

Assignment Project Exam Help

Machine learning lecture slides

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Regression I: Linear regression

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- ▶ Statistical model for regression
- ▶ College GPA example
- ▶ Ordinary least squares for linear regression
- ▶ The expected mean squared error

Assignment Project Exam Help

- ▶
- ▶
- ▶
- ▶

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Figure 1: Galton board

Real-valued predictions

- ▶ Example: Galton board
- ▶ Physical model: hard
- ▶ Statistical model: final position of ball is random
 - ▶ Normal (Gaussian) distribution with mean μ and variance σ^2

<https://eduassistpro.github.io>
 $\pi\sigma$

- ▶ Goal: predict final position accurately, i.e.
(also called squared error)

$$(\text{prediction} - \text{outcome})^2$$

- ▶ Outcome is random, so look at expected squared loss (also called mean squared error)

Optimal prediction for mean squared error

- ▶ Predict $\hat{y} \in \mathbb{R}$; true final position is Y (random variable) with mean $\mathbb{E}(Y) = \mu$ and variance $\text{var}(Y) = \mathbb{E}[(Y - \mathbb{E}(Y))^2] = \sigma^2$.
- ▶ Square error is $(\hat{y} - Y)^2$.
- ▶ Bias-variance decomposition:

<https://eduassistpro.github.io>

$$y - \mu \quad \sigma^2.$$

Add WeChat [edu_assist_pro](#)

- ▶ This is true for any random variable assumption.
- ▶ So optimal prediction is $\hat{y} = \mu$.
- ▶ When parameters are unknown, can estimate from related data, ...
- ▶ Can also do an analysis of a plug-in prediction ...

Statistical model for regression

- ▶ Setting is same as for classification except:
 - ▶ Label is real number, rather than $\{0, 1\}$ or $\{1, 2, \dots, K\}$
 - ▶ Care about squared loss, rather than whether prediction is correct
 - ▶ Mean squared error of f :

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Optimal prediction function for regression

- If (X, Y) is random test example, then

optimal prediction function is

Assignment Project Exam Help

$$f^*(x) = \mathbb{E}[Y \mid X = x]$$

function
<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Test MSE (1)

- ▶ Just like in classification, we can use test data to estimate $\text{mse}(\hat{f})$ for a function \hat{f} that depends only on training data.

▶ IID model

$(X_1, Y_1), \dots, (X_n, Y_n), (X'_1, Y'_1), \dots, (X'_m, Y'_m), (X, Y)$ are iid

<https://eduassistpro.github.io>

- ▶ Predictor \hat{f} is based only on training exa
- ▶ Hence, **test examples are independent** (important!)
- ▶ We would like to estimate $\text{mse}(\hat{f})$

Add WeChat edu_assist_pro

Test MSE (2)

- ▶ Test MSE $\text{mse}(\hat{f}, T) = \frac{1}{m} \sum_{i=1}^m (\hat{f}(X'_i) - Y'_i)^2$
- ▶ By law of large numbers, $\text{mse}(\hat{f}, T) \rightarrow \text{mse}(\hat{f})$ as $m \rightarrow \infty$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Example: College GPA

- ▶ Data from 750 Dartmouth students' College GPA
 - ▶ Mean: 2.46
 - ▶ Standard deviation: 0.746
- ▶ Assume this data is iid sample from the population of Dartmouth students (false)

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Predicting College GPA from HS GPA (1)

- ▶ Students represented in data have High School (HS) GPA
 - ▶ Maybe HS GPA is predictive of College GPA?

Assignment Project Exam Help

▶ Data: $S := ((x_1, y_1) \dots, (x_n, y_n))$

- ▶ x_i is HS GPA of i -th student

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

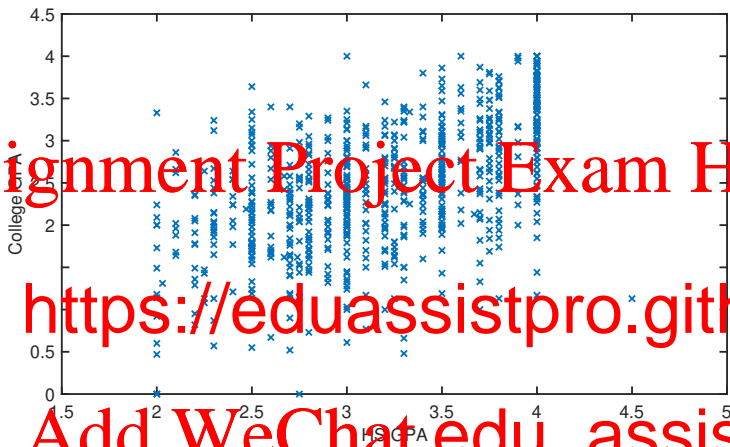


Figure 3: Plot of College GPA vs HS G

Predicting College GPA from HS GPA (2)

