as:

$$\theta_j \leftarrow \theta_j (1 - \alpha \lambda) - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

Lasso Regularization

\hookrightarrow regularization

Feature Selection

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^d |\theta_j|$$

$$\theta_j = \theta_j - \frac{\alpha}{n} \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) \mathbf{x}_j^{(i)} - \alpha \lambda \text{sign}(\theta_j)$$

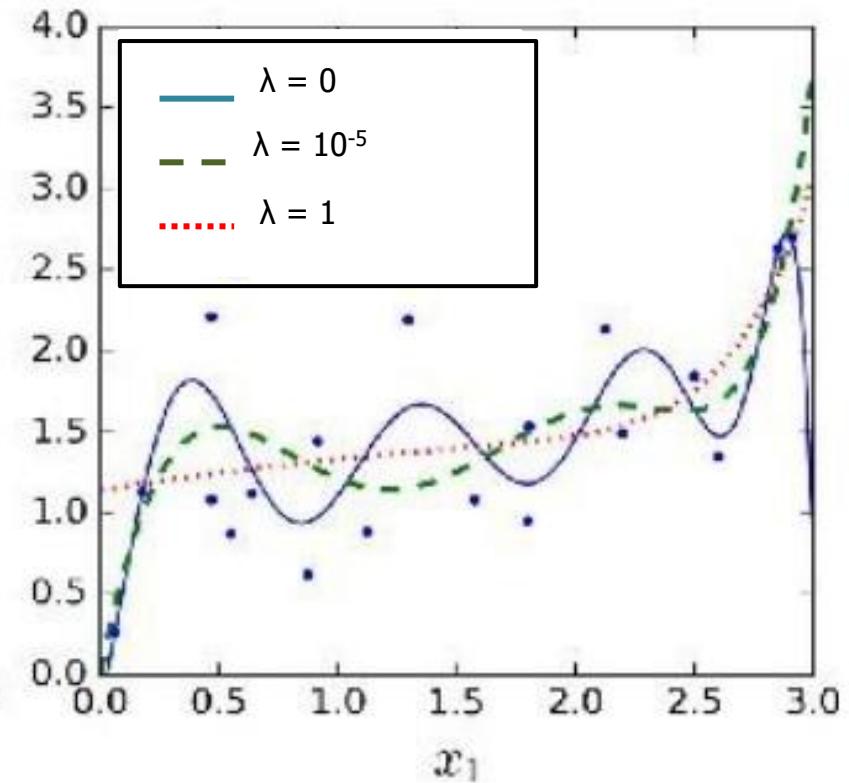
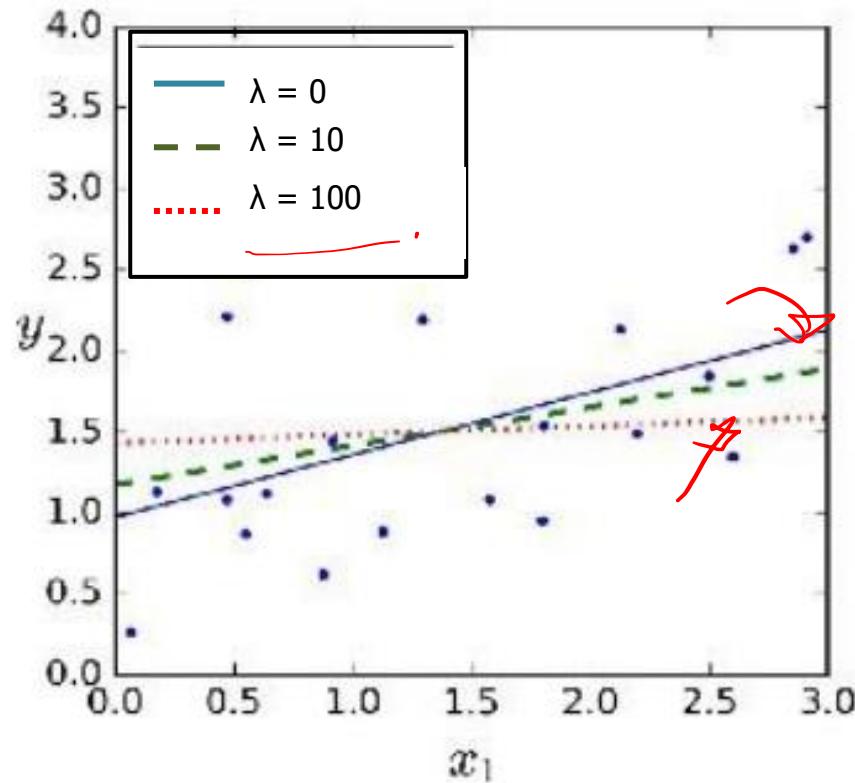
if assoc.
with some
feature $\Rightarrow 0$

$-\lambda$ if $\theta_j > 0$

$+\lambda$ if $\theta_j < 0$

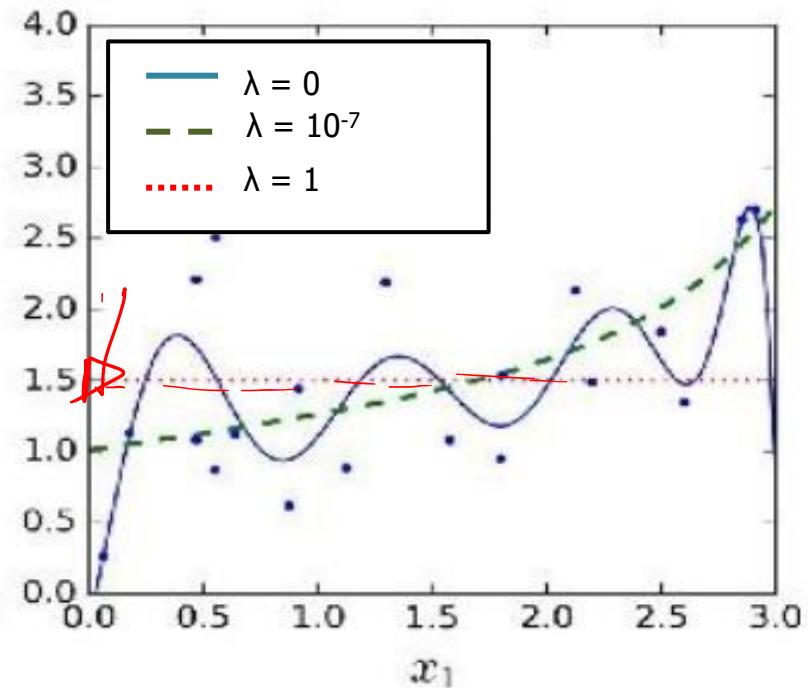
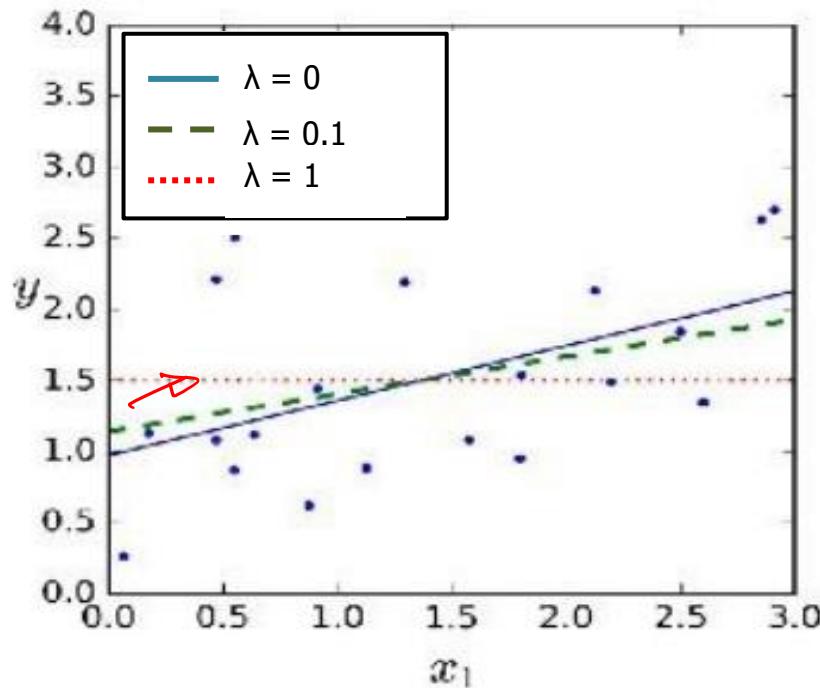
Ridge Vs Lasso Regularization