- ▶ First attempt:

- ▶ Define intervals of possible HS GPAs:

$(0.00, 0.25]$, $(0.25, 0.50]$, $(0.50, 0.75]$, ...

ge

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$\hat{f}(x) := \begin{cases} \hat{\mu}_{(0.25, 0.50]} \\ \hat{\mu}_{(0.50, 0.75]} \end{cases}$$

- ▶ (What to do about an interval I that contains no student's HS GPA?)

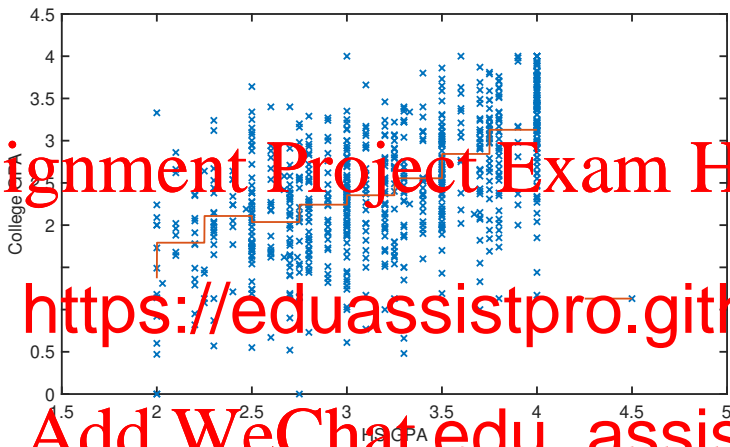


Figure 4: Plot of mean College GPA vs binne

Predicting College GPA from HS GPA (3)

- Define

$$\text{mse}(f, S) := \frac{1}{|S|} \sum_{(x, y) \in S} (f(x) - y)^2,$$

the mean squared error of predictions made by f on examples

<https://eduassistpro.github.io>

$$\text{mse}(\hat{f}, S) = 0.376$$

$\sqrt{\text{mse}(\hat{f}, S)} = 0.613 < 0.710$ (Add WeChat [edu_assist_pro](#))

- Piece-wise constant function \hat{f} is an improvement over the constant function (i.e., just predicting the mean 2.46 for all x)!

Predicting College GPA from HS GPA (4)

- ▶ But \hat{f} has some quirks.
- ▶ E.g., those with HS GPA between 2.50 and 2.75 are predicted to have a lower College GPA than those with HS GPA between 2.25 and 2.50.
- ▶ E.g., something unusual with the student who has HS GPA of

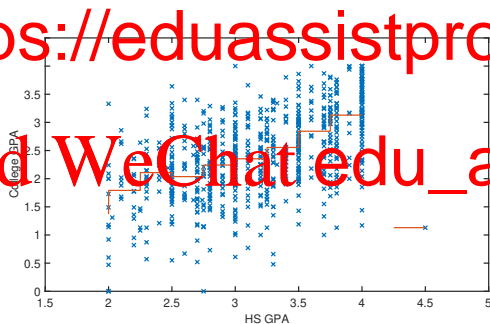


Figure 5: Plot of mean College GPA vs binned HS GPA

Least squares linear regression (1)

- Suppose we'd like to only consider functions with a specific functional form, e.g., a linear function:

Assignment Project Exam Help

$$f(x) = mx + \theta$$

<https://eduassistpro.github.io>

m

prediction of College GPA.