w
 l_1
 l_2



Ridge

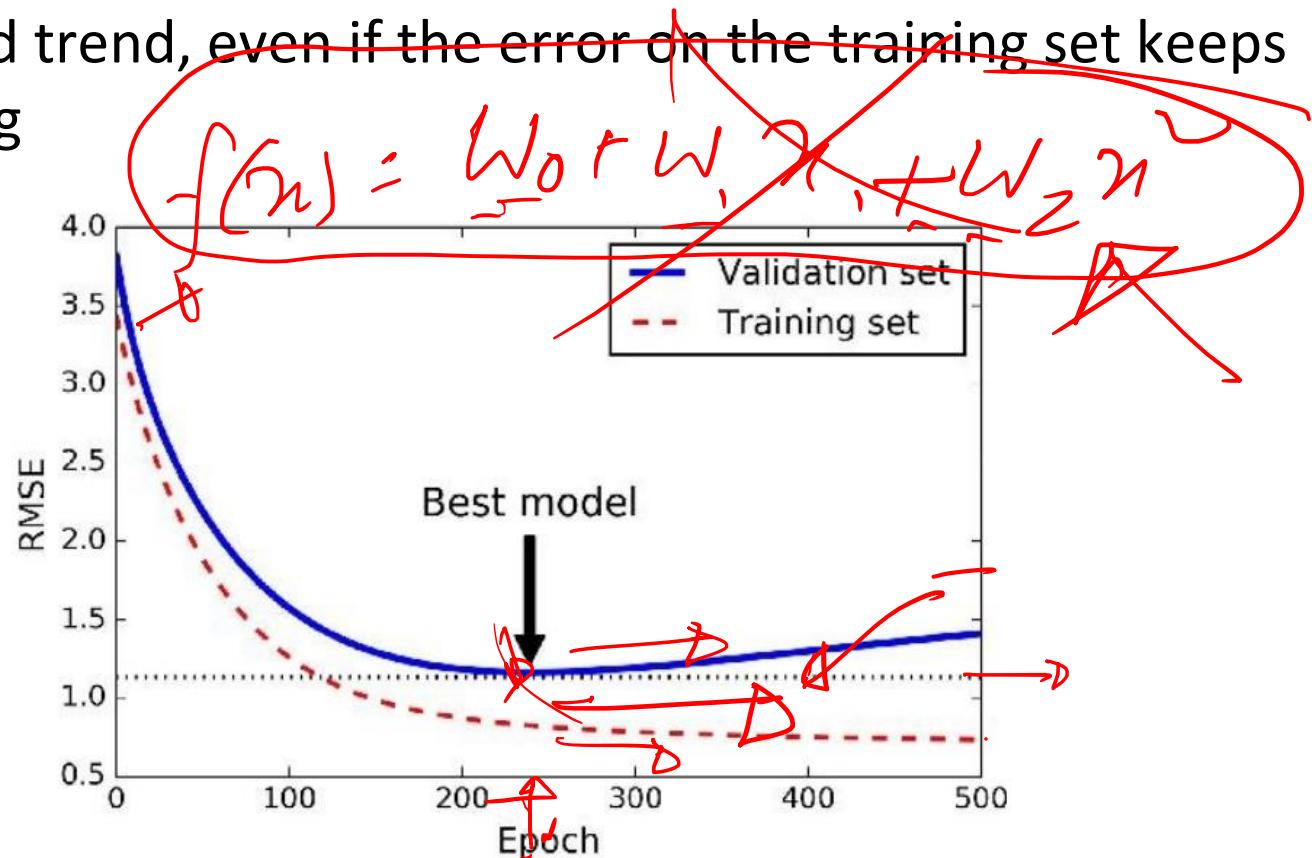
Ridge Vs Lasso Regularization



Lasso

Early Stopping

- Do Not Over train to prevent overfitting
- Stop training once error on the validation set starts showing an upward trend, even if the error on the training set keeps decreasing



How Good is the Fit ?

- Coefficient of Determination

$$R^2 = \frac{SSM}{SST} = \frac{\sum_i [f(x_i) - \bar{y}]^2}{\sum_i [y_i - \bar{y}]^2}$$

*predicted av ...
 av. target output
 org. target output*

- R^2 is related to correlation coefficient

$$r = \frac{\sigma_{xy}}{\sqrt{\sigma_{xx}\sigma_{yy}}}.$$

$$R^2 = \frac{\sigma_{xy}^2}{\sigma_{xx}\sigma_{yy}}.$$

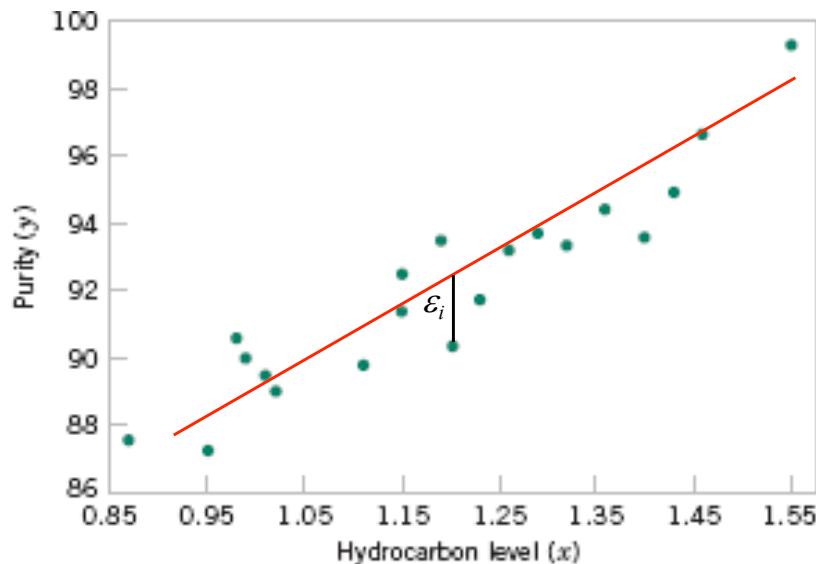
of data point

$$\text{Adjusted } R^2 = 1 - \left(\frac{N-1}{N-d} \right) (1 - R^2),$$

*Polynomial
 order of the*

Simple linear regression

Based on the scatter diagram, it is probably reasonable to assume that the mean of the random variable Y is related to X by the following **simple linear regression model**:



Response Regressor or Predictor
 $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ $i = 1, 2, \dots, n$
 Intercept Slope Random error
 $\varepsilon_i \sim N(0, \sigma^2)$

where the slope and intercept of the line are called **regression coefficients**.

- The case of simple linear regression considers a single regressor or predictor x and a dependent or response variable Y.

Regression coefficients

$$y - \beta_0 + \beta_1 x$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}$$

$$S_{xy} = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}$$

$$\begin{aligned} \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{S_{xy}}{S_{xx}} \end{aligned} \quad \left. \begin{array}{l} \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ \text{Fitted (estimated) regression model} \end{array} \right\}$$

Caveat: regression relationships are valid only for values of the regressor variable within the range of the original data. Be careful with extrapolation.

Estimation of variance

- Using the fitted model, we can estimate value of the response variable for given predictor

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

- Residuals: $r_i = y_i - \hat{y}_i$
- Our model: $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, i=1,\dots,n, \text{Var}(\varepsilon_i) = \sigma^2$
- Unbiased estimator (MSE: Mean Square Error)

$\hat{\sigma}^2 = \text{MSE} = \frac{\sum_{i=1}^n r_i^2}{n-2}$

$r_i = \hat{y}_i - y_i \rightarrow \text{target}$

$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

Key Point

- the coefficients

$$\hat{\beta}_1 \text{ and } \hat{\beta}_0$$

are both calculated from data, and they are subject to error.

- if the true model is $y = \beta_1 x + \beta_0$, $\hat{\beta}_1$ and $\hat{\beta}_0$ are point estimators for the true coefficients
- we can talk about the ``accuracy'' of $\hat{\beta}_1$ and $\hat{\beta}_0$

Assessing linear regression model

- Test hypothesis about true slope and intercept

$$\beta_1 = ?, \quad \beta_0 = ?$$

- Construct confidence intervals

$\alpha = 0.007$
 $= 0.993$

$$\beta_1 \in [\hat{\beta}_1 - a, \hat{\beta}_1 + a] \quad \beta_0 \in [\hat{\beta}_0 - b, \hat{\beta}_0 + b] \quad \text{with probability } 1 - \alpha$$

- Assume the errors are normally distributed

$$\varepsilon_i \sim N(0, \sigma^2)$$

Properties of Regression Estimators

slope parameter β_1

$$E(\hat{\beta}_1) = \beta_1$$

expected value of $\hat{\beta}_1$

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

intercept parameter β_0

$$E(\hat{\beta}_0) = \beta_0$$

expected value of $\hat{\beta}_0$

$$V(\hat{\beta}_0) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]$$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i \right)^2}{n}$$

Standard errors of coefficients



- We can replace σ^2 with its estimator $\hat{\sigma}^2$...

$$\hat{\sigma}^2 = MSE = \frac{\sum_{i=1}^n r_i^2}{n-2}$$

$$r_i = y_i - \hat{y}_i \quad \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

- Using results from previous page, estimate the

$$se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \quad \text{and} \quad se(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}$$

Hypothesis test in simple linear regression

- we wish to test the hypothesis whether the slope equals a constant $\beta_{1,0}$

$$H_0: \beta_1 = \beta_{1,0}$$

$$H_1: \beta_1 \neq \beta_{1,0}$$

- e.g. relate **ads** to **sales**, we are interested in study whether or not increase a \$ on ads will increase \$ $\beta_{1,0}$ in sales?
- $\text{sale} = \beta_{1,0} \text{ ads} + \text{constant?}$



A related and important question...



- whether or not the slope is zero ?

$$H_0: \beta_1 = 0$$

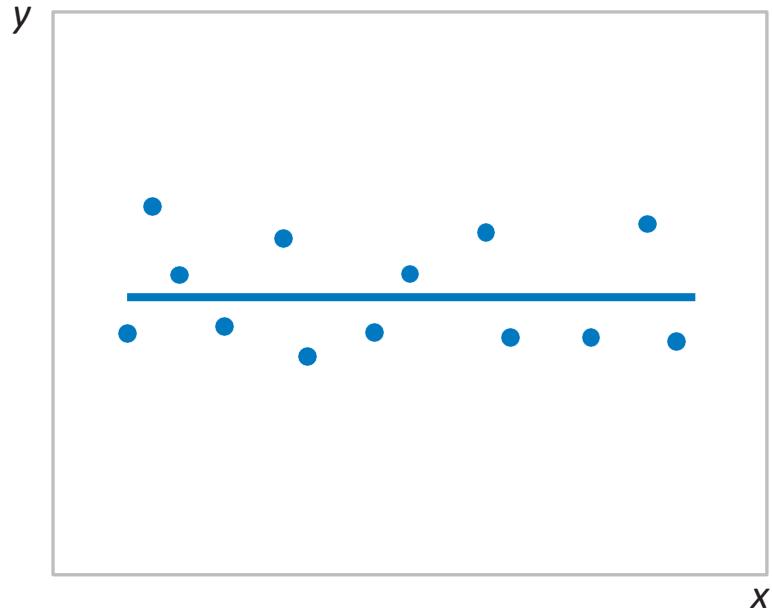
$$H_1: \beta_1 \neq 0$$

- if $\beta_1 = 0$, that means Y does not depend on X, i.e.,
- Y and X are **independent**
- In the advertisement example, does **ads increase sales?** or no effect?