Add WeChat edu_assist_pr

Least squares linear regression (2)

- ▶ What is the linear function with smallest MSE on $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R} \times \mathbb{R}$? This is the problem of

least squares linear regression

Assignment Project Exam Help

- ▶ Find $(m, \theta) \in \mathbb{R}^2$ to minimize

<https://eduassistpro.github.io>

- ▶ Also called ordinary least squares ____

Add WeChat edu_assist_pr

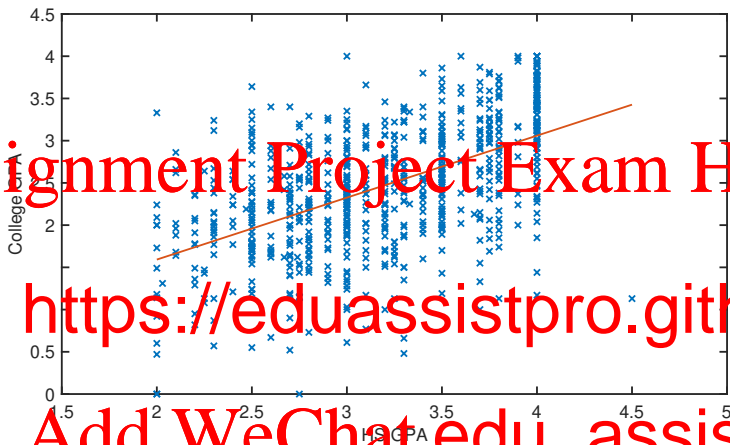


Figure 6: Plot of least squares linear regression

Computing OLS (1)

- Derivatives equal zero conditions ([normal equations](#)):

$$\frac{\partial}{\partial \theta} \left\{ \frac{1}{n} \sum_{i=1}^n (mx_i + \theta - y_i)^2 \right\} = \frac{2}{n} \sum_{i=1}^n (mx_i + \theta - y_i) = 0$$

$$\frac{\partial}{\partial m} \left\{ \frac{1}{n} \sum_{i=1}^n (mx_i + \theta - y_i)^2 \right\} = \frac{2}{n} \sum_{i=1}^n x_i (mx_i + \theta - y_i) = 0.$$

- Define

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} := \frac{1}{n} \sum_{i=1}^n y_i,$$

$$\overline{xy} := \frac{1}{n} \sum_{i=1}^n x_i y_i, \quad \bar{y} := \frac{1}{n} \sum_{i=1}^n y_i,$$

so system can be re-written as

$$\bar{x}m + \theta = \bar{y}$$

$$\bar{x}^2 m + \bar{x} \theta = \overline{xy}.$$

- Write in matrix notation:

$$\begin{bmatrix} 1 & 1 \\ \frac{1}{x^2} & \frac{1}{x} \end{bmatrix} \begin{bmatrix} \hat{m} \\ \hat{\theta} \end{bmatrix} = \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix}$$

<https://eduassistpro.github.io>

- Solution: $(\hat{m}, \hat{\theta}) \in \mathbb{R}^2$ given by

$$\hat{m} := \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2}, \quad \hat{\theta} := - \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2}$$

Computing OLS (3)

- ▶ Catch: The above solution only makes sense if $\overline{x^2} - \bar{x}^2 \neq 0$, i.e., the variance of the x_i 's is non-zero.

Assignment Project Exam Help

<https://eduassistpro.github.io>

- ▶ If $\overline{x^2} - \bar{x}^2 = 0$, then the matrix defining the equations is singular.

Add WeChat edu_assist_pro

- ▶ In general, “derivative equals zero” is only a necessary condition for a solution to be optimal; not necessarily a sufficient condition!

Assignment Project Exam Help

<https://eduassistpro.github.io>

- ▶ **Theorem** Every solution to the normal eq is an optimal solution to the least squares linear r

Add WeChat: edu_assist_pro

Decomposition of expected MSE (1)

- ▶ Two different functions of HS GPA for predicting College GPA.
 - ▶ What makes them different?
 - ▶ We care about prediction of College GPA for student we haven't seen before based on their HS GPA.

Assignment Project Exam Help

- ▶
- ▶
- ▶

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Decomposition of expected MSE (2)