Significance of regression

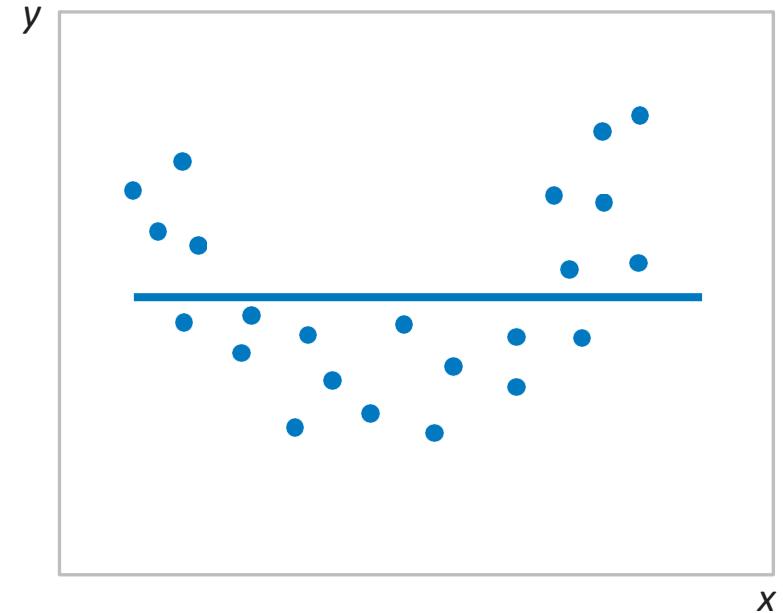
The screenshot shows a Google search results page for the query "beginner yoga classes". The results are categorized by type: Web, Images, Maps, Videos, News, and More. Several results are highlighted with red boxes:

- Laura Yoga Studio** (646) 702-4596
www.laurayoga.com
Great for beginners. Get the first 3 classes free! Call now.
- Youth Yoga Classes**
www.yogakids.com
Yoga for all ages! We offer modern facilities and reasonable rates
Yoga Kids Inc. – 610 McKenzie Blvd. Denver, CO
- Yoga Accessories**
www.yogaaccessories.com
Experts or **beginners**, we have everything you need for **yoga**.
- Yoga Yoga Denver**
www.yogayogadenviers.com
Yoga classes in denver. New to **Yoga**? Start here! Mommy & baby **yoga**!
Map & directions to studio · Rent our Space · Energy/Exchange opportunities
- Yoga Basics: Your guide to the Practice of Yoga**
www.yogabasics.com
- Hot Yoga Classes**
www.yogabears.com/hotyoga
Dynamic, fun and cost effective!
Special: 10 classes for \$100
- Yoga for beginners**
www.vinashiyoga.com
Burn calories and find peace.
Small classes. First week free!
(354) 555-0111 - Directions
- Lilac Yoga Studio**
www.lilacyogadenvier.com
Try our popular **yoga** sessions
Limited time \$100 for 10!



(a)

- H_0 not rejected



(b)

- H_0 rejected

Use t-test for slope


$$\hat{\sigma}_e^2 \approx \hat{\sigma}^2 = \frac{\sum e_i^2}{N-2}$$

Under H_0

slope parameter β_1

$$E(\hat{\beta}_1) = \beta_{1,0}$$

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

$$\hat{\beta}_1 \sim N\left(\beta_{1,0}, \frac{\sigma^2}{S_{xx}} \right)$$

- Under H_0 , test statistic

$$T_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{\hat{\sigma}^2 / S_{xx}}} \rightarrow D$$

~ t distribution with
 $n-2$ degree of freedom

- Reject H_0 if

$$|t_0| > t_{\alpha/2, n-2}$$

(two-sided test)

Use t-test for intercept

- Use a similar form of test

$$H_0: \beta_0 = \beta_{0,0}$$

$$H_1: \beta_0 \neq \beta_{0,0}$$

- **Test statistic** $T_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{\sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}} = \frac{\hat{\beta}_0 - \beta_{0,0}}{se(\hat{\beta}_0)}$

Under H_0 , $T_0 \sim t$ distribution with $n-2$ degree of freedom

- Reject H_0 if $|t_0| > t_{\alpha/2, n-2}$

Confidence interval

- we can obtain confidence interval estimates of slope and intercept
- width of confidence interval is a measure of the overall quality of the regression

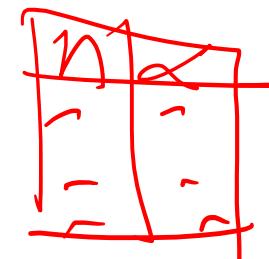
		true parameter
slope	intercept	
$T_0 = \frac{\hat{\beta}_1 - \boxed{\beta_{1,0}}}{\sqrt{\hat{\sigma}^2 / S_{xx}}}$	$T_0 = \frac{\hat{\beta}_0 - \boxed{\beta_{0,0}}}{\sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}}$	
$\sim t$ distribution with $n-2$ degree of freedom	$\sim t$ distribution with $n-2$ degree of freedom	!0

Confidence intervals

Conf. interval α

t-distribution

bivariate



a $100(1 - \alpha)\%$ confidence interval on the slope β_1

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

a $100(1 - \alpha)\%$ confidence interval on the intercept β_0

$$\begin{aligned} \hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \\ \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} \end{aligned}$$

Example: house selling price and annual taxes

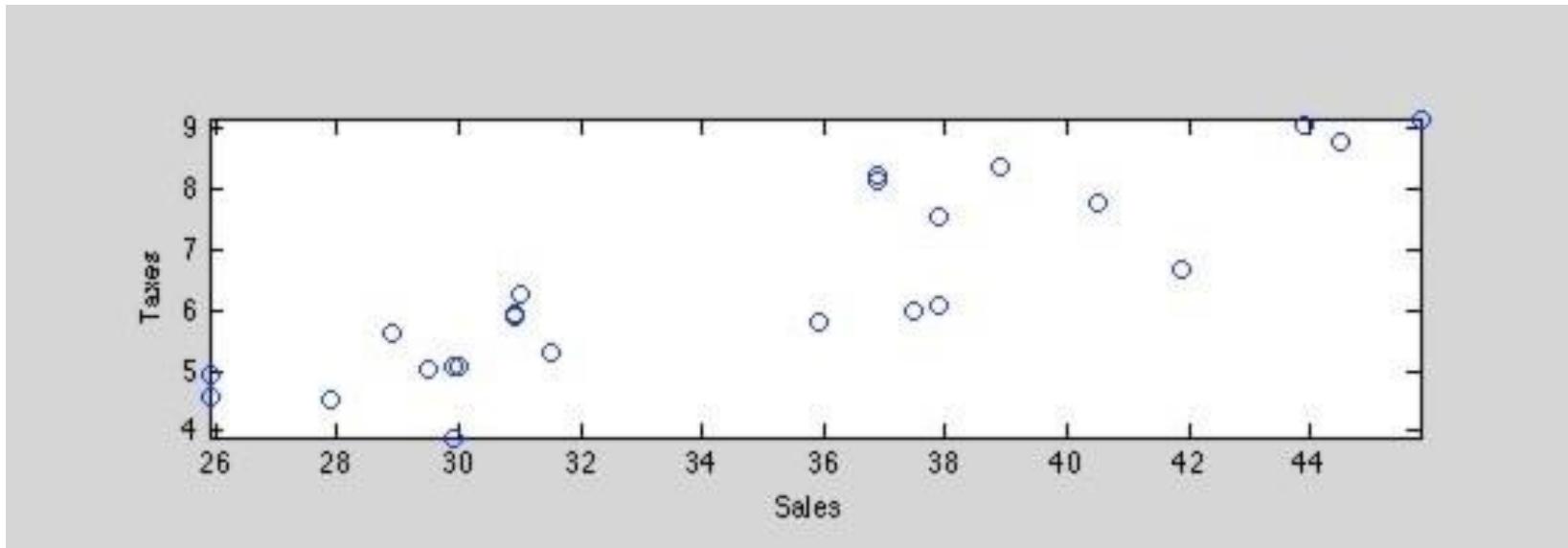
Sale Price/1000	Taxes (Local, School), County/1000	Sale Price/1000	Taxes (Local, School), County/1000
25.9	4.9176	30.0	5.0500
29.5	5.0208	36.9	8.2464
27.9	4.5429	41.9	6.6969
25.9	4.5573	40.5	7.7841
29.9	5.0597	43.9	9.0384
29.9	3.8910	37.5	5.9894
30.9	5.8980	37.9	7.5422
28.9	5.6039	44.5	8.7951
35.9	5.8282	37.9	6.0831
31.5	5.3003	38.9	8.3607
31.0	6.2712	36.9	8.1400
30.9	5.9592	45.8	9.1416



Independent variable X: SalePrice

Dependent variable Y: Taxes

- qualitative analysis



Calculate correlation

$$\hat{\rho} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \times \sum_{i=1}^n (y_i - \bar{y})^2}} = 0.8760$$

Independent variable Y: SalePrice

Dependent variable X: Taxes

$$\underline{n = 24} \quad \underline{\bar{x} = 34.6125} \quad \underline{\bar{y} = 6.4049}$$

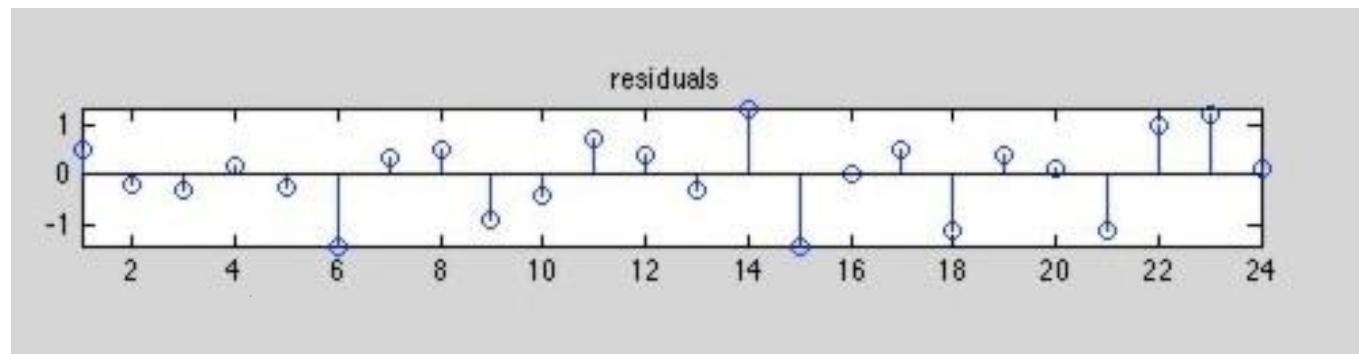
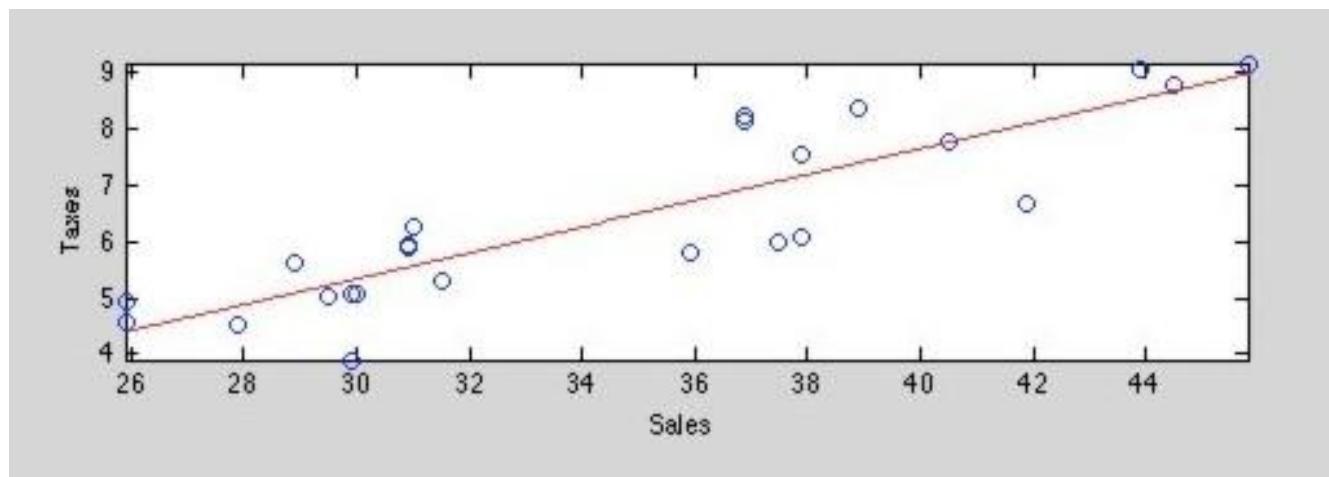
$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = 829.0462$$

$$S_{xy} = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) = 191.3612$$

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{191.3612}{829.0462} = 0.2308$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 6.4049 - 0.2308 \times 34.6125 = -1.5837$$

Fitted simple linear regression model $\hat{y} = -1.5837 + 0.2308x$



$$\sum_{i=1}^n r_i^2$$

- residuals: $\hat{\sigma}^2 = MSE = \frac{\sum_{i=1}^n r_i^2}{n-2} = 0.6088$

- standard error of regression coefficients
-

$$se(\hat{\beta}_1) = \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} = \sqrt{\frac{0.6088}{829.0462}} = 0.0271$$

$$se(\hat{\beta}_0) = \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]} = \sqrt{0.6088 \left[\frac{1}{24} + \frac{34.6125^2}{829.0462} \right]} = 0.9514$$

- test

Confidence 95%

Test $H_0: \beta_1 = 0$ using the t-test; use $\alpha = 0.005$

- calculate test statistics

$t_{\alpha/2}$

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2/S_{xx}}} = \frac{\hat{\beta}_1}{se(\hat{\beta}_1)} = \frac{0.2308}{0.0271} = 8.5166$$

- threshold

$$t_{\alpha/2, n-2} = t_{0.0025, 22} = 3.119$$

- value of test statistic is greater than threshold
-  reject H_0

- construct confidence interval for slope parameter

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}$$

$$t_{\alpha/2, n-2} = t_{0.0025, 22} = 3.119$$

$$0.2308 - 3.119 \times 0.0271 \leq \beta_1 \leq 0.2308 + 3.119 \times 0.0271$$

~~0.2308 - 3.119 × 0.0271 ≤ β₁ ≤ 0.2308 + 3.119 × 0.0271~~

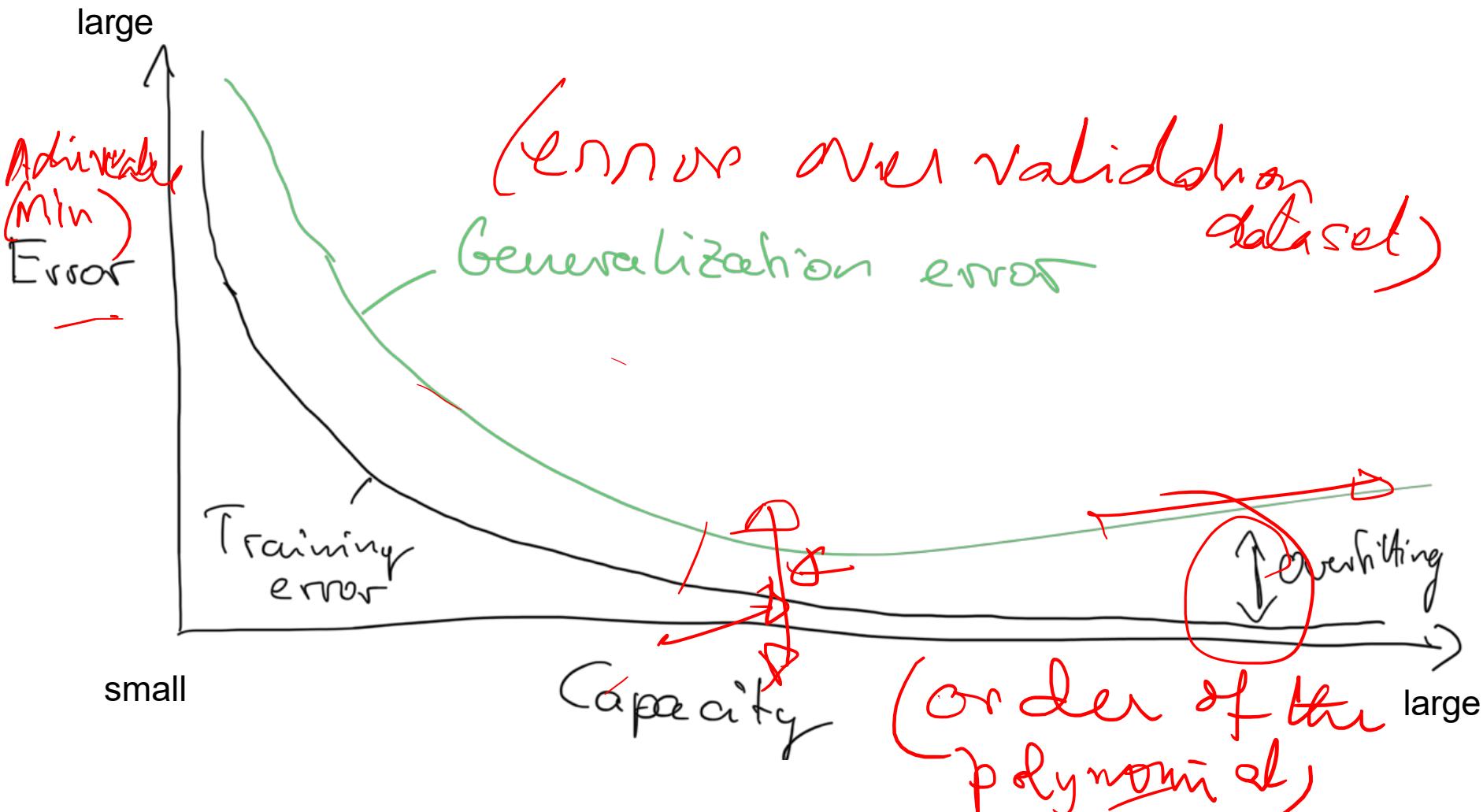
$$0.14631 \leq \beta_1 \leq 0.3153$$

99.5%

Bias-Variance Decomposition and Bias-Variance Trade- off

(and how it related to overfitting and underfitting)

Overfitting and Underfitting



Bias-Variance Decomposition



- Decomposition of the loss into bias and variance help us understand learning algorithms, concepts are related to underfitting and overfitting
- Helps explain why ensemble methods might perform better than single models

Bias-Variance Intuition

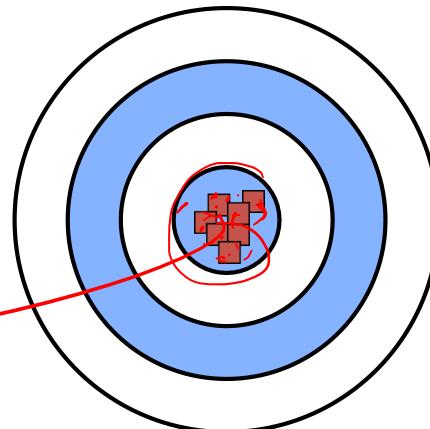


Low Variance
(Precise)

High Variance
(Not Precise)