Assignment Project Exam Help

$$\begin{aligned}\mathbb{E}[\text{mse}(\hat{f})] &= \mathbb{E}[\mathbb{E}[(f(X) - Y)^2 \mid \hat{f}]] \\ &= \mathbb{E}[\mathbb{E}[(f(X) - \mathbb{E}[Y \mid X])^2 \mid \hat{f}]] \\ &= \mathbb{E}[\text{var}(Y \mid X) + \mathbb{E}[(f(X) - \mathbb{E}[Y \mid X])^2 \mid \hat{f}]] \\ &= \mathbb{E}[\text{var}(Y \mid X) + \text{var}(\hat{f}(X) \mid X) + \mathbb{E}[(f(X) - \mathbb{E}[Y \mid X] - \hat{f}(X) + \hat{f}(X))^2 \mid X]] \\ &= \underbrace{\mathbb{E}[\text{var}(Y \mid X)]}_{\text{unavoidable error}} + \underbrace{\mathbb{E}[\text{var}(\hat{f}(X) \mid X)]}_{\text{variability of } \hat{f}} + \underbrace{\mathbb{E}[\mathbb{E}[(f(X) - \mathbb{E}[Y \mid X] - \hat{f}(X) + \hat{f}(X))^2 \mid X]]}_{\text{approximation error of } \hat{f}}\end{aligned}$$

<https://eduassistpro.github.io>

Add WeChat: edu_assist_pro

Decomposition of expected MSE (3)

- ▶ First term quantifies inherent unpredictability of Y (even after seeing X)
- ▶ Second term measures the “variability” of \hat{f} due to the random nature of training data. Depends on:

<https://eduassistpro.github.io>

- ▶ Third term quantifies how well a function p fitting procedure can approximate the reg after removing the “variability” of

Multivariate linear regression (1)

- ▶ For Dartmouth data, also have SAT Score for all students.
 - ▶ Can we use both predictor variables (HS GPA and SAT Score) to get an even better prediction of College GPA?
 - ▶ Binning approach: instead of a 1-D grid (intervals), consider a 2-D grid (squares).

<https://eduassistpro.github.io>

for some $(m_1, m_2) \in \mathbb{R}^2$ and θ

Add WeChat edu_assist_pr

Multivariate linear regression (2)

- The general case: a (homogeneous) linear function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ of the form

Assignment Project Exam Help
 $f(x) = x^T w$
for some $w \in \mathbb{R}^d$.

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Multivariate ordinary least squares (1)

- ▶ What is the linear function with smallest MSE on $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$?

▶ Find $w \in \mathbb{R}^d$ to minimize

$$\frac{1}{n} \sum_{i=1}^n (w^\top x_i - y_i)^2$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Multivariate ordinary least squares (2)

- In matrix notation:

Assignment Project Exam Help

where

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

- If we put vector $v \in \mathbb{R}^d$ in the context, it is treated as a column vector by default!
- If we want a row vector, we write

- Therefore

$$Aw - b = \frac{1}{\sqrt{n}} \begin{bmatrix} x_1^\top w - y_1 \\ \vdots \\ x_n^\top w - y_n \end{bmatrix}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Figure 7: Geometric picture of least squares linear regression

Multivariate normal equations (1)

- ▶ Like the one-dimensional case, optimal solutions are characterized by a system of linear equations (the “derivatives equal zero” conditions) called the normal equations:

$$\frac{\partial}{\partial w_d} (w)$$

$$\frac{\partial}{\partial w_d} \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

which is equivalent to

$$A^T A w =$$

Multivariate normal equations (2)

- ▶ If $A^T A$ is non-singular (i.e., invertible), then there is a unique solution given by

Assignment Project Exam Help

$$\hat{w} := (A^T A)^{-1} A^T b.$$



<https://eduassistpro.github.io>



Add WeChat edu_assist_pr

Theorem: Every solution to the normal eq
optimal solution to the least squares linear r

Algorithm for least squares linear regression

- ▶ How to solve least squares linear regression problem?
 - ▶ Just solve the normal equations, a system of d linear equations in d unknowns.
 - ▶ Time complexity (naïve) of Gaussian elimination algorithm: $O(d^3)$.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Classical statistics view of OLS (1)

- ▶ Normal linear regression model

- ▶ Model training examples $(X_1, Y_1), \dots, (X_n, Y_n)$ as iid random variables taking values in $\mathbb{R}^d \times \mathbb{R}$ where

$$Y \mid X = x \sim \mathcal{N}(x^\top w, \sigma^2)$$

- ▶ <https://eduassistpro.github.io>
problem of finding the maximum likelihood

Add WeChat edu_assist_pr

Classical statistics view of OLS (2)

- ▶ Suppose your data really does come from a distribution in this statistical model, say, with parameters w and σ^2 .

Assignment Project Exam Help

- ▶ Then the function with smallest MSE is the linear function $f^*(x) = x^T w$, and its MSE is $\text{mse}(f^*) = \sigma^2$.

- ▶ So estimating w is a sensible idea! (Plug-in principle...)