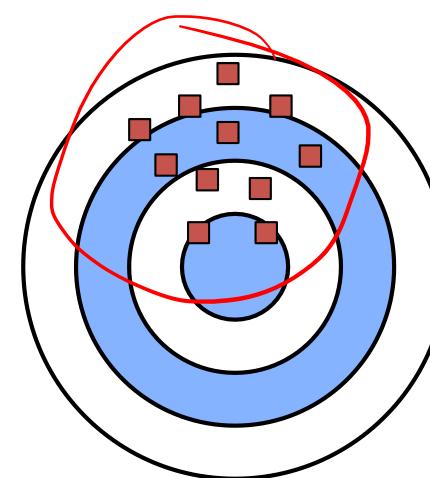
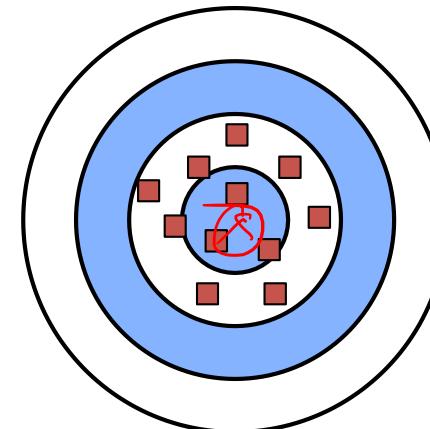
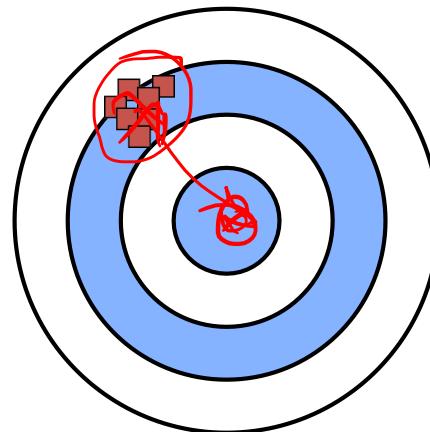
$(\hat{\beta}_0, \hat{\beta}_1)$

Low Bias
(Accurate)

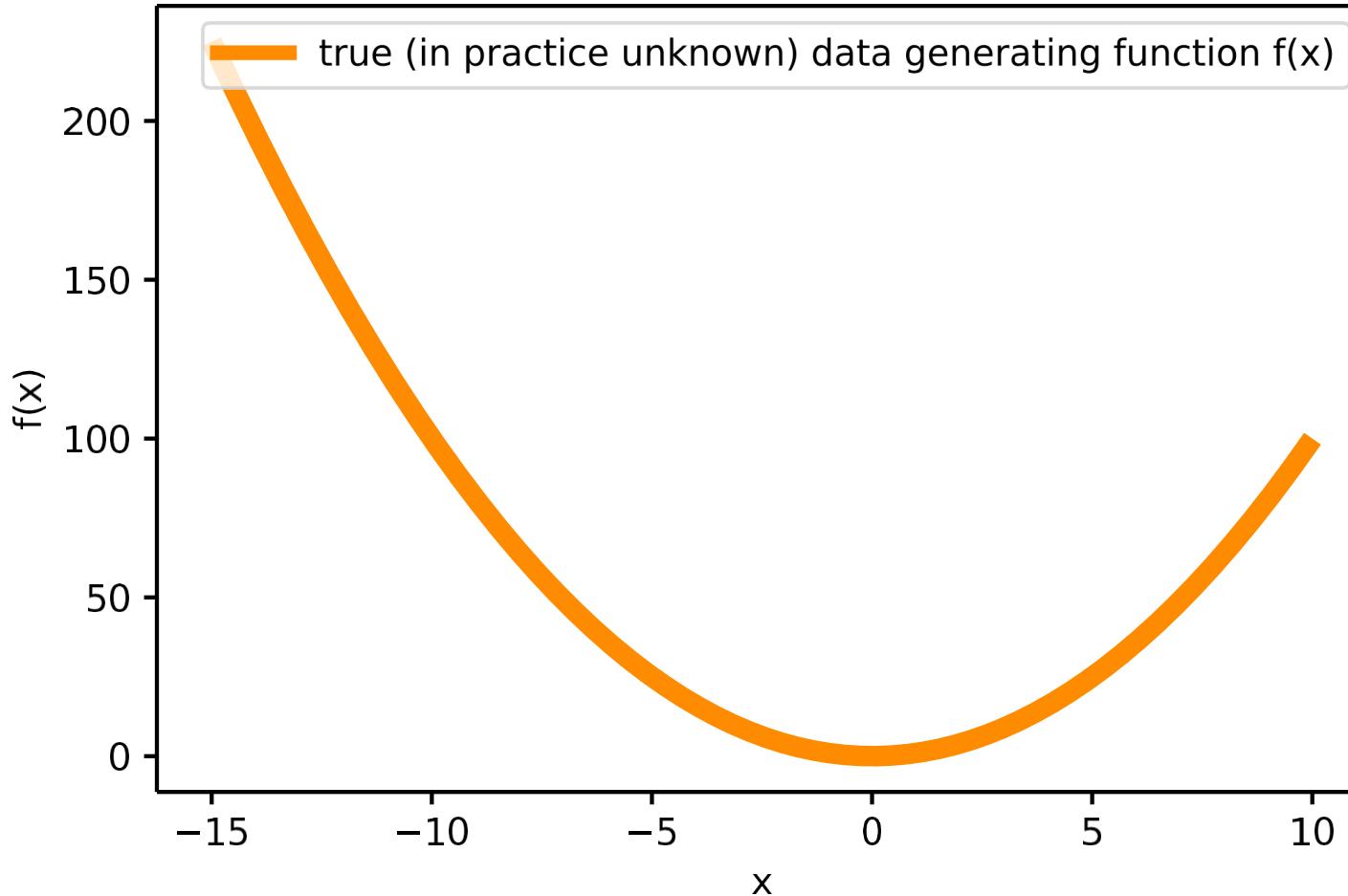


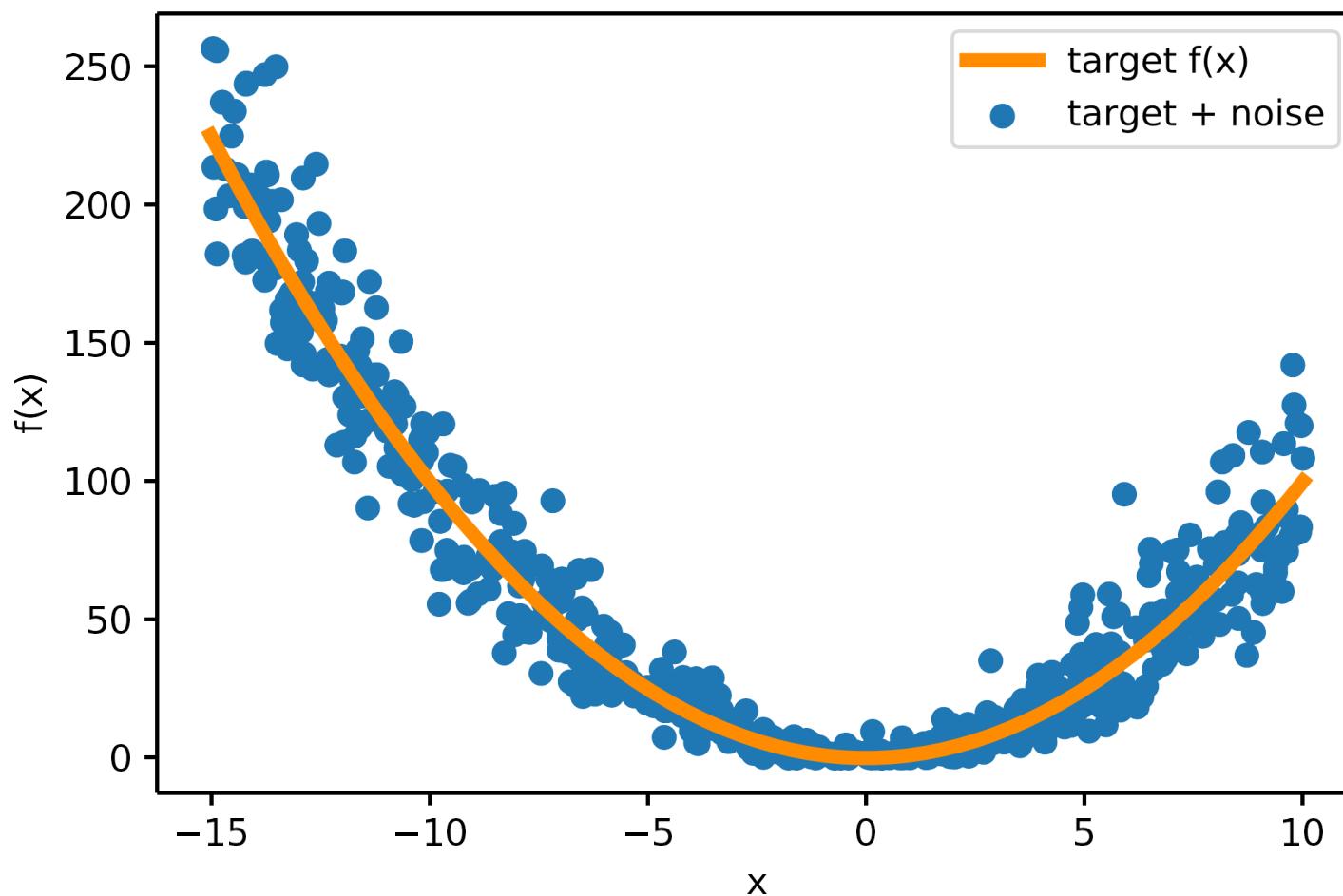
$(\hat{\beta}_0, \hat{\beta}_1)$

High Bias
(Not Accurate)

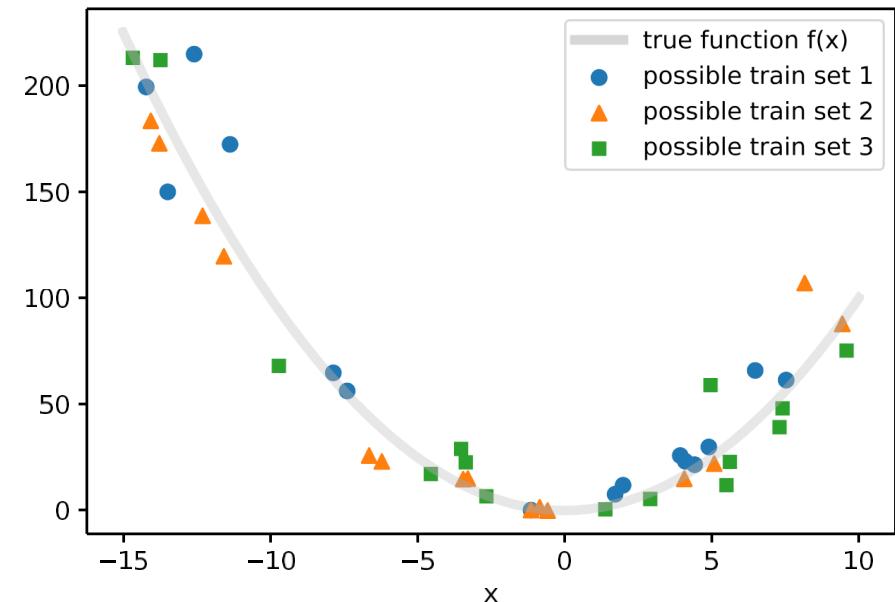
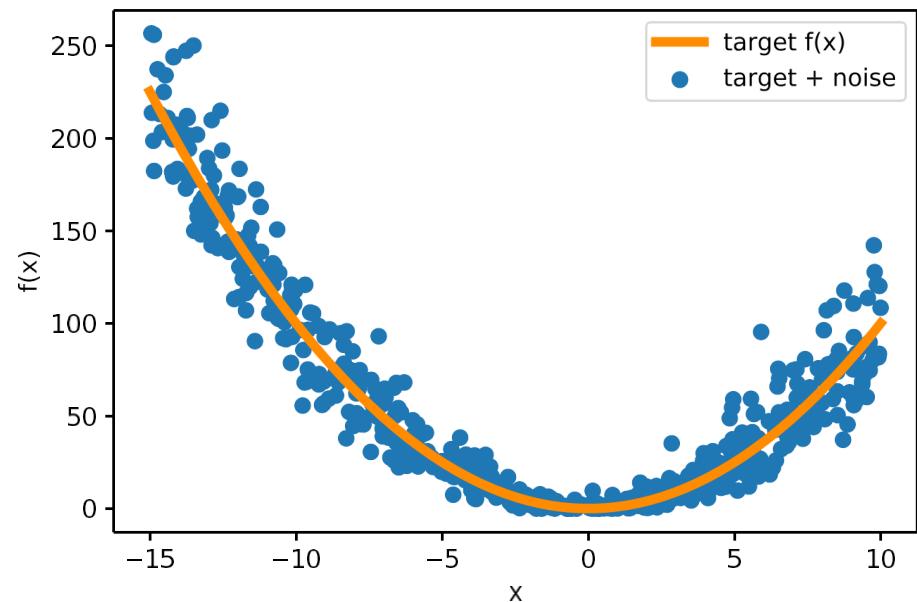


Bias-Variance Intuition

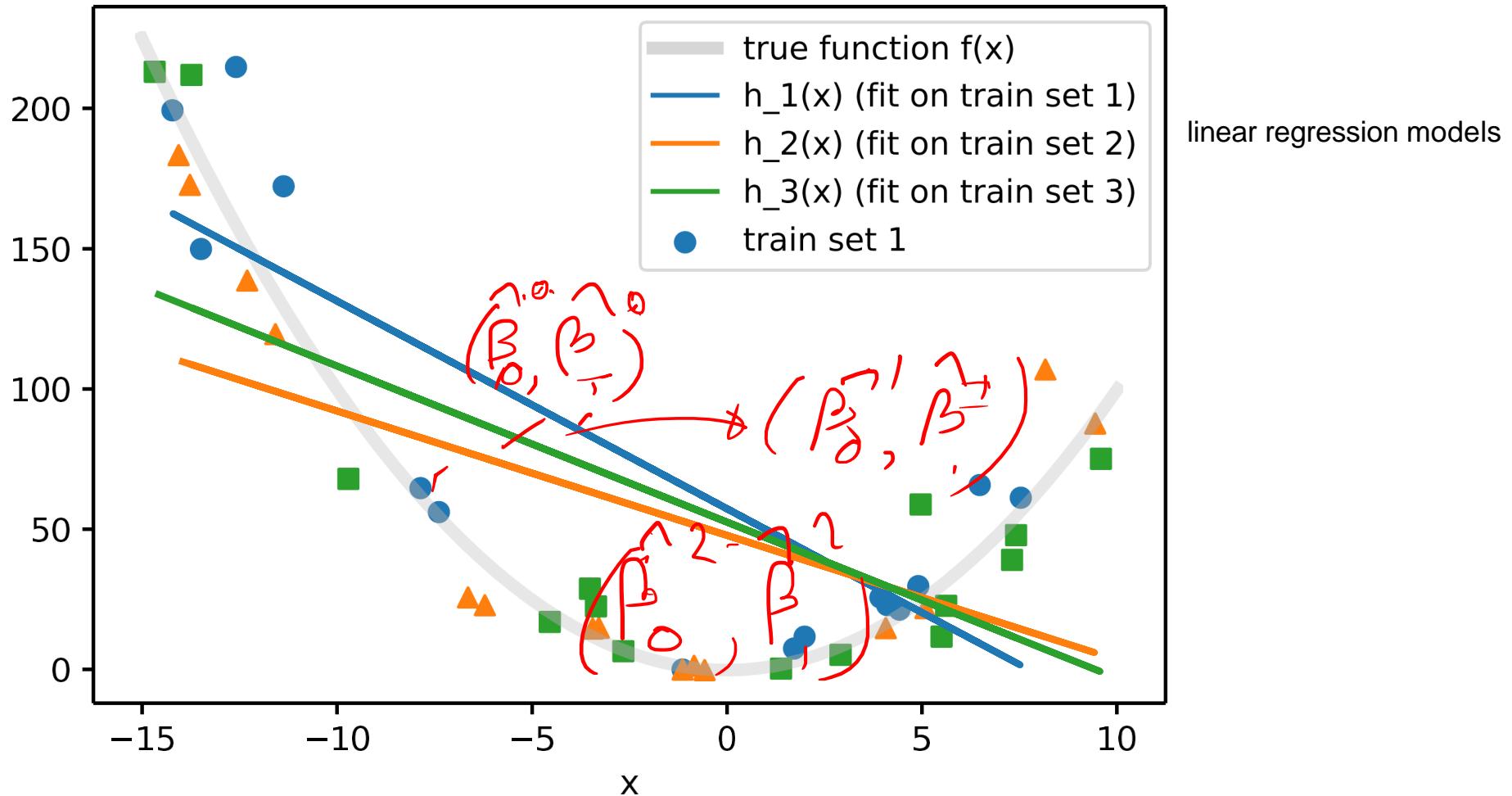




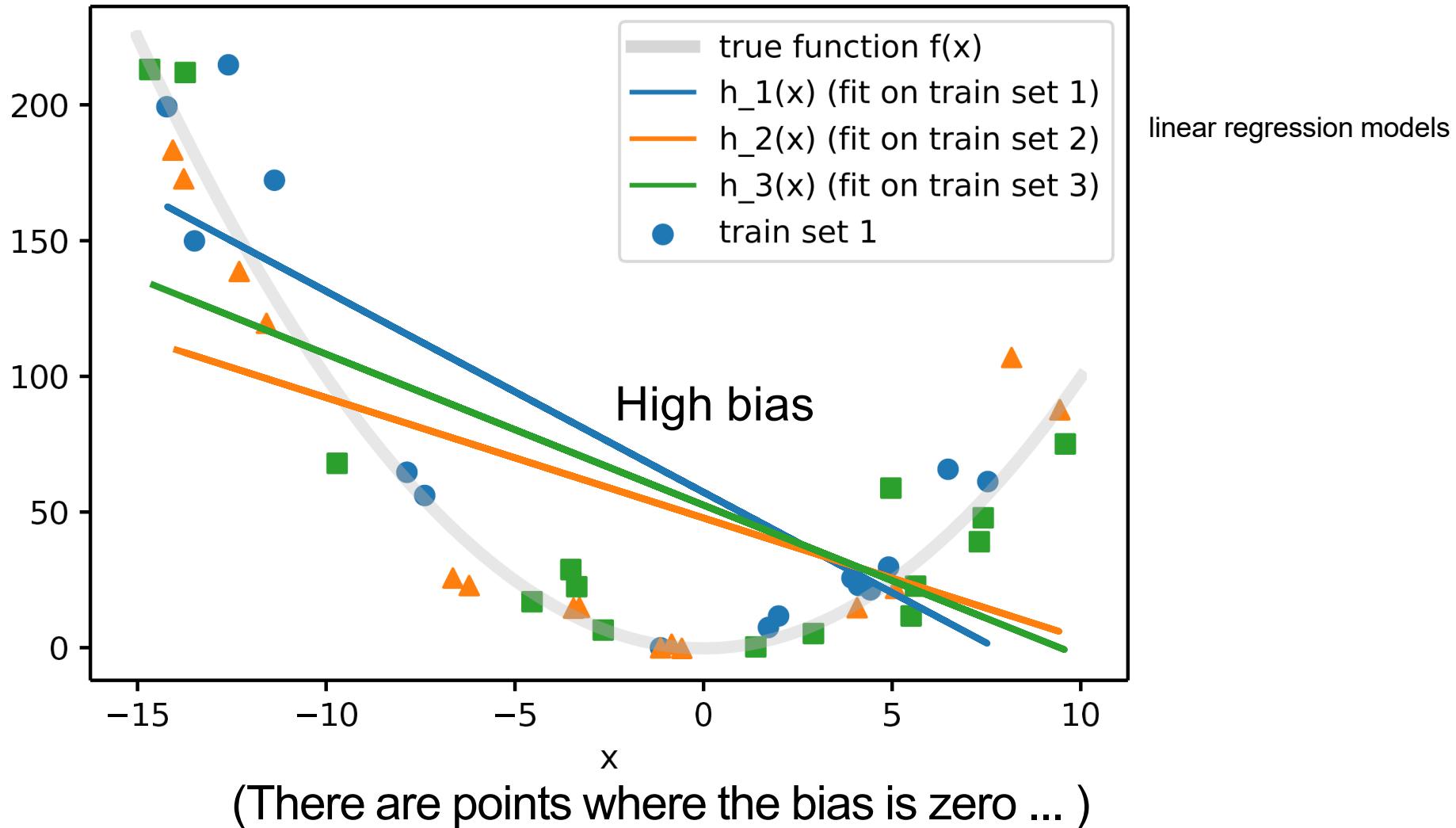
Bias-Variance Intuition



Bias-Variance Intuition



Bias-Variance Intuition



Terminology



Point estimator $\hat{\theta}$ of some parameter



(could also be a function, e.g., the hypothesis is an estimator of some target function)