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Statistical learning view of OLS (1)

- ▶ IID model: $(X_1, Y_1), \dots, (X_n, Y_n), (X, Y) \sim_{\text{iid}} P$ are iid random variables taking values in $\mathbb{R}^d \times \mathbb{R}$
 - ▶ (X, Y) is the (unlabeled) 'test' example
- ▶ Goal: find a (linear) function $w \in \mathbb{R}^d$ with small MSE

- ▶ <https://eduassistpro.github.io>
 $w \in \mathbb{R}^d$,
since it is an expectation (e.g., integral) with respect to unknown distribution P

Add WeChat edu_assist_pro

Statistical learning view of OLS (2)

- ▶ However, we have an iid sample $S := ((X_1, Y_1), \dots, (X_n, Y_n))$.
- ▶ We swap out P in the definition of $\text{mse}(f)$, and replace it with the empirical distribution on S .

Assignment Project Exam Help

n

—

<https://eduassistpro.github.io>

on the

i -th training example.

- ▶ Resulting objective function is

Add WeChat edu_assist_pro

$$\mathbb{E}[(\tilde{X}^\top w - \tilde{Y})^2] = \frac{1}{n} \sum_{i=1}^n (X_i^\top w - Y_i)^2$$

where $(\tilde{X}, \tilde{Y}) \sim P_n$.

Statistical learning view of OLS (3)

- ▶ In some circles:

- ▶ (True/population) risk of w : $\mathcal{R}(w) := \mathbb{E}[(X^\top w - Y)^2]$

- ▶ Empirical risk of w : $\hat{\mathcal{R}}(w) := \frac{1}{n} \sum_{i=1}^n (X_i^\top w - Y_i)^2$

- ▶ This is another instance of the plug-in principle!

e

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Statistical learning view of OLS (4)

- ▶ This is not specific to linear regression; also works for other types of functions, and also other types of prediction problems, including classification.
- ▶ For classification:

<https://eduassistpro.github.io>

- ▶ Procedure that minimizes empirical risk:

Empirical risk minimization (ERM)

Add WeChat edu_assist_pro

Upgrading linear regression (1)

- ▶ Make linear regression more powerful by being creative about features

Assignment Project Exam Help

- ▶ We are forced to do this if x is not already provided as a vector of numbers

- ▶ ion φ

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Upgrading linear regression (2)

- ▶ Examples:

- ▶ Affine feature expansion, e.g., $\varphi(x) = (1, x)$, to accommodate intercept

- ▶ Standardization, e.g., $\varphi(x) = (x - \mu)/\sigma$ where (μ, σ^2) are (estimates of) the mean and variance of the feature value x

<https://eduassistpro.github.io>

- ▶ Polynomial expansion, e.g.,

- $\varphi(x) = (1, x_1, \dots, x_d, x_1^2, \dots, x_d^2)$

- ▶ Headless neural network, $\varphi(x) =$

- $N: \mathbb{R}^d \rightarrow \mathbb{R}^k$ is a map computed by a in neural network

- ▶ (Later, we'll talk about how to “learn” N .)

Example: Taking advantage of linearity

- ▶ Example: y is health outcome, x is body temperature
 - ▶ Physician suggests relevant feature is (square) deviation from normal body temperature $(x - 98.6)^2$
 - ▶ What if you didn't know the magic constant 98.6? (Apparently it is wrong in the US anyway)

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Example: Binning features

- ▶ Dartmouth data example, where we considered intervals for the HS GPA variable:

Assignment Project Exam Help
 $(0.00, 0.25], (0.25, 0.50], (0.50, 0.75], \dots$

- ▶ ear
- ▶ <https://eduassistpro.github.io>
- ▶ $\varphi(x)^\top w = w_j$ if x is in the j -th i

Add WeChat edu_assist_pr

Effect of feature expansion on expected MSE

$$\begin{aligned} \mathbb{E}[\text{mse}(\hat{f})] &= \underbrace{\mathbb{E}[\text{var}(Y | X)]}_u + \underbrace{\mathbb{E}[\text{var}(\hat{f}(X) | X)]}_{\text{approximation error}} + \underbrace{\mathbb{E}[(\mathbb{E}[\hat{f}(X) | X] - \mathbb{E}[Y | X])^2]}_{\text{error of } \hat{f}} \end{aligned}$$

- ▶ <https://eduassistpro.github.io>
(approximation error)
- ▶ But maybe at the cost of increasing the second (variability)

Performance of OLS (1)

- ▶ Study in context of IID model
- ▶ $(X_1, Y_1), \dots, (X_n, Y_n), (X, Y)$ are iid, and assume $\mathbb{E}[XX^T]$ is invertible (WLOG).
- ▶ Let w^* denote the minimizer of $\text{mse}(w)$ over all $w \in \mathbb{R}^a$.

Assignment Project Exam Help

<https://eduassistpro.github.io>

- ▶ How much larger is $\text{mse}(\hat{w})$ comp

Add WeChat edu_assist_pr

Performance of OLS (2)

- **Theorem:** In the IID model, the OLS solution \hat{w} satisfies

Assignment Project Exam Help

$$\frac{\mathbb{E}[\text{mse}(\hat{w})] - \text{mse}(w^*)}{\text{mse}(w^*)} \rightarrow \frac{\text{tr}(\text{cov}(\varepsilon W))}{n}, \text{ where } W = \mathbb{E}[XX^\top]^{-1/2}X \text{ and } \varepsilon = Y - X^\top w^*.$$

- <https://eduassistpro.github.io>

Add WeChat edu_assist_pro

which is more typically written as

$$\mathbb{E}[\text{mse}(\hat{w})] \rightarrow \left(1 + \frac{d}{n}\right) \text{mse}(w^*).$$

Linear algebraic view of OLS (1)

- ▶ Write $A = \begin{bmatrix} \uparrow & & \uparrow \\ a_1 & \cdots & a_d \\ \downarrow & & \downarrow \end{bmatrix}$
 - ▶ $a_j \in \mathbb{R}^n$ is j -th column of A
 - ▶ Span of a_1, \dots, a_d is $\text{range}(A)$, a subspace of \mathbb{R}^n

Assignment Project Exam Help

▶ <https://eduassistpro.github.io>ding

Add WeChat edu_assist_pr

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

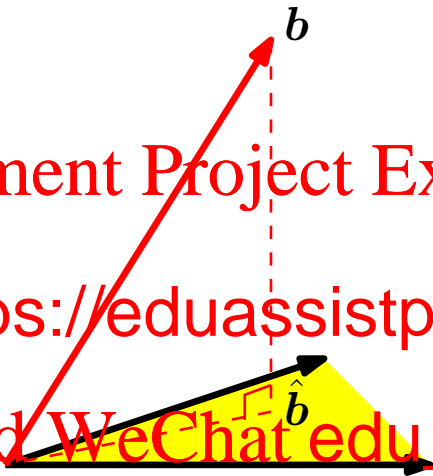


Figure 8: Orthogonal projection of b onto $\text{range}(A)$

Linear algebraic view of OLS (2)

- ▶ Solution \hat{b} is orthogonal projection of b onto $\text{range}(A)$

- ▶ \hat{b} is unique

- ▶ Residual $b - \hat{b}$ is orthogonal to \hat{b}

- ▶ To get w from \hat{b} , solve $Aw = \hat{b}$ for w .

- ▶ If $\text{rank}(A) < d$ (always the case if $n < d$), then infinitely-many

- ▶ <https://eduassistpro.github.io> $n \geq d$,

Add WeChat edu_assist_pr

Over-fitting (1)

- In the IID model, over-fitting is the phenomenon where the true risk is much worse than the empirical risk.

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Over-fitting (2)

- ▶ Example:

- ▶ $\varphi(x) = (1, x, x^2, \dots, x^k)$, degree- k polynomial expansion

- ▶ Dimension is $d = k + 1$

- ▶ Any function of $\leq k + 1$ points can be interpolated by polynomial of degree $\leq k$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

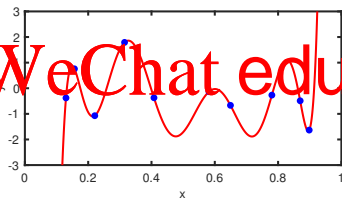


Figure 9: Polynomial interpolation

- ▶ Recall plug-in principle
 - ▶ Want to minimize risk with respect to (unavailable) P ; use P_n instead
- ▶ What if we can't regard data as iid from P ?

$\frac{1}{n}$ $\frac{1}{n}$

<https://eduassistpro.github.io>

- ▶ How to implement plug-in principle?

Add WeChat edu_assist_pr