Terminology



Point estimator $\hat{\theta}$ of some parameter

(could also be a function, e.g., the hypothesis is

~~$E(\theta)$~~ ¹²³
~~an estimator of some target function~~)

Bias = $E[\hat{\theta}] - \theta$

+ true value

General Definition

$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

$$\text{Var}_{\hat{\vartheta}} = E[\hat{\vartheta}^2] - (E[\hat{\vartheta}])^2$$

$$\text{Var}[\hat{\vartheta}] = E[(E[\hat{\vartheta}] - \hat{\vartheta})^2]$$

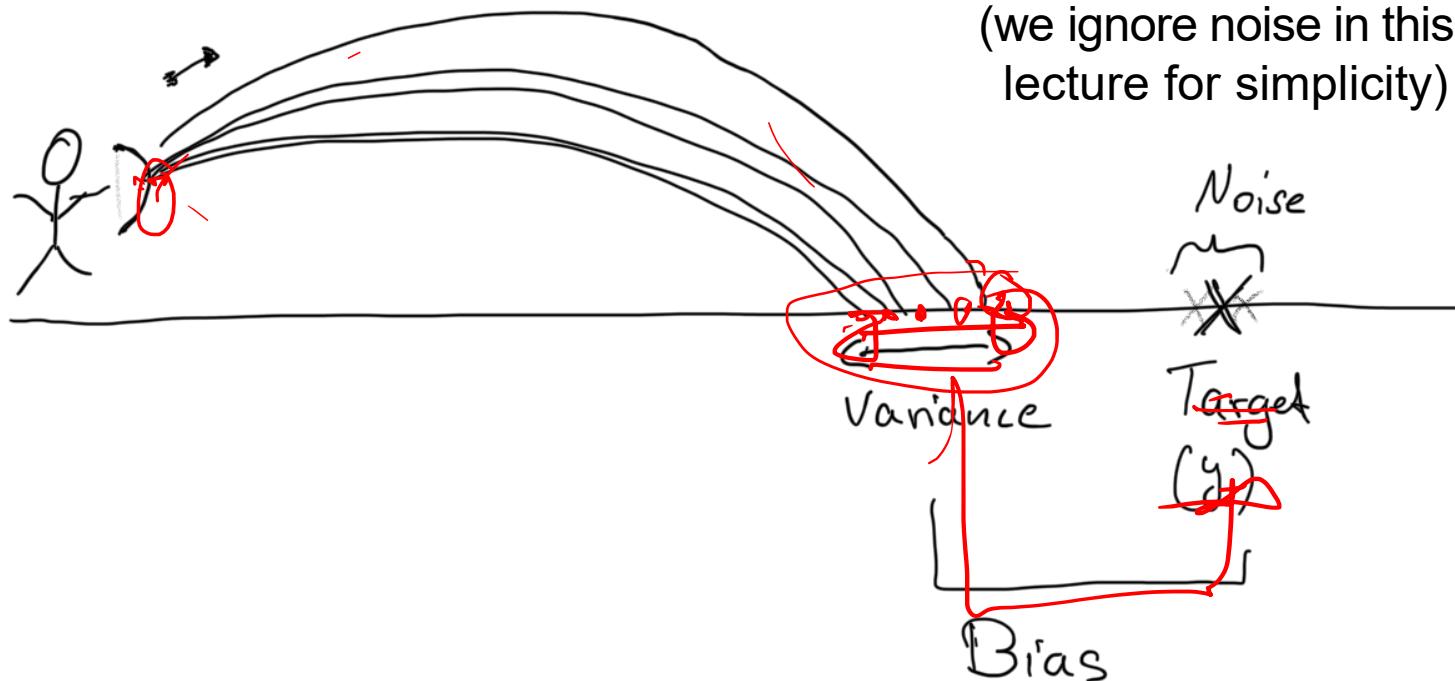
Terminology



$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

$$\text{Var}[\hat{\vartheta}] = E [(E[\hat{\vartheta}] - \hat{\vartheta})^2]$$

Intuition



Terminology



$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

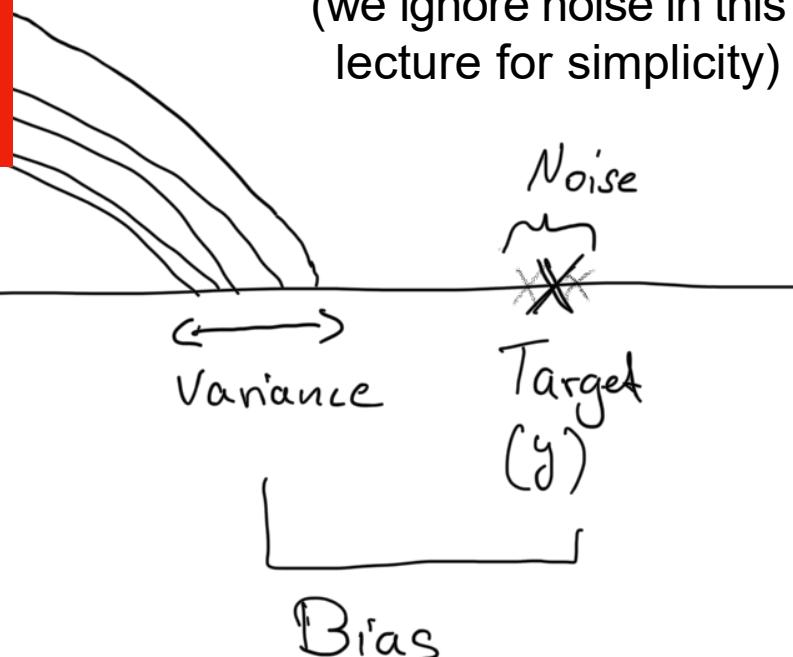
Bias is the difference between the average estimator from different training samples and the true value.

(The expectation is over the training



$$\text{Var}[\hat{\vartheta}] = E [(E[\hat{\vartheta}] - \vartheta)^2]$$

(we ignore noise in this lecture for simplicity)



Terminology



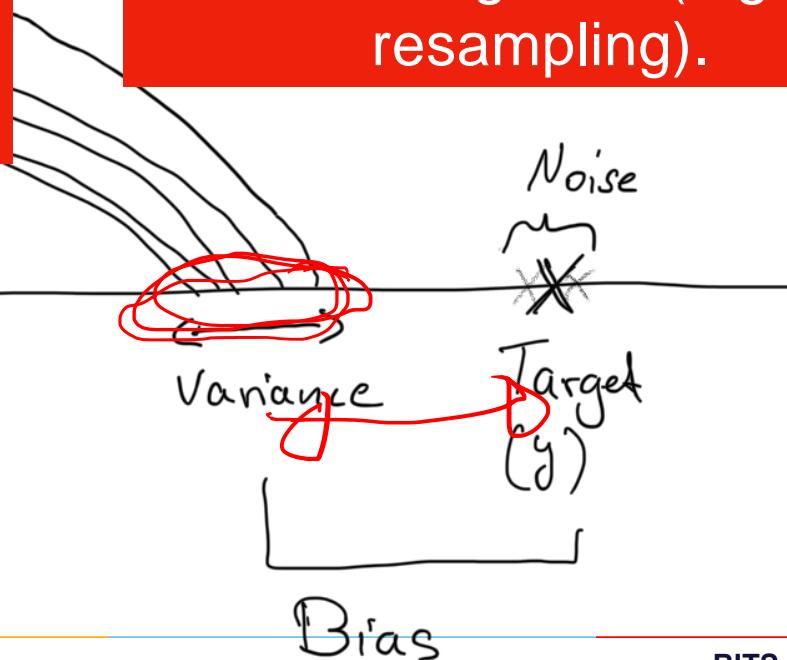
$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

Bias is the difference between the average estimator from different training samples and the true value.

(The expectation is over the

$$\text{Var} [\hat{\vartheta}] = E [(E[\hat{\vartheta}] - \hat{\vartheta})^2]$$

The variance provides an estimate of how much the estimate varies as we vary the training data (e.g., by resampling).



Bias-Variance Decomposition

Loss = Bias² + Variance + Noise

Bias-Variance of the Squared Error

$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

$$\text{Var}[\hat{\vartheta}] = E[\hat{\vartheta}^2] - (E[\hat{\vartheta}])^2$$

$$\text{Var}[\hat{\vartheta}] = E[(E[\hat{\vartheta}] - \hat{\vartheta})^2]$$

"ML Notation" for Squared Error Loss

$y = f(x)$ target

$\hat{y} = \hat{f}(x) = h(x)$ prediction

for simplicity, we ignore the noise term

$S = (y - \hat{y})^2$ squared error

(Next slides: the expectation is over the training data, i.e., the average estimator from different training samples)

Bias-Variance of the Squared Error

"ML Notation" for

$y = f(x)$ target

$\hat{y} = \hat{f}(x) = h(x)$ prediction

**Squared Error
Loss**

$S = (y - \hat{y})^2$ squared error

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 - 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

Bias-Variance of the Squared Error

$$S = (y - \hat{y})^2$$

$$(y - \hat{y})^2 = (y - E[\hat{y}] + E[\hat{y}] - \hat{y})^2$$

$$= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2 - 2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})$$

???

$$\begin{aligned}
 E[2(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] &= 2E[(y - E[\hat{y}])(E[\hat{y}] - \hat{y})] \\
 &= 2(y - E[\hat{y}])E[(E[\hat{y}] - \hat{y})] \\
 &= 2(y - E[\hat{y}])(E[E[\hat{y}]] - E[\hat{y}]) \\
 &= 2(y - E[\hat{y}])(E[\hat{y}] - E[\hat{y}]) \\
 &= 0
 \end{aligned}$$

Bias-Variance of the Squared Error

$$S = (y - \hat{y})^2$$

$$\begin{aligned}
 (y - \hat{y})^2 &= (y - E[\hat{y}]) + E[\hat{y}] - \hat{y})^2 \\
 &= (y - E[\hat{y}])^2 + (E[\hat{y}] - \hat{y})^2
 \end{aligned}$$

Bias Variance

$$E[S] = E[(y - \hat{y})^2]$$

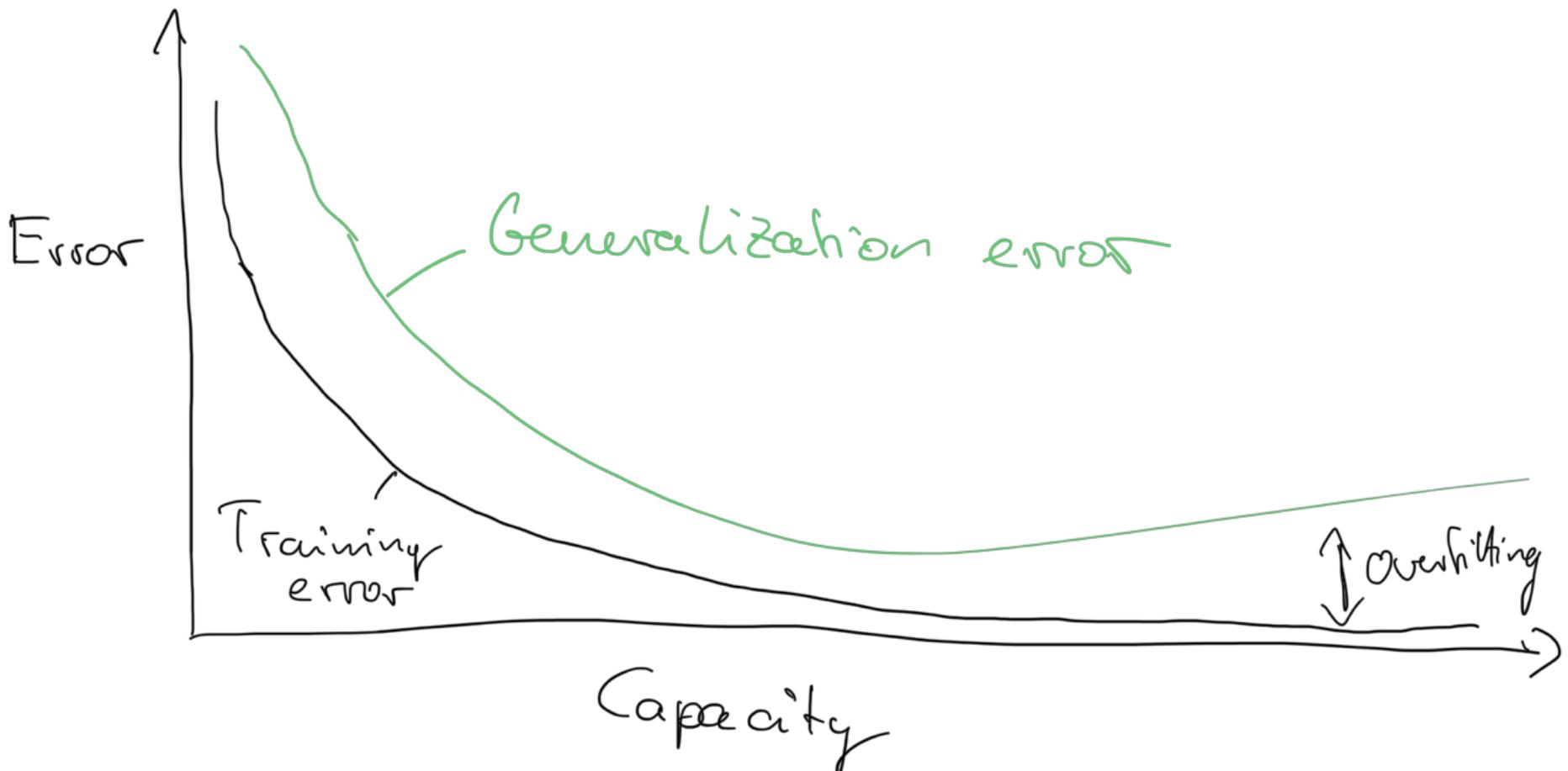
$$\begin{aligned}
 E[(y - \hat{y})^2] &= (y - E[\hat{y}])^2 + E[(E[\hat{y}] - \hat{y})^2] \\
 &= \text{Bias}^2 + \text{Var}
 \end{aligned}$$

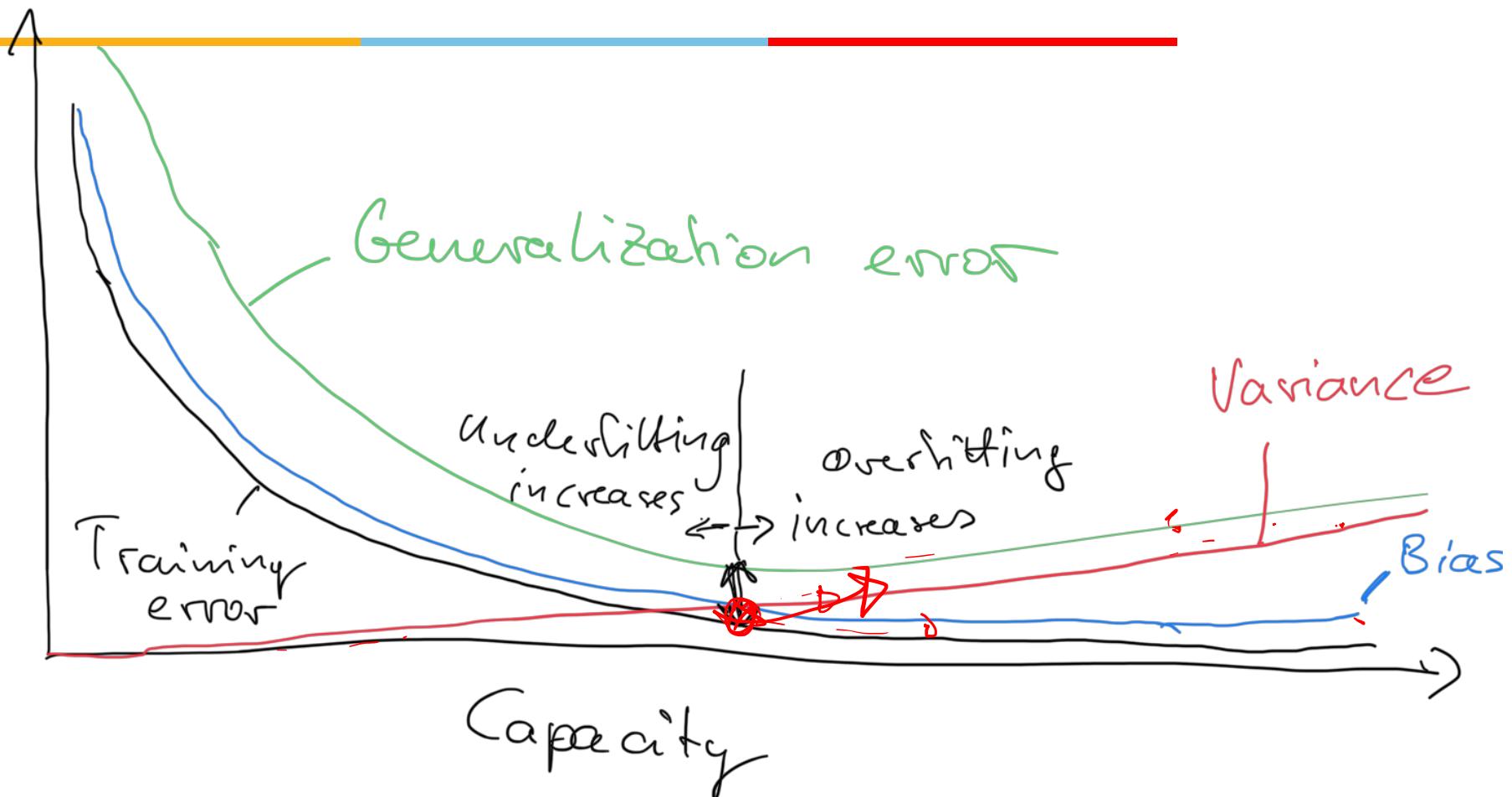
$$\text{Bias}[\hat{\vartheta}] = E[\hat{\vartheta}] - \vartheta$$

$$\text{Var}[\hat{\vartheta}] = E[\hat{\vartheta}^2] - (E[\hat{\vartheta}])^2$$

$$\text{Var}[\hat{\vartheta}] = E[(E[\hat{\vartheta}] - \hat{\vartheta})^2]$$

Now, how is this related to overfitting and underfitting?







Thank you